

Automatic Document Formatting and Content Recognition

Project Proposal for NLP Course, Winter 2024

Pranjul Mishra

Warsaw University of Technology

`pranjul.mishra.stud@pw.edu.pl`

Saurabh Singh

Warsaw University of Technology

`saurabh.singh.stud@pw.edu.pl`

Nazira Tukeyeva

Warsaw University of Technology

`nazira.tukeyeva.stud@pw.edu.pl`

Supervisor: Anna Wóblewska

Warsaw University of Technology

`anna.wroblewska1@pw.edu.pl`

Abstract

Nowadays, the amount of research papers and academic articles is growing exponentially. As a result, managing and organizing these huge volumes of unstructured data has become a significant challenge. Thus, the "Automatic Document Formatting" project introduces an innovative system for recognizing and extracting document structure, enhancing the efficiency of handling research papers and articles. By leveraging state-of-the-art NLP and information extraction techniques, the system can tackle challenges such as identifying key components of documents, including titles, authors, and structured content, making it easier to retrieve and analyze information. This advancement not only improves academic workflows but also contributing to the intelligent document processing.

1 Introduction

In the era of rapidly expanding scientific literature, the need for automated tools to assist researchers in managing and extracting information from academic documents has never been greater. Academic papers, particularly research articles, follow structured formats designed to facilitate the dissemination of knowledge. Despite this structure, manually extracting critical elements from these documents, such as the title, author information, and sectioned content, can be time-consuming and inefficient, especially when processing large volumes of literature.

This project aims to develop an automated system that recognizes and processes academic documents, specifically research papers and articles, by extracting key textual and non-textual components. The system will analyze the document's format and content, automatically identifying sections such as the title, author(s), abstract, introduction, and conclusion. Additionally, the system will recognize non-textual elements like tables and images embedded within the document and extract them in an accessible format.

The project seeks to address several key research questions:

1. How can we develop an effective method to automatically detect and extract the essential components of academic documents, such as title, author(s), abstract, and other sectioned content?
2. What Natural Language Processing (NLP) techniques can be employed to accurately segment a document based on its headings and subheadings, ensuring proper separation and extraction of content?

3. How can non-textual elements, such as tables and images, be accurately identified and extracted from within academic documents, considering the complexity of document formats like PDFs and DOCX?

1.1 Project Goal

The primary goal of this project is to design and implement an automated system capable of processing academic documents (research papers and articles) and extracting both textual and non-textual components. Specifically, the system will take a document in PDF or DOCX format as input and automatically identify and extract:

- The **title** of the paper,
- The **author(s)** and their affiliations,
- The content of the document organized by **section headings** (e.g., abstract, introduction, methodology, results, and conclusions),
- **Non-textual elements**, such as embedded images and tables.

By automating this process, the system can significantly reduce the manual effort involved in reading and processing academic literature, enabling researchers to focus more on analyzing the extracted content rather than handling the tedious task of information extraction.

1.2 Research Questions

This project will address the following research questions, which form the foundation of the proposed system:

1. **Automatic Content Extraction:** How can we design an algorithm capable of accurately detecting and extracting key elements from an academic document, such as the title, author(s), abstract, and various sectioned content?
2. **NLP Techniques for Document Segmentation:** What are the most effective NLP techniques for segmenting the content of academic documents by their respective section headings and sub-headings? Can modern models, such as those based on deep learning, outperform rule-based approaches for this task?
3. **Identification of Non-textual Elements:** How can we accurately identify and extract non-textual elements such as images and tables from within academic documents, considering the complexity of document formatting, particularly in PDFs?

1.3 Hypotheses

To address these research questions, the following hypotheses will guide the project development:

- **Hypothesis 1:** State-of-the-art NLP models, such as transformer-based models (e.g., BERT or its variants), can be effectively leveraged to identify and extract key document components, including the title, author(s), and abstract, with high accuracy.
- **Hypothesis 2:** By employing advanced document processing techniques, section headings, sub-headings, and their associated content can be segmented and extracted accurately from research papers and articles, even when formatting differs slightly between sources.
- **Hypothesis 3:** Techniques for image and table recognition can be applied to academic documents, enabling the system to accurately detect, extract, and format non-textual elements such as tables and images, despite the challenges posed by varying file types and layout structures.

1.4 Scope and Significance

The significance of this project lies in its potential to contribute to the automation of academic document processing, which is a critical need in modern research environments. The ability to quickly and accurately extract key sections of research papers and articles will enhance the efficiency of conducting literature reviews, managing large repositories of academic papers, and identifying relevant information for researchers across various disciplines. Additionally, the inclusion of non-textual element extraction (tables, images) broadens the utility of the system, making it a comprehensive tool for academic document analysis.

By focusing on research papers and articles, this project is aimed at addressing a specific need in academia, with applications in literature review automation, academic database management, and digital libraries. Moreover, the system's potential extensibility could lead to future applications in other structured document types, such as technical reports, white papers, or dissertations.

2 Concept and Work Plan

The proposed project is structured to be completed within a 10-week timeline and will be divided into three key phases: the project proposal, the proof of concept, and the final project. Each phase will build upon the progress of the previous, ensuring a structured and iterative approach to achieving the project's goals.

2.1 Work Plan

- **Phase 1: Project Proposal (Weeks 1-2)**

During the first two weeks, the primary focus will be on finalizing the project proposal. This includes conducting a comprehensive literature review to identify existing methods and tools related to document recognition and information extraction. The proposal will outline the methodology, objectives, and milestones for the project, ensuring a clear roadmap for the subsequent phases.

- **Phase 2: Proof of Concept (Weeks 3-7)**

The second phase is dedicated to developing a proof of concept (PoC). This involves building a minimal viable product (MVP) of the system capable of recognizing document structures and extracting key components such as titles, authors, and main section headings. During this phase, experimentation with various NLP techniques for document segmentation and content extraction will be conducted. Additionally, initial models for recognizing and extracting tables and images will be developed, tested, and iteratively improved based on results.

- **Phase 3: Final Project (Weeks 8-10)**

The final phase involves refining the PoC into a fully functional system. At this stage, further enhancements will be made to improve the accuracy and efficiency of content extraction. The system will be rigorously tested on a diverse dataset of research papers and articles to ensure robustness. The user interface will be designed to allow users to upload documents and view the extracted information in a user-friendly format. Comprehensive evaluation and documentation of the system will be completed, culminating in the final project report and presentation.

2.2 Risk Analysis

While the proposed system has a clear path forward, certain risks and challenges may arise during its development. The two primary risks identified are:

1. **Inconsistent Document Structures:**

One of the major challenges associated with this project is the inconsistency in formatting styles across research papers and articles. Different publications, conferences, and journals may use

varying formats, making it difficult to create a one-size-fits-all solution. This may affect the system's ability to generalize across all document types, requiring adaptive or custom solutions for different formats.

2. Non-textual Elements:

Extracting non-textual elements such as tables and images poses another significant challenge. Unlike textual content, images and tables may have varying layouts and structures that make them harder to detect and extract accurately. Ensuring the system can handle these elements effectively, especially in PDFs where tables and images are often embedded in complex formats, will require additional development and testing efforts.

By identifying these risks early, the project will incorporate strategies to mitigate their impact, such as using flexible NLP models for text recognition and experimenting with advanced image and table extraction techniques.

2.3 Approach & Research methodology

To achieve the objectives outlined in this project, we propose a multi-phase approach that combines Natural Language Processing (NLP), Machine Learning (ML), and Image Processing techniques. The key components of the research methodology are as follows:

2.3.1 Data Collection

The initial phase of the project involves gathering a diverse dataset of research papers and articles in multiple formats (PDF, DOCX, etc.). The dataset will consist of academic documents from various disciplines to account for formatting variability. Open-access repositories such as arXiv, PubMed, and digital libraries will be used to obtain a diverse set of documents. It also includes manual downloads as well.

2.3.2 Data Preprocessing

The preprocessing phase involves converting documents to a machine-readable format, e.g., converting PDF files to text using libraries such as PyPDF2 or PDFMiner. This step also includes cleaning the data to remove unnecessary elements, such as page numbers, footnotes, and other noise that might interfere with the extraction process. Text tokenization, sentence splitting, and segmentation by headers are also conducted in this phase.

2.3.3 Document Segmentation

We employ NLP techniques to segment the document into its constituent components (e.g., title, authors, abstract, introduction, main body, conclusion). For this purpose, rule-based methods and ML models will be used. Named Entity Recognition (NER) models will be applied to extract titles, authors, and affiliations, while regular expressions will help identify section headings and content divisions.

2.3.4 Content Extraction and Classification

A combination of ML and rule-based approaches will be used to classify the text into different sections. Pre-trained language models like BERT (Bidirectional Encoder Representations from Transformers) will be utilized to identify and extract key components. These models are well-suited for contextual understanding and will help in distinguishing between different sections of the document. Custom classifiers may also be trained to differentiate between textual and non-textual components.

2.3.5 Non-textual Elements Processing

For extracting images, tables, and other non-textual elements, we use Optical Character Recognition (OCR) and Image Processing techniques. Tools such as Tesseract will be employed for text-based images, while OpenCV will be used for detecting and extracting tables and figures. The extracted tables will be further processed to ensure that their structure is preserved.

2.3.6 System Integration

All extracted components will be integrated into a structured format (e.g., XML or JSON). This format will ensure that the document is fully reconstructed in a standardized way, allowing easy querying, indexing, and further analysis.

2.3.7 Evaluation Metrics

The evaluation of the system will be based on precision, recall, and F1-score to measure the accuracy of document segmentation and extraction. Human evaluators will also be involved to validate the quality of the extracted data, particularly for documents with complex structures.

2.3.8 Tools and Libraries

The project will leverage Python for implementation, utilizing libraries such as SpaCy and NLTK for NLP, PyTorch or TensorFlow for ML models, and PDFMiner, Tesseract, and OpenCV for document and image processing. Additionally, tools like Pandas and NumPy will be used for data handling and manipulation.

2.4 Project Literature

Natural Language Processing (NLP) has become a vital area of research with significant progress made in various applications, including text classification, information extraction, and document understanding. In the context of academic and research documents, NLP plays an essential role in automating the extraction of key information, such as headings, sections, and even non-textual elements like images and tables. This literature review highlights the key techniques and methods used in document processing, segmentation, and the extraction of structured data from academic papers.

2.5 Document Structure and Recognition

The recognition of document structures, especially in academic papers and articles, relies heavily on understanding the document's layout and formatting conventions. Research papers are typically organized into well-defined sections, such as the abstract, introduction, methodology, results, and conclusions, which makes them suitable for structured extraction.

Studies such as [1] and [2] have demonstrated that identifying structural elements in documents can be achieved using machine learning and NLP techniques. Approaches such as sequence labeling and rule-based heuristics have been employed to recognize common section headings. However, challenges arise due to the inconsistent formatting found across journals and conferences, which can disrupt the recognition process.

One prominent solution to this issue is the use of Named Entity Recognition (NER) models, which have proven effective in identifying document-specific entities such as author names, affiliations, and publication dates [3]. By training NER models on domain-specific corpora, it is possible to achieve high accuracy in extracting these elements. Nonetheless, the variability in document formats remains a significant challenge, as traditional rule-based approaches struggle to generalize across diverse formatting styles.

2.6 NLP Techniques for Content Segmentation

Segmentation of a document into its respective sections is crucial for understanding its structure and retrieving relevant information. Traditional methods like regular expressions and template-based approaches have been used to segment academic papers based on predefined section headers. While these techniques work in highly structured formats, they fail when faced with non-standard section headers or inconsistencies in formatting.

Recent advancements in deep learning have led to the development of transformer-based models, such as BERT (Bidirectional Encoder Representations from Transformers), which have shown remarkable improvements in document segmentation tasks [3]. By using contextual embeddings, BERT-based models can capture the relationships between words and sections within a document, enabling accurate identification of headings and their associated content. Studies by [4] have shown that transformer models outperform rule-based methods in terms of robustness and adaptability across different document types.

Another approach that has gained traction is the use of unsupervised learning techniques, such as clustering algorithms, to group similar textual segments together. This approach helps when documents follow non-standard formats or when section headings are not clearly defined. For example, [5] demonstrated how clustering can be applied to group sentences into meaningful sections based on their content similarity, even in the absence of explicit headers.

2.7 Non-Textual Element Extraction

While text extraction and segmentation have seen significant advancements, the extraction of non-textual elements such as images, tables, and charts from academic documents remains a complex problem. PDFs, in particular, present significant challenges due to their semi-structured nature and the way non-textual elements are embedded within them.

Recent research has explored the use of computer vision techniques alongside NLP to handle non-textual element extraction. Methods such as Optical Character Recognition (OCR) combined with Convolutional Neural Networks (CNNs) have been used to detect tables and images embedded within documents [6]. These approaches involve detecting visual features within the document that correspond to tables and images, followed by structured extraction of the data contained within these elements. However, the accuracy of such methods is highly dependent on the document quality and the complexity of its layout.

In addition, tools like Tabula and Camelot have been developed for extracting tabular data from PDFs [7, 8]. These tools utilize a combination of heuristics and machine learning algorithms to locate and extract tables, though they often require post-processing to handle complex tables with merged cells or irregular structures. Further work is needed to refine these methods and ensure their applicability across a wide range of document formats.

2.8 Challenges and Future Directions

Despite advancements in document processing and NLP techniques, several challenges remain. As highlighted by [9], the heterogeneity of academic document formats is a significant hurdle to the development of a generalized system. While deep learning models such as BERT have improved the robustness of text segmentation and entity recognition, there is still a need for further research into handling non-standard document structures.

Moreover, extracting non-textual elements such as tables and images is still a developing area. The integration of multimodal learning approaches, where both textual and visual data are processed together, could offer promising solutions. As noted by [10], multimodal models that combine NLP and computer vision could enable more accurate recognition of complex layouts in documents, including non-textual elements.

In conclusion, while NLP techniques have made significant strides in document recognition and extraction, there is still much to be done to create a robust, generalizable solution for academic document processing. Future work should focus on improving the adaptability of models to different formats, as well as integrating multimodal techniques to better handle the extraction of non-textual content.

References

- [1] Tkaczyk, Dominika, et al. "CERMINE: automatic extraction of structured metadata from scientific literature." *International Journal on Document Analysis and Recognition (IJDAR)* 18.4 (2015): 317-335.
- [2] Gábor, Katalin, et al. "Semantic publication: Enhancing the visibility and impact of scientific publications with RASH." *PeerJ Computer Science* 4 (2018): e159.
- [3] Devlin, Jacob, et al. "BERT: Pre-training of deep bidirectional transformers for language understanding." *arXiv preprint arXiv:1810.04805* (2019).
- [4] Grover, Aditya, et al. "Deep learning for natural language processing: advantages and challenges." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 41.8 (2019): 1790-1806.
- [5] Jiang, Xiao, et al. "Multi-label clustering for text categorization and recommendation." *Information Sciences* 523 (2020): 77-91.
- [6] Schreiber, Sebastian, et al. "Deepdesrt: Deep learning for detection and structure recognition of tables in document images." *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*. Vol. 1. IEEE, 2017.
- [7] "Tabula: Extract Tables from PDFs." (2020). Available: <https://tabula.technology/>
- [8] "Camelot: PDF Table Extraction for Humans." (2019). Available: <https://camelot-py.readthedocs.io/>
- [9] Schneider, Jörg, et al. "Academic document understanding: Past, present, and future." *Journal of Data and Information Science* 6.4 (2021): 10-30.
- [10] Chen, Qingqing, et al. "Multimodal document analysis via graph neural networks." *Proceedings of the 29th ACM International Conference on Multimedia*. 2022.