



SCHOOL OF COMPUTATION, INFORMATION  
AND TECHNOLOGY - INFORMATICS

TECHNICAL UNIVERSITY OF MUNICH

Bachelor's Thesis in Information Systems

**Evaluating and Enhancing Location-Aware  
Visual Document Segmentation for Oncology  
Guidelines**

**Matteo Felipe Merz**



SCHOOL OF COMPUTATION, INFORMATION  
AND TECHNOLOGY - INFORMATICS

TECHNICAL UNIVERSITY OF MUNICH

Bachelor's Thesis in Information Systems

# **Evaluating and Enhancing Location-Aware Visual Document Segmentation for Oncology Guidelines**

## **Evaluierung und Erweiterung positioneller visueller Dokumentensegmentierung für Onkologie-Richtlinien**

Author:	Matteo Felipe Merz
Supervisor:	Prof. Dr. Florian Matthes
Advisor:	Jonas Gottal, M.Sc.; Juraj Vladika, M.Sc.
Submission Date:	22.02.2026

I confirm that this bachelor's thesis in information systems is my own work and I have documented all sources and material used.

---

Location, Submission Date

---

Author

# AI Assistant Usage Disclosure

## Introduction

Performing work or conducting research at the Chair of Software Engineering for Business Information Systems (sebis) at TUM often entails dynamic and multi-faceted tasks. At sebis, we promote the responsible use of *AI Assistants* in the effective and efficient completion of such work. However, in the spirit of ethical and transparent research, we require all student researchers working with sebis to disclose their usage of such assistants.

For examples of correct and incorrect AI Assistant usage, please refer to the original, unabridged version of this form, located at [this link](#).

## Use of *AI Assistants* for Research Purposes

**I have used AI Assistant(s) for the purposes of my research as part of this thesis.**

Yes      No

**Explanation:**

I confirm in signing below, that I have reported all usage of AI Assistants for my research, and that the report is truthful and complete.

---

Location, Date

---

Author

## Acknowledgments

# Abstract

# Kurzfassung

# Contents

<b>Acknowledgments</b>	<b>iv</b>
<b>Abstract</b>	<b>v</b>
<b>Kurzfassung</b>	<b>vi</b>
<b>1. Introduction</b>	<b>1</b>
1.1. Problem Statement . . . . .	1
1.2. Objectives . . . . .	2
<b>2. Foundations</b>	<b>4</b>
2.1. Oncology guideline documents . . . . .	4
2.2. Vision-Language Models . . . . .	4
2.3. Document Parsing . . . . .	4
2.3.1. Modular Pipeline Systems . . . . .	5
2.3.2. End-to-End VLM models . . . . .	6
2.4. Document Chunking . . . . .	7
2.5. Retrieval-Augmented Generation . . . . .	7
2.5.1. Tokenization . . . . .	7
2.5.2. Text Embeddings . . . . .	7
2.5.3. Semantic Similarity . . . . .	8
2.5.4. Architecture of Retrieval-Augmented Generation systems . . . . .	8
2.5.5. Indexing . . . . .	9
2.6. Source Attribution . . . . .	9
2.6.1. Bounding Boxes . . . . .	9
2.6.2. Intersection over Union . . . . .	9
<b>3. Methodology</b>	<b>10</b>
3.1. Pipeline overview . . . . .	10
3.2. Data representations . . . . .	11
3.2.1. ParsingBoundingBox . . . . .	11
3.2.2. ParsingResultType . . . . .	11
3.2.3. ParsingResult . . . . .	12
3.2.4. ChunkingResult . . . . .	12
3.2.5. Chunk . . . . .	12
3.3. Parsing module . . . . .	13
3.3.1. Unstructured.io . . . . .	13



3.3.2.	Docling . . . . .	14
3.3.3.	MinerU . . . . .	14
3.3.4.	Gemini 2.5 Flash . . . . .	15
3.3.5.	LlamaParse . . . . .	15
3.3.6.	Google Document AI Layout Parser . . . . .	15
3.4.	Post-processing methods . . . . .	15
3.5.	Chunking module . . . . .	15
3.5.1.	Fixed-Size Chunking . . . . .	17
3.5.2.	Recursive Character Chunking . . . . .	17
3.5.3.	Breakpoint-based Semantic Chunking . . . . .	17
3.5.4.	Hierarchical Chunking . . . . .	18
3.6.	Evaluation Framework . . . . .	18
3.6.1.	Document Layout Analysis evaluation . . . . .	18
3.6.2.	Content Parsing Evaluation . . . . .	19
3.6.3.	Chunking Evaluation . . . . .	20
<b>4.</b>	<b>Results</b>	<b>21</b>
<b>5.</b>	<b>Discussion</b>	<b>23</b>
<b>6.</b>	<b>Conclusion</b>	<b>24</b>
<b>A.</b>	<b>General Addenda</b>	<b>25</b>
	<b>List of Figures</b>	<b>28</b>
	<b>List of Tables</b>	<b>29</b>
	<b>Acronyms</b>	<b>30</b>
	<b>Bibliography</b>	<b>31</b>

# 1. Introduction

## 1.1. Problem Statement

Clinical practice guideline documents are fundamental to efficient and reliable treatment of various illnesses [1]. Oncology guidelines are a subgroup of these documents, revolving around the treatment of various forms of cancer [2]. Guideline documents not only aid doctors in deciding on the optimal treatment options but also support patients in understanding their illness. In recent years, due to advancements in technology, novel therapeutics, and personalized medicine, clinical guidelines have drastically increased in size and complexity [3]. As of 2019, the average oncology guideline published by the National Comprehensive Cancer Network (NCCN) was 198 pages long, showing an annual increase of 7.5 percent over the previous 23 years [3]. This increase of complexity forces medical personnel to invest more time in order to be able to provide optimal care for cancer patients. Especially for individual practitioners this additional strain might become unsustainable if complexity continues increasing [3].

The Aidvice project proposes to address this problem by leveraging recent advantages in artificial intelligence (AI). Specifically, the project revolves around the development of a retrieval-augmented generation (RAG) based knowledge assistant [4]. RAG is an emerging paradigm which addresses a fundamental problem of traditional large language models (LLMs) [5]. While LLMs excel at many natural language processing (NLP) tasks, they are prone to ‘hallucinations’ and inaccurate answers, when sought information goes beyond the model’s training data [6]. This provides a major obstacle for the usage of LLMs in the medical field, where accurate and reliable answers are of the highest priority [7]. RAG mitigates these drawbacks by retrieving additional context from an external knowledge source which the LLM can take advantage of during answer generation [8, 5].

The efficiency of such a RAG system is fundamentally constrained by the quality and relevance of the context retrieved from the knowledge base [9, 10]. The quality of the construction of the knowledge base directly determines whether the correct context can be retrieved for a given query [6]. Therefore, the construction of the knowledge base out of the oncology guidelines is a critical aspect of the project. Additionally, as the project has a clear focus on verifiability and traceability, there is an additional requirement to provide visual source attribution. This means that retrieved passages need to include accurate positional information, giving visual confirmation to the practitioner about the origin of the retrieved context [11].

The guideline documents are stored in the unstructured portable document format (PDF). In order to be further processed for the knowledge base, they first need to be transformed into a machine-readable structured data format through a process called document parsing

(DP) [12]. The inherent structure of the guidelines poses multiple challenges for this process, such as complex tables, varying layouts and occasional formatting errors.

As LLMs are constrained by the size of their context window, it is not feasible to store the entire guideline documents as individual entries in the knowledge base [6]. Therefore, the documents need to be split up into smaller text chunks that fit into the model's context window [6]. This process is called document chunking [13].

During retrieval the model identifies the most relevant passages in the knowledge base based on their semantic similarity to the user's query [5]. If the stored text chunks are too long, important information might be lost between irrelevant details [9, 10]. On the other hand, storing too short text chunks can result in important statements being broken up into multiple chunks and losing their meaning. In order to maximize the quality of the retrieved chunks, both the chunk size as well as the document chunking strategy, used to decide where to split up the oncology guidelines, need to be optimized [10].

Additionally, established implementations of popular document chunking techniques do not fulfill the requirement of visual source attribution at the granularity required by the Aidvice project [14, 15, 16, 17]. Therefore there is a need for the development of a novel solution, that addresses this issue.

## 1.2. Objectives

This study addresses three fundamental research questions regarding the data preparation for a RAG based knowledge assistant for oncology guidelines. Each of the following research questions addresses a specific aspect of the evaluation and improvement of the document segmentation process required for the construction of the knowledge base.

- **RQ1:** How are the challenges introduced by oncology guidelines reflected in established benchmarks for document parsing?
- **RQ2:** Which metrics are most useful to measure the effectiveness of document parsing and chunking methods?
- **RQ3:** How can current segmentation methods be adapted or expanded on to fulfill the requirements of a RAG based knowledge assistant with visual source attribution?

### [(TODO: explain title)]

To identify the challenges posed by the oncology guidelines to the DP process, we perform a qualitative analysis identifying the characteristics of oncology guidelines relevant during the DP process, such as their formatting, layout and common types of structural elements. We then identify established document benchmarks and datasets which contain documents that most closely resemble these characteristics. Through this analysis, we underline the transferability of results achieved on these datasets to our application, while identifying unrepresented characteristics which require manual comparisons. This approach allows the evaluation of various DP techniques on established benchmarks, without the availability of a dedicated oncology guideline benchmark.

In order to evaluate the effectiveness of both document parsing and chunking techniques, we identify various metrics used in existing literature. We then perform a comparative analysis of the identified metrics, evaluating their suitability for our application. Based on this analysis, we select a set of metrics which are most suitable for the evaluation of various document parsing and chunking techniques.

Finally, we propose a novel solution for the visual source attribution requirement of the Aidvice project. We adapt and expand on existing document chunking techniques in order to provide accurate positional information for each text chunk. By introducing a universal dataformat for the output of the DP implementations, we enable the direct comparison of various document parsing techniques using the benchmarks and metrics identified in RQ1 and RQ2. Through this evaluation we identify promising combinations of document parsing and chunking techniques for the creation of the knowledge base of the Aidvice project, while providing a modular framework for future experiments and improvements.

## 2. Foundations

### 2.1. Oncology guideline documents

Qualitative analysis of the properties of oncology guideline documents.

### 2.2. Vision-Language Models

LLMs are inherently confined to processing exclusively text-based data. This limitation restricts their applicability in complex, real-world scenarios, where understanding and combining data from multiple modalities is crucial [18, 19]. Vision-language model (VLM) are a class of models which respond to these limitations by combining visual and textual processing capabilities into a single architecture [18]. These models find applications involving both the comprehension and generation of multimodal content, such as image captioning, and visual question answering [18].

The typical architecture of VLMs is depicted in [\[\*\*\(TODO: figure\)\*\*\]](#). The architecture consists of a text encoder, an image encoder, an image-text

[\[\*\*\(TODO: Maybe need to add explanations for transformers?\)\*\*\]](#)

### 2.3. Document Parsing

Also known as document content extraction, DP aims to convert unstructured and semi-structured documents into structured, machine readable data formats [12, 20]. During this process elements such as headings, tables, and figures are extracted from the document while preserving their structural relationships. DP is crucial for many document-related tasks, providing access to previously unavailable information sources. Especially for LLMs, where leveraging additional training data is crucial for enhancing the model's factual accuracy and knowledge grounding, DP plays an important role [20, 21]. With the emergence of the RAG paradigm, DP has also been critical in the creation of the knowledge database, as important information is often stored inside file formats which can not directly be processed by machines [22]. While DP is used for converting a range of document formats into machine-readable content, we will focus solely on the parsing of PDF documents for the purposes of this thesis, as this is the datatype that the oncology guidelines are stored as.

Converting PDF documents is particularly challenging due to their variable formatting, lack of standardization and focus on visual characteristics [22]. The format not only includes digitally-born files but also includes photographed and scanned documents. Therefore DP

systems need to be able to adapt to a wide range of different layouts, image qualities and document types, such as academic papers, invoices or presentation slides [20, 23]. While there are many tools and implementations available for DP [22, 21, 24, 25], most of them can be categorized into either modular pipeline systems or end-to-end VLM models.

### 2.3.1. Modular Pipeline Systems

Modular pipeline systems employ various different modules in a sequential order to perform DP. This modular design enables the targeted optimization of individual components and flexibility integration of new modules and techniques [26]. Additionally, by making use of lightweight models and integrating parallelization, pipeline systems can reach efficient parsing speeds [20]. While different formations are possible, most implementations consist of three different stages [12].

**Document Layout Analysis (DLA):** According to Q. Zhang, B. Wang, V. S.-J. Huang, et al. [12], DLA refers to the identification of the structural elements of a document, such as paragraphs, section headers, tables, figures, and mathematical equations, as well as their respective bounding boxes [12, 27]. There are two types of methods for performing DLA. Unimodal methods focus purely on visual features of the document in order to identify structural elements [12, 28]. Notably, convolutional neural networks (CNN)- and transformer-based methods adapt models initially designed for object detection tasks, such as the YOLO [29] and DETR [30] families of models, to accurately identify structural elements in document images [12, 27]. Hereby, transformer-based methods excel at capturing global relationships between structural elements at the cost of computational intensivity and expensive pretraining [12]. The second type of DLA methods are multi-modal methods. Additionally to the visual representations, multi-modal methods also rely on the content and position of the pages' textual content, performing DLA using a VLM [28, 31]. This approach allows more granular classifications and the analysis of highly complex layouts [12, 28].

**Content Extraction:** To extract the content of the identified structural elements different recognizers are applied to the element's region based on its classification [12, 21, 22]. For textual elements, such as paragraphs or section headings, the textual content is identified using optical character recognition (OCR). OCR engines use techniques from computer vision in order to identify and extract text from images [12, 32]. Popular OCR engines include EasyOCR [33] and the Tesseract OCR engine [34]. Additionally to extracting content using OCR, DP implementation often provide specific recognizers for additional element types [12, 21]. Most commonly this includes a specific model for table structure recognition, referring to the extraction of table content as a structured file format, such as HTML, XML or Markdown [12, 22, 21, 35]. Other options for class-specific recognizers include mathematical formula recognition and chart recognition [21, 24, 12].

**Relation Integration:** During relation integration the identified elements are combined into the final output format. During this stage, rule-based methods and specialized AI models may be employed, for example for filter out duplicate or unwanted elements or correcting the reading order of the document [12, 21, 22]. Depending on the chosen output format, this process might lead to the loss of information, such as the loss of bounding box information

for an output in Markdown format [16].

Systems following the modular pipeline approach also have some inherent drawbacks. Mainly, due to handling the parsing of each structural element independently of each other, pipeline systems fail to capture information about the global context of the document, leading to semantic loss [23]. Additionally, because of the sequential nature of the pipeline approach, errors from different stages propagate through the pipeline [23, 24].

### 2.3.2. End-to-End VLM models

Due to recent recent advancements in VLM architectures, end-to-end VLM models have emerged as a promising alternative to traditional pipeline-based approaches. Research such as General OCR Theory (GOT) have demonstrated the ability of VLMs to perform OCR with a high accuracy, while being able to extract tables, charts or mathematical formulas with a singular model [36]. Contrary to pipeline-based methods, VLM-based approaches are able to generate structured outputs directly from the input document, addressing the error propagation problem of modular pipelines [26]. Additionally, these models demonstrate advantages in understanding the structure and hierarchy of complex documents [12]. VLM-based approaches can be divided into two further subcategories:

**General-Purpose VLMs:** General purpose VLMs are not trained solely for document-centric tasks, but are still able to show promising results for DP, due to their large parameter count and extensive training data [24, 19]. However, these models are often either proprietary or require extensive computational resources [24]. Additionally, they often struggle with documents that follow more complex layouts or contain densely packed text blocks [24].

**Domain-Specific VLMs** Domain-specific VLMs are trained and optimized specifically for optimal DP [24, 12, 23]. In recent years, there has been promising developments towards domain-specific VLMs that encapsulate DLA, content extraction and relation integration into a single model [37]. These models are able to achieve state-of-the-art performance on document parsing benchmarks, while being a fraction of the size of general-purpose VLMs [37, 24]. However, as VLMs are not bound to the stages of traditional pipeline systems, there has also been additional research regarding models optimized for the direct generation of structured outputs, most notably Markdown [23]. However, this approach inherently leads to the loss of information, such as positional information for the extracted elements, which inherently are not included in the Markdown output, making this class of models unsuitable for the purposes of this research [37, 23].

Recently, there has also been research towards multi-stage VLM-based approaches [24, 26]. These models use one or more VLMs in multiple stages, aiming to encapsulate the computational efficiency of pipeline approaches with the improved accuracy and structure understanding of VLM-based methods [24]. However, especially when multiple VLMs are in use, these approaches come with a further increase in complexity and computational requirements and may show suboptimal performance in tasks such as reading order inference compared to single-stage VLM-based approaches [24, 37]. Current challenges regarding the development of VLM-based approaches are the risks of ‘hallucinations’, especially on longer documents [24, 16], as well as their high computational requirements and compared to

modular pipeline systems [16].

## 2.4. Document Chunking

[(TODO: Other usages of chunking than RAG)]

Content:

- Definition
- Role in RAG systems
- Different Types: Rule based, File dependent

## 2.5. Retrieval-Augmented Generation

While LLMs have extensive general domain knowledge due to their enormous corpora of training data, compiled from various open-domain sources [38], they struggle with tasks that require domain-specific knowledge which they did not encounter during training [6]. This can lead to ‘hallucinations’ and inaccuracies, as the model tries to synthesize a matching answer based on its domain-wise irrelevant training data [6, 8]. RAG addresses this limitation, extending the usage of LLMs to applications requiring extensive knowledge in a specific domain [5]. This is achieved by retrieving information from an external knowledge source comprised of application-relevant text passages, supplying additional context to the LLM during answer generation [5, 6].

### 2.5.1. Tokenization

Tokenization refers to the segmentation of text into subword units called tokens [39]. Tokens are the fundamental text representation for most NLP tasks. With a granularity located between characters and words, tokens can retain linguistic meaning while also being able to represent arbitrary text with a relatively concise vocabulary [39]. Using tokenization any given text can essentially be represented as a list of integers, with each integer being the identifier to a specific token in the tokenizer’s dictionary [40]. During training, the tokenizer creates its dictionary by finding character pairings that occur with the highest frequency in the training data [40]. Additionally, with the multitude of different techniques for modern subword tokenization [41, 42, 43], the same input text can lead to drastically different outputs depending on the specific tokenizer and training data. Therefore, tokenizers always need to match the NLP models they are used with.

### 2.5.2. Text Embeddings

Content:

- Definition



- Types of embeddings

### 2.5.3. Semantic Similarity

Content:

- Definition
- Role in RAG

### 2.5.4. Architecture of Retrieval-Augmented Generation systems

While there are many advanced and extended versions of RAG systems, for this study we will focus on the standard Naive RAG architecture as depicted in Figure 2.1 [6]. Naive RAG is based on the original RAG architecture proposed by Lewis, Perez, Piktus, et al. [5]. Naive RAG systems consist of two modules:

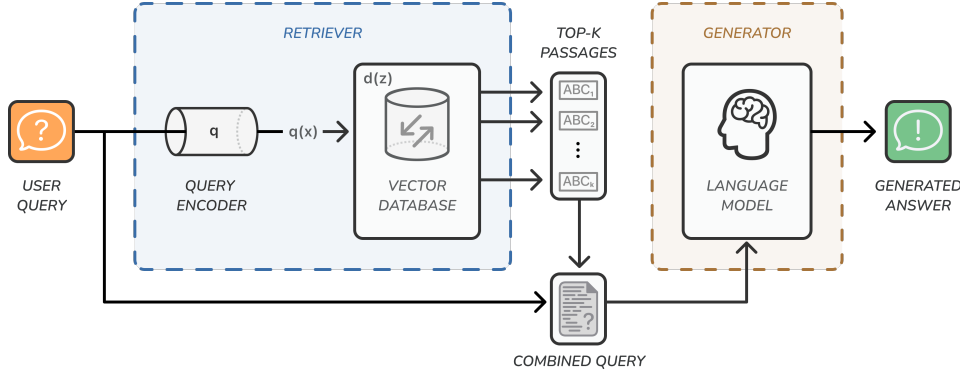


Figure 2.1.: Architecture of the Naive RAG system.

**Retriever:** The retriever module consists of a query encoder and an external knowledge base [5]. It is responsible for retrieving relevant context from the knowledge base, based on the user’s query [6]. The module is based on the bi-encoder architecture, with the query encoder  $q$  and document encoder  $d$  encoding texts into a shared embedding space [5, 44]. The knowledge base is a vector database consisting of application-specific text passages  $z$ . Each passage is stored in the database as a vector embedding  $d(z)$ , encoded through the document encoder  $d$  [5]. To identify the relevant passages for a query  $x$ ,  $x$  is first transformed into a vector embedding  $q(x)$  using the retriever’s query encoder [6]. Based on the similarity scores between the query embedding and the stored chunk embeddings, the top- $k$  documents  $z$  with the highest similarity scores, are then retrieved from the database [5, 6].

**Generator:** The generator module is responsible for synthesizing the final answer based on the user’s query and the passages retrieved by the retriever [5]. Firstly, the original query  $x$  and the retrieved passages  $z$  are combined into a single input query [6]. The LLM is then tasked with generating the final answer  $y$ , conditioned on this combined input [6, 5].

### 2.5.5. Indexing

In order to apply RAG systems to knowledge-intensive tasks in a specific domain, the external knowledge base needs to be created from relevant data sources. This process is referred to as Indexing [6]. Indexing begins with the preparation of the data sources into short text passages [6]. For the purpose of this study we will refer to this process as document segmentation. Document segmentation includes both DP, the conversion of unstructured documents, such as PDFs and images, into structured data, as well as chunking, the splitting of this data into smaller text passages called chunks. Chunking is a necessary step for RAG systems, as both LLMs and encoders are limited in the number of tokens that fit into their context window [6]. Furthermore, indexing includes the encoding of these chunks into vector embeddings. Both the embeddings and the original chunks are then stored as key-value pairs in a vector database, allowing fast and frequent searches during retrieval [6].

## 2.6. Source Attribution

### 2.6.1. Bounding Boxes

- Definition
- Formats of Bounding boxes
- What do we mean when we say bounding boxes in the thesis? (normalized, ltrb)

### 2.6.2. Intersection over Union

According to the definition from Kaur and Singh [45], the intersection over union (IoU) between two bounding boxes  $BB_a$  and  $BB_b$  is defined as described in Equation 2.1. The IoU can take on any value between 0 and 1, where a value of 0 means that there is no overlap between the two bounding boxes, and a value of 1 means that the two bounding boxes are identical. In the context of object detection, IoU is commonly used to evaluate the accuracy of predicted bounding boxes against ground truth bounding boxes [45].

$$IoU(BB_a, BB_b) = \frac{\text{Area of intersection of } BB_a \text{ and } BB_b}{\text{Area of union of } BB_a \text{ and } BB_b} \quad (2.1)$$

## 3. Methodology

### 3.1. Pipeline overview

In order to be able to compare different DP implementations and chunking strategies, we developed a modular document segmentation pipeline. The system transforms the raw PDF documents into chunks through a two-stage process. Due to this modular approach, the pipeline allows to easily swap out both the used DP implementation and chunking strategy, while maintaining a unified interface for both modules. The architecture of the data processing pipeline is illustrated in Figure 3.1.

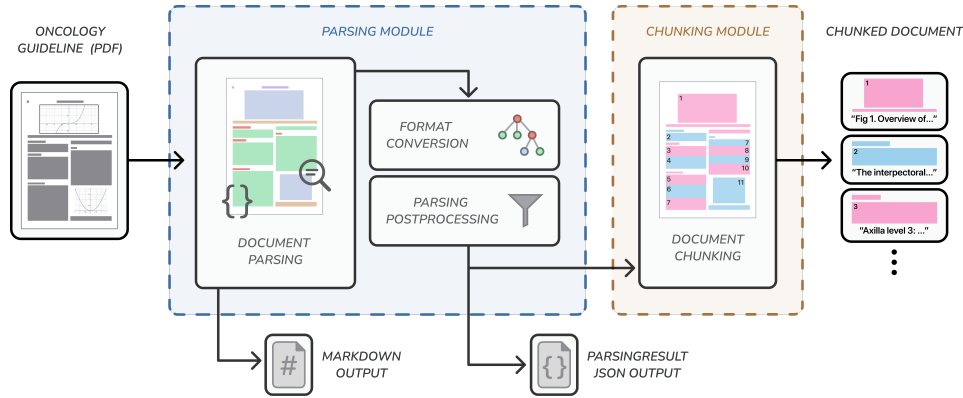


Figure 3.1.: Overview of the document segmentation pipeline architecture. The pipeline consists of two main modules: the parsing module and the chunking module. The parsing module transforms unstructured PDF documents into the structured ParsingResult format while extracting its textual content as a Markdown file. The chunking module then splits the ParsingResult into smaller chunks suitable for RAG applications, while maintaining positional information for visual source attribution.

**Parsing module:** The parsing module is the first step of the pipeline. It provides a unified interface for eight different DP implementations. Firstly, document elements and their structural information are extracted from the document using the underlying DP implementation. The output then gets converted into a standardized data format and undergoes further post-processing steps, such as filtering unwanted element types, to prepare the data for the chunking module. Additionally, the structured data representation of the document gets persisted to the file system as a lossless JavaScript Object Notation (JSON) serialization as well as a lossy serialization to Markdown format.

**Chunking module:** The chunking module is responsible for the splitting of the structured data into smaller chunks using one of multiple available document chunking strategies. We adapt four established strategies to provide accurate positional information in the form of bounding boxes, enabling visual source attribution for each chunk. To achieve this we propose a novel approach to the document chunking paradigm, enabling traceability of the chunk’s content on the token-level. Through this approach, we are able to determine more accurate chunk bounding boxes compared to established methods.

## 3.2. Data representations

As of the time of writing, every DP implementation defines their own datatypes, making comparisons very complex. In order to consolidate multiple different DP implementations into our modular pipeline and evaluate them against each other, the output of each implementation needs to be transformed into a uniform format. For this purpose, we define our own universal datatypes to be used in the data processing pipeline.

### 3.2.1. ParsingBoundingBox

The ParsingBoundingBox (Figure 3.2) datatype adds additional fields to a traditional bounding box tuple. Since the guideline documents usually contain multiple pages, we include a page field indicating the page number that the bounding box belongs to. The coordinates of the bounding box are saved as a normalized tuple in the  $(x_{\min}, y_{\min}, x_{\max}, y_{\max})$  format. Additionally, the spans field allows for saving more granular bounding boxes such as boxes for every line of text inside the bounding box.

```
class ParsingBoundingBox:
    page: int
    left: float
    top: float
    right: float
    bottom: float
    spans: list[ParsingBoundingBox]
```

Figure 3.2.: Python implementation of the ParsingBoundingBox datatype.

### 3.2.2. ParsingResultType

One problem with comparing multiple DP methods is their lack of universal categories. This leads to two specific issues.

Firstly, the naming of element types differs between implementations. For example, a paragraph gets classified as NarrativeText by the Unstructured.io framework [25], while Docling [22] names the same category as simply TEXT.

Secondly, some methods provide classifications, which are not provided by others. One example for this is the addition of a `ref_text` element type in the MinerU implementation [21, 46], referring to an entry in a bibliography. For our purposes, we aggregated all possible categories from the tested implementations and provided mappings to the universal types for each DP implementation. The full list of `ParsingResultTypes` can be found in Table A.

### 3.2.3. ParsingResult

The `ParsingResult` (Figure 3.3) is the output of the parsing module. It is inspired by the `DoclingDocument` datatype from the Docling library [22]. Mainly, `ParsingResult` adopts the tree structure proposed by multiple DP implementations to represent hierarchical relationships between document elements.

[(TODO: Explain difference to doclingdoc: no different classes, easier bounding boxes, only tree structure with no group or text lists)]

```
class ParsingResult:
    id: str
    type: ParsingResultType
    content: str
    geom: list[ParsingBoundingBox]
    parent: ParsingResult | None
    children: list[ParsingResult]
    metadata: dict
    image: str
```

Figure 3.3.: Python implementation of the `ParsingResult` datatype.

The `geom` attribute contains a list of bounding boxes related to the element. This allows elements, such as tables, to span multiple pages by containing multiple bounding boxes. Using the `parent` and `children` attributes of the `ParsingResult`, the hierarchy of the input document is represented as a tree of `ParsingResult` nodes. This way, hierarchical relationships between the elements can be represented. For example, a section header has all text paragraphs contained in it's section as children.

### 3.2.4. ChunkingResult

The `ChunkingResult` (Figure 3.4) is the output of the chunking module. It provides a wrapper around a list of `Chunks`, adding a `metadata` field for information about the document and the preceeding processes.

### 3.2.5. Chunk

[(TODO: write)]

```
class ChunkingResult:
    chunks: list[Chunk]
    metadata: dict

class Chunk:
    id: str
    content: str
    metadata: dict
    geom: list[ParsingBoundingBox]
```

Figure 3.4.: Python implementation of the ChunkingResult and Chunk datatypes.

### 3.3. Parsing module

**[(TODO: rewrite)]** The parsing module processes an incoming PDF document in a multi-step process. Firstly, the document is passed to the internal DP implementation using the `_parse` function. Through the `_transform` function, the raw output is then converted into the `ParsingResult` format. Additionally, the `_get_md` function extracts the Markdown representation of the document from the raw output of the internal parser. Finally, additional metadata about the parsing process, such as the duration of the internal parsing operation, are added to the root `ParsingResult`. After performing additional post-processing steps, both the Markdown as well as the lossless serialization of the `ParsingResult` in JSON output are then persisted to the file system. As the raw types returned by the internal parser differ based on the specific implementation, each of the abstract functions mentioned above need to be implemented by the specific parser class. An overview of the abstract functions to be implemented can be found in Table 3.3.

Function	Input	Output
<code>_parse</code>	PDF document	Custom data format
<code>_transform</code>	Custom data format	<code>ParsingResult</code>
<code>_get_md</code>	Custom data format	Markdown string

Table 3.1.: Abstract functions of the Parsing Module to be fulfilled by the implementation.

#### 3.3.1. Unstructured.io

Unstructured.io is a prominent provider for DP, offering both a cloud-based API as well as an open-source library. For our study, we will focus on the open-source library version of Unstructured.io [25, 35]. While the developers themselves explicitly highlight that the open-source library is not suited for large-scale production environments [35], its inclusion within the documentation of popular RAG frameworks, such as Langchain [14] and LlamaIndex [15] make it a popular choice for a first point of contact with DP. Therefore, we will regard the open-source library as a baseline for the compared implementations. Unstructured.io follows a modular pipeline approach. Specifically, the implementation uses YOLOX, a unimodal vision transformer, to perform DLA [35, 47]. The library also includes a specialized model for table structure recognition [35].

### 3.3.2. Docling

Docling, which was developed by IBM in 2022, is one of the most popular available open-source DP libraries [22, 16]. Docling particularly stands out from other DP implementations through its permissive MIT license. To achieve this, Docling relies primarily on custom models instead of using third-party software, which are often not as permissive [22]. Docling offers two different approaches for DP:

**Parsing pipeline:** Docling’s processing pipeline consists of three components: a PDF backend called DoclingParse, an internal model pipeline containing multiple AI models, and a post-processing stage. Firstly, the PDF backend extracts useful information from the document using both contained programmatic information as well as OCR techniques. This includes bounding boxes for every text element inside the document. The internal model pipeline then performs both the DLA as well as content extraction steps. Hereby, Docling provides their own models for table structure recognition with the TableFormer model [48] as well as for DLA with their Heron model [27]. Heron is derived from RT-DETR [49], a unimodal vision transformer, and retrained on DocLayNet [50], Docling’s own dataset for DLA. During DLA identified bounding boxes are compared and intersected with bounding boxes retrieved from the PDF backend in order to provide more accurate localization [16]. During post-processing the recognized elements are then combined into the DoclingDocument datatype [16].

**Granite Docling:** Granite Docling is an end-to-end VLM for DP. It belongs to the group of domain-specific VLM models, specifically build for document understanding and conversion [51]. The model is very compact, consisting of around 258 million parameters [52, 51]. With this model, Docling proposes the DocTags data format, a structured format designed for representing both text and structure of the document through extensible markup language (XML)-style tags [52].

### 3.3.3. MinerU

Another popular choice for open-source on-device DP is the MinerU framework. Similar to Docling, MinerU also offers both a pipeline as well as a VLM-based approach for DP [21, 24, 46].

**Parsing pipeline:** MinerU extends the traditional processing pipeline through a pre- and postprocessing stage. In the preprocessing stage unprocessable files are filtered out and metadata about the document is extracted using the PyMuPDF library [21, 53]. This metadata includes the language of the document, the document’s page dimensions and the identification of scanned documents [21]. The pipeline then uses models from the DP model library PDF-Extract-Kit for DLA and content extraction [21, 46, 54, 55]. For content extraction, special models for formula and table recognition are employed by the pipeline. The model used for DLA is a finetuned version of LayoutLMv3, a multi-modal model [21, 31]. During the final postprocessing stage, overlapping elements are cleaned up and the reading order of the document elements is inferred using a segmentation algorithm [21].

**VLM:** With MinerU2.5, the implementation’s offerings were expanded by a multi-stage

VLM-based DP approach. This approach employs a 1.2 billion parameter VLM to perform DP in a two-stage approach [24]. Firstly, the model is used to perform DLA on the document, identifying elements and their reading order. In the second stage, the same model is applied again on individual image crops of the page element and is tasked to extract the content from the crop [24].

#### 3.3.4. Gemini 2.5 Flash

Gemini 2.5 Flash is a closed-source proprietary model developed by Google with strong multi-modal capabilities across text, vision and audio [56]. While Google offers a more capable model in the form of Gemini 2.5 Pro, we follow the sentiment from Niu, Z. Liu, Gu, et al. [24], that DP tasks “typically exhibit relatively low dependency on large-scale language models” (p.7) and both models similar results on various image understanding benchmarks [56], instead opting to rely on the cheaper, faster Gemini 2.5 Flash model for our study. Gemini 2.5 Flash belongs to the group of general-purpose VLMs and, due to its closed-source nature, is only accessible through an application programming interface (API). The Gemini family of models received additional training in order to provide improved accuracy on object detection and image segmentation tasks [56, 57]. We follow the documentation provided by Google on harnessing Gemini’s image understanding capabilities [57] to formulate a prompt, that takes advantage of this additional training for the DP task. The full prompt is available in A.

#### 3.3.5. LlamaParse

#### 3.3.6. Google Document AI Layout Parser

### 3.4. Post-processing methods

### 3.5. Chunking module

The chunking module provides an interface to perform document chunking on the ParsingResult tree representation of the PDF document. We propose a novel solution aimed at increasing the traceability of the chunk content to its constituent ParsingResults. Prominent implementations of document chunking methods, such as the ones found in the RAG frameworks LlamaIndex [15] and Langchain [14], treat the chunking process as a transformation of a single string, typically the Markdown representation of the document, into multiple smaller strings. Following this approach, resulting chunks lose their direct connection to the underlying structural elements, complicating source attribution.

Domain-specific implementations, such as Docling’s Hybrid and Hierarchical Chunkers [22, 16], improve upon this by including a list of elements associated with the resulting chunk. Additionally, they take advantage of the document’s hierarchy by including relevant section headers in the chunk text. However, these implementations do not distinguish between partially and fully included elements. This results in chunk bounding boxes which always



contain the entire element, regardless of how much of it's text is actually included in the text of the chunk. Furthermore, while the content of the section headers is included in the chunk, their corresponding bounding boxes are not retained.

Our solution addresses these limitations by enabling traceability at the highest relevant granularity for RAG chunking. While in traditional text segmentation, characters are the smallest instance that a text can be broken up at, the same can not be said about AI models. AI models are not directly constrained by the character length of their input and are instead limited by the amount of tokens that can fit into their context window. This means that the smallest instance that a text should be broken up into in the context of RAG is a token. To enable the tracing of every token inside the chunk to a token inside an element, we introduce the RichToken. A RichToken (Figure 3.5) is an object that wraps around a single token and it's original text, the id of the element that the token belongs to, and it's corresponding index inside that element. When the module performs the chunking operation on a ParsingResult, the specific chunking implementation processes the tree and returns the content of each chunk as a list of RichTokens. The module then creates a Chunk object using the RichTokens. By aggregating the content of the RichToken, the module can reconstruct the content of the entire Chunk, while maintaining the relationships to the constructing ParsingResults. Since the module is aware of the specific tokens that are included in the Chunk, we can determine tighter chunk bounding boxes using the span-level bounding boxes added to the elements in the DP post-processing.

```
class RichToken:
    element_id: str
    token_idx: int
    token: int
    text: str
```

Figure 3.5.: Python implementation of the RichToken datatype.

How the RichTokens are grouped together into chunks is to be decided by the specific chunking strategy. The module provides a `_tokenize` function, which transforms a given ParsingResult node into RichTokens using the `all-MiniLM-L6-v2` sentence transformer as a tokenizer. To prevent distinct document elements from merging into a single text block during chunk creation, `\n` delimiters are appended to the content of the ParsingResult nodes before the chunking process begins. These delimiters act as textual representations of the breakpoints between document elements. Additionally, to avoid creating chunks, which are too big for the embedding module's context window, a limit  $N$  for the maximum amount of tokens per chunk can be set on the chunking module. We provide implementations for four different document chunking strategies.

### 3.5.1. Fixed-Size Chunking

Fixed-size chunking segments the document into chunks of the chunking module’s maximum chunk length  $N$  with an overlap of  $O$  tokens. It provides a simple and computationally efficient way to perform document chunking, while producing chunks of a constant size [13]. However, fixed-size chunking has no regard for the content of the document, which may result in chunk borders inside a single word or sentence [58].

The strategy traverses the `ParsingResult` tree in reading order (e.g., depth-first), creating a queue of the document’s `RichTokens` in the progress. When the queue reaches a length larger than  $N$ , the first  $N$  tokens inside the queue are combined into a chunk while  $N - O$  tokens are removed from the queue. This sliding window approach, leaving  $O$  tokens inside the queue after the chunk is created, aims to maintain some contextual continuity between the chunks, leading to improved recall during the retrieval phase [13, 58].

### 3.5.2. Recursive Character Chunking

Recursive character chunking encapsulates a similar approach to fixed-size chunking. However, instead of naively setting the border of a chunk at a fixed token count, recursive character chunking utilizes an hierarchical list of delimiters (e.g., paragraphs, sentences, words) to define chunk boundaries [58, 17, 14]. When the `RichToken` queue exceeds the maximum chunk length  $N$ , the strategy splits the tokens using the highest priority delimiter. If there still exists a split which is larger than  $N$ , the process recurses on the oversized split with the next delimiter in the list. This ‘coarse-to-fine’ approach preserves logical groupings while avoiding unnecessary fragmentation [58]. Similar to fixed-size chunking, recursive character chunking also incorporates sliding window chunking with an overlap of  $O$  tokens between adjacent chunks. The delimiters used in this implementation are adapted from `LangChain`’s `RecursiveCharacterTextSplitter` [14, 17], with punctuation added to better identify sentence endings, as suggested by B. Smith and Troynikov [59]. This results in the following delimiters: `[\n\n\n", "\n\n", "\n.\n", "\n!\n", "\n?\n", "\n \n", "\n"]`.

### 3.5.3. Breakpoint-based Semantic Chunking

Breakpoint-based semantic chunking separates the document at the sentence level, inserting breakpoints in between sentences to denote chunk borders [13]. Instead of relying on

For this strategy, the semantic distances between the embeddings of every adjacent sentence pair are calculated, with a high distance indicating a topical shift between the sentences. If the distance is larger than the  $Q$ -th percentile of all distances, a breakpoint is inserted.

### 3.5.4. Hierarchical Chunking

## 3.6. Evaluation Framework

### 3.6.1. Document Layout Analysis evaluation

The goal of the DLA evaluation is to assess the correctness of the bounding boxes and type labels produced by the parsing module [60]. While there are multiple datasets available for this task [61, 50, 20], we will use the PubLayNet dataset [62] for our evaluation. While many datasets focus on evaluating DLA on a range of different document types such as forms, invoices or handwritten documents, PubLayNet consists solely of medical scientific articles [20, 62]. This format closely resembles the format of the oncology guideline documents which makes it a suitable choice for this evaluation. Comprised of over 360.000 automatically annotated document pages collected from PubMed Central Open Access (PMCOA), PubLayNet is one of the largest datasets for DLA. As the dataset in its entirety is no longer publicly available and far too large for the purposes of this thesis, we will use *publaynet-mini*, a small subset of 500 pages of the original dataset for this evaluation [63]. As seen in Table 3.2, the subset contains around 5000 ground truth annotations for elements from 5 different classes.

Category	Annotations
Text	3,676
Title	1,000
List	73
Table	128
Figure	172
<b>Total</b>	<b>5,049</b>

Table 3.2.: Distribution of ground truth annotations across the different element types contained in the *publaynet-mini* subset of the PubLayNet dataset.

To assess the performance of different predictors on object detection tasks such as DLA, average precision (AP) is the most commonly used metric [64]. Previous evaluations of DLA models on the PubLayNet dataset also use a version of this metric [65]. However, AP relies on the predictor’s confidence values, indicating how confident the predictor is about a predicted bounding box and class label. As most of the DP implementations provide ‘hard-predictions’, which do not contain any confidence values, the AP is not a viable metric for the purposes of this thesis [66, 67].

For this reason, we will compare the implementations based on their achieved F1 score. Similarly to AP, this metric takes into account two important measures for object detectors: precision and recall [68]. According to Padilla, Passos, Dias, et al. [68], “Precision is the ability of a model to identify only relevant objects. [...] Recall is the ability of a model to find all relevant cases [...]” (p. 9). In order to calculate their values, firstly the detected bounding boxes (DTBBs) are classified into true positives (TPs) and false positives (FPs). A DTBB is

classified as a TP if there exists a ground truth bounding box (GTBB) from the same class, so that their IoU is greater than a given threshold. One GTBB can not be matched to multiple DTBBs. If there does not exist a GTBB that fulfils these criterions, the DTBB is classified as a FP. Any GTBBs which were not matched to a DTBB are classified as false negatives (FNs). Following the definition from Padilla, Passos, Dias, et al. [68] for a model that, on a dataset with  $G$  GTBBs, outputs  $N$  DTBBs, out of which  $S$ , ( $S \leq N$ ) are TPs, precision and recall can be formulated as shown in Equation 3.1 and Equation 3.2.

$$\text{Pr} = \frac{\sum_{n=1}^S \text{TP}_n}{\sum_{n=1}^S \text{TP}_n + \sum_{n=1}^{N-S} \text{FP}_n} = \frac{\sum_{n=1}^S \text{TP}_n}{\text{all detections}} \quad (3.1)$$

$$\text{Re} = \frac{\sum_{n=1}^S \text{TP}_n}{\sum_{n=1}^S \text{TP}_n + \sum_{n=1}^{G-S} \text{FN}_n} = \frac{\sum_{n=1}^S \text{TP}_n}{\text{all ground truths}} \quad (3.2)$$

The F1 score is the weighted harmonic mean between precision and recall and is calculated as defined in Equation 3.3 [68]. The F1 score is calculated for a single class at a set IoU threshold. Selecting a higher threshold will lead to a stricter metric as predictions need to be more precise to be counted as a TP [68]. A F1 score calculated at an IoU threshold  $T\%$  is commonly referred to as  $F1@T$  [69].

$$F_1 = 2 \frac{\text{Pr} \cdot \text{Rc}}{\text{Pr} + \text{Rc}} \quad (3.3)$$

For scenarios with multiple classes, such as the PubLayNet dataset, the Macro F1 score can be used to assess the overall performance of the predictor [70]. The Macro F1 score is the mean of the single class F1 scores. For a dataset with  $M$  different classes, the calculation of the Macro F1 score is described in Equation 3.4. Hereby,  $N_{:j}$  and  $G_{:j}$  denote the DTBBs and GTBBs belonging to elements of class  $j$  [70].

$$F_{1_{\text{Macro}}}(N, G) = \frac{1}{M} \sum_{j=1}^M F_1(N_{:j}, G_{:j}) \quad (3.4)$$

The DP implementations will be evaluated on both their single-class and Macro F1 scores. Specifically, their (Macro)  $F1@50$  and  $F1@50 : 95$  will be compared against each other.  $F1@50:95$  refers to the mean of the F1 values calculated at 10 evenly spaced IoU thresholds between 0.5 and 1.0 and is inspired by the primary challenge metric found in the MS COCO dataset [71]. This rewards implementations, which provide more accurate bounding boxes [71].  $F1@50$  is chosen, as a threshold of 50% is one of the most commonly used threshold values for metrics in object detection [68]. To calculate these metrics, the `faster-coco-eval` package is used to determine the recall and precision values at the IoU thresholds [72].

### 3.6.2. Content Parsing Evaluation

OmniDocBench Content:

- Data Creation

- Evaluation Modes (End2End *to* Evaluate Markdown output (Content))
- Usages in other papers
- State of the art evaluations (if I find any, most for all kinds of documents)
- Types of Documents (English and Chinese, different kinds: Scientific Paper. . . ) all single page
- Subset for the thesis: English Scientific Papers

#### 3.6.3. Chunking Evaluation

Chroma Evaluation

## 4. Results

	text	title	list	table	figure	all
unstructured_io	0.8123	0.8032	0.1269	0.9337	0.6138	0.6216
docling	0.8687	0.8775	<b>0.8022</b>	0.9530	0.5951	<b>0.8063</b>
docling_granite	0.6737	<u>0.6309</u>	0.6329	0.9001	0.1811	0.5296
mineru_pipeline	0.8735	0.9558	0.4771	0.9784	<b>0.6534</b>	0.6528
mineru_vlm	<b>0.9119</b>	0.8822	0.5702	0.9796	0.2508	0.5941
llamaparse	0.7711	0.6370	<u>0.0000</u>	<u>0.6831</u>	<u>0.0000</u>	<u>0.4110</u>
document_ai	<u>0.5823</u>	<b>0.9789</b>	<u>0.0000</u>	<b>0.9911</b>	<u>0.0000</u>	0.5043
gemini	0.8242	0.7619	0.1271	0.8725	0.6530	0.6153

(a) F1@50

	text	title	list	table	figure	all
unstructured_io	0.7583	0.6029	0.0682	0.8707	0.4987	0.5095
docling	<b>0.8206</b>	0.6311	<b>0.7406</b>	0.9143	0.4952	<b>0.6650</b>
docling_granite	0.6241	0.4302	0.5694	0.8497	0.1608	0.4382
mineru_pipeline	0.8097	0.6170	0.4331	0.9407	0.5436	0.5342
mineru_vlm	0.8032	0.4461	0.4906	0.9409	0.2034	0.4338
llamaparse	0.7240	<u>0.3057</u>	<u>0.0000</u>	<u>0.6594</u>	<u>0.0000</u>	<u>0.3073</u>
document_ai	<u>0.5304</u>	<b>0.6593</b>	<u>0.0000</u>	<b>0.9669</b>	<u>0.0000</u>	0.3928
gemini	0.7347	0.5002	0.0760	0.7636	<b>0.5822</b>	0.4890

(b) F1@50:95

Figure 4.1.: F1 scores of the evaluated DP implementations on the PubLayNet dataset. (a) contains the F1@50 scores, (b) contains the F1@50:95 scores. Scores are reported per element type. The column all reports the respective Macro F1 score for each implementation. Highest values are bolded and smallest values are underlined. Higher values are preferred.

---

Method	Parsing		Transformation	
	mean	std	mean	std
docling	2.1146	1.3147	0.0017	0.0024
docling_granite	16.7379	<u>63.5661</u>	0.0017	0.0034
document_ai	<b>1.6445</b>	<b>0.4928</b>	<u>1.6354</u>	<u>0.2689</u>
gemini	12.5033	11.1618	0.0013	0.0025
llamaparse	37.0956	41.6822	0.7354	0.1782
mineru_pipeline	15.4215	20.2052	0.0020	0.0059
mineru_vlm	<u>42.0182</u>	35.4030	0.0017	0.0012
unstructured_io	4.9745	5.5285	<b>0.0008</b>	<b>0.0005</b>

---

Table 4.1.: Mean and standard deviation of the DP implementation’s parsing and transformation times per page. Times are reported in seconds. Fastest times are bolded and slowest times are underlined. Lower values are preferred.

## 5. Discussion



## 6. Conclusion

- How to improve tables
- How to include figures
- Make bounding boxes more granular *to* token level instead of line level

## A. General Addenda

Listing A.1: Prompt to apply the Gemini 2.5 Flash model to DP tasks. Gemini models are trained to output coordinates from 0 to 1000, with the origin at the left-top corner of the image. In order to maximize the accuracy of detected bounding boxes, 'box\_2d', the key used in Google's official documentation, is used to denote the bounding box tuples in the output JSON.

```
<system_role>
You are an expert Document Layout Analysis AI. Your goal is to perfectly transcribe
and segment PDF documents into structured data.
</system_role>

<task_description>
Analyze the provided document image. Identify every layout element, its bounding box,
its category, and its textual content.
</task_description>

<categories>
Classify each element into exactly one of these categories:
section_header, text, formula, list_item, ref_item, table, image, caption,
page_header, page_footer, watermark

Rules for Categorization:
- Use "section_header" for titles and headings. Infer hierarchy based on content and
font size/boldness.
- Use "image" for charts, diagrams, or photos.
- Use "unknown" if the element is ambiguous.
</categories>

<bounding_boxes>
1. Format: [y0, x0, y1, x1] (Top-Left to Bottom-Right). You MUST provide the
coordinates in this exact order.
2. Success conditions:
- The bounding box MUST enclose the entire layout element while minimizing
unnecessary white space.
- If a character belongs to the content ALL of its pixels MUST BE CONTAINED inside
the bounding box.
3. Page Index: The current page is "page_number": {*}.
</bounding_boxes>

<extraction_rules>
```

- 
- **Text Fidelity:** Extract text EXACTLY as it appears. Do NOT fix spelling or grammar. You MAY use any formatting that is available for a standard Markdown document.
  - **Character Escaping:** You MUST escape any special characters that can break the final JSON output. Also you must escape any quotation marks.
  - **Reading Order:** Sort elements by natural human reading order.
  - **Special Formatting:**
    - image: Content must be an empty string "".
    - formula: Content must be LaTeX.
    - table: Content must be a Markdown table representation. TABLE CONTENT MUST NOT BREAK THE JSON FORMAT!
    - list\_item, ref\_item: Content MUST be a valid Markdown list. You MUST replace alternative bullet point symbols with "-". Ordered lists must start with their numbering followed by ".".
    - section\_header: You MUST NOT use Markdown header formatting. You MUST add a "heading\_level" field (int). Infer the level by checking the content for any numbering and analyzing the font size and styling of the header.

</extraction\_rules>

<output\_schema>

Do not return any additional text with the result.

Return a SINGLE JSON object with this exact structure:

```
{
  "layout_elements": [
    {
      "category": "string_(from_list)",
      "heading_level": integer (include only for headers),
      "content": "string",
      "bbox": {
        "page_number": integer,
        "box_2d": bounding_box (list[integer]) (SINGLE bounding box)
      }
    }
  ]
}
```

YOU MUST ENSURE THAT YOUR OUTPUT IS A VALID JSON OBJECT!

</output\_schema>

Classification	Description
ROOT	The top-level node containing the entire document structure
<b>TEXTS</b>	
TITLE	The specific main title of the document
PARAGRAPH	Standard body text content
SECTION_HEADER	Section headings or subheaders within the text body
FOOTNOTE	Explanatory notes usually placed at the bottom of a page/text
<b>LISTS</b>	
LIST	A container node for a list of items
LIST_ITEM	An individual item within a list
REFERENCE_LIST	A container node for a list of reference items
REFERENCE_ITEM	An individual item within a reference list
<b>FIGURES AND TABLES</b>	
CAPTION	Descriptive text immediately accompanying a table or figure
FIGURE	Graphical elements, diagrams, or pictures
TABLE	A container node for tabular data
DOC_INDEX	A tabular node containing the TOC
TABLE_ROW	A horizontal row within a table
TABLE_CELL	An individual cell containing data within a table row
<b>MISCELLANEOUS</b>	
PAGE_FOOTER	Repeating page footer (page numbers, copyright, etc.)
KEY_VALUE	A specific key-value pair
PAGE_HEADER	Repeating header found at the top of pages (e.g., journal name)
KEY_VALUE_AREA	A distinct region grouped by key-value pairs (e.g., article info)
FORM_AREA	A region indicating form content (e.g., text-fields)
FORMULA	A mathematical formula
WATERMARK	A watermark from the publishing organization
<b>FALLBACK</b>	
UNKNOWN	Parser cannot determine the element type
MISSING	Parser returns a classification for which no mapping exists

Table A.1.: Complete list of classifications permitted to be returned by a DP implementation. Each implementation provides a mapping from their native output classifications to the standard set defined here. Some ParsingResultTypes may only be returned from a subset of these implementations.

# List of Figures

- 2.1. Naive RAG . . . . . 8
- 3.1. Document Segmentation Pipeline . . . . . 10
- 3.2. ParsingBoundingBox . . . . . 11
- 3.3. ParsingResult . . . . . 12
- 3.4. ChunkingResult . . . . . 13
- 3.5. RichToken . . . . . 16
- 4.1. F1 scores on the PubLayNet dataset . . . . . 21

# List of Tables

3.1. Parsing module abstract functions . . . . .	13
3.2. Distribution of ground truth annotations across the different element types contained in the publaynet-mini subset of the PubLayNet dataset. . . . .	18
4.1. Parsing and transformation times per page . . . . .	22
A.1. ParsingResultType . . . . .	27

# Acronyms

**AI** artificial intelligence. 1, 7, 13, 14

**AP** average precision. 16

**CNN** convolutional neural networks. 7

**DLA** document layout analysis. 6, 7, 12, 13, 16

**DP** document parsing. 1–3, 6, 7, 9–14, 16–19, 23

**DTBB** detected bounding box. 16, 17

**FN** false negative. 16

**FP** false positive. 16

**GTBB** ground truth bounding box. 16, 17

**IoU** intersection over union. 5, 16, 17

**LLM** large language model. 1, 2, 4, 6

**NCCN** National Comprehensive Cancer Network. 1

**NLP** natural language processing. 1, 4, 5

**OCR** optical character recognition. 7, 13

**PDF** portable document format. 1, 6, 9, 12, 13

**RAG** retrieval-augmented generation. 1, 2, 4, 6, 9, 12–14

**TP** true positive. 16, 17

**VLM** vision-language model. 7, 13

# Bibliography

- [1] E. Steinberg, S. Greenfield, D. M. Wolman, M. Mancher, and R. Graham. *Clinical practice guidelines we can trust*. national academies press, 2011.
- [2] National Comprehensive Cancer Network. *About Clinical Practice Guidelines*. NCCN. URL: <https://www.nccn.org/guidelines/guidelines-process/about-nccn-clinical-practice-guidelines> (visited on 01/29/2026).
- [3] B. H. Kann, S. B. Johnson, H. J. Aerts, R. H. Mak, and P. L. Nguyen. “Changes in length and complexity of clinical practice guidelines in oncology, 1996-2019”. In: *JAMA Network Open* 3.3 (2020), e200841–e200841.
- [4] Chair of Software Engineering for Business Information Systems, TUM School of Computation, Information and Technology, Technical University of Munich. *AI-Based Knowledge Assistant for Cancer Care (Aidvice)*. 2025. URL: <https://www.cs.cit.tum.de/sebis/research/natural-language-processing/ai-based-knowledge-assistant-for-cancer-care-aidvice/> (visited on 01/28/2026).
- [5] 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*. 2021. arXiv: 2005.11401 [cs.CL]. URL: <https://arxiv.org/abs/2005.11401>.
- [6] Y. Gao, Y. Xiong, X. Gao, K. Jia, J. Pan, Y. Bi, Y. Dai, J. Sun, M. Wang, and H. Wang. *Retrieval-Augmented Generation for Large Language Models: A Survey*. 2024. arXiv: 2312.10997 [cs.CL]. URL: <https://arxiv.org/abs/2312.10997>.
- [7] S. Wang, F. Zhao, D. Bu, and et al. “LINS: A general medical Q&A framework for enhancing the quality and credibility of LLM-generated responses”. In: *Nature Communications* 16.1 (Oct. 2025), p. 9076. DOI: 10.1038/s41467-025-64142-2. URL: <https://doi.org/10.1038/s41467-025-64142-2>.
- [8] J. Hladěna, K. Šteflovíč, P. Čech, K. Štekerová, and A. Žváčková. “The Effect of Chunk Size on the RAG Performance”. In: *Software Engineering: Emerging Trends and Practices in System Development*. Ed. by R. Silhavy and P. Silhavy. Cham: Springer Nature Switzerland, 2025, pp. 317–326. ISBN: 978-3-032-00712-4.
- [9] T. Chen, H. Wang, S. Chen, W. Yu, K. Ma, X. Zhao, H. Zhang, and D. Yu. *Dense X Retrieval: What Retrieval Granularity Should We Use?* 2024. arXiv: 2312.06648 [cs.CL]. URL: <https://arxiv.org/abs/2312.06648>.
- [10] L. Müller, J. Holstein, S. Bause, G. Satzger, and N. Kühl. *Data Quality Challenges in Retrieval-Augmented Generation*. 2025. arXiv: 2510.00552 [cs.AI]. URL: <https://arxiv.org/abs/2510.00552>.



- [11] X. Ma, S. Zhuang, B. Koopman, G. Zuccon, W. Chen, and J. Lin. *VISA: Retrieval Augmented Generation with Visual Source Attribution*. 2024. arXiv: 2412.14457 [cs.IR]. URL: <https://arxiv.org/abs/2412.14457>.
- [12] Q. Zhang, B. Wang, V. S.-J. Huang, J. Zhang, Z. Wang, H. Liang, C. He, and W. Zhang. *Document Parsing Unveiled: Techniques, Challenges, and Prospects for Structured Information Extraction*. 2025. arXiv: 2410.21169 [cs.MM]. URL: <https://arxiv.org/abs/2410.21169>.
- [13] R. Qu, R. Tu, and F. Bao. “Is semantic chunking worth the computational cost?” In: *Findings of the Association for Computational Linguistics: NAACL 2025*. 2025, pp. 2155–2177.
- [14] H. Chase. *LangChain*. <https://github.com/langchain-ai/langchain>. Oct. 2022. URL: <https://github.com/langchain-ai/langchain>.
- [15] J. Liu. *LlamaIndex*. Nov. 2022. DOI: 10.5281/zenodo.1234. URL: [https://github.com/jerryjliu/llama\\_index](https://github.com/jerryjliu/llama_index).
- [16] N. Livathinos, C. Auer, M. Lysak, A. Nassar, M. Dolfi, P. Vagenas, C. B. Ramis, M. Omenetti, K. Dinkla, Y. Kim, S. Gupta, R. T. de Lima, V. Weber, L. Morin, I. Meijer, V. Kuropiatnyk, and P. W. J. Staar. *Docling: An Efficient Open-Source Toolkit for AI-driven Document Conversion*. 2025. arXiv: 2501.17887 [cs.CL]. URL: <https://arxiv.org/abs/2501.17887>.
- [17] LangChain Inc. *LangChain Docs: Splitting recursively*. URL: [https://docs.langchain.com/oss/python/integrations/splitters/recursive\\_text\\_splitter](https://docs.langchain.com/oss/python/integrations/splitters/recursive_text_splitter) (visited on 01/27/2026).
- [18] A. Ghosh, A. Acharya, S. Saha, V. Jain, and A. Chadha. *Exploring the Frontier of Vision-Language Models: A Survey of Current Methodologies and Future Directions*. 2025. arXiv: 2404.07214 [cs.CV]. URL: <https://arxiv.org/abs/2404.07214>.
- [19] S. Bai, K. Chen, X. Liu, J. Wang, W. Ge, S. Song, K. Dang, P. Wang, S. Wang, J. Tang, et al. “Qwen2. 5-vl technical report”. In: *arXiv preprint arXiv:2502.13923* (2025).
- [20] L. Ouyang, Y. Qu, H. Zhou, J. Zhu, R. Zhang, Q. Lin, B. Wang, Z. Zhao, M. Jiang, X. Zhao, J. Shi, F. Wu, P. Chu, M. Liu, Z. Li, C. Xu, B. Zhang, B. Shi, Z. Tu, and C. He. *OmniDocBench: Benchmarking Diverse PDF Document Parsing with Comprehensive Annotations*. 2024. arXiv: 2412.07626 [cs.CV]. URL: <https://arxiv.org/abs/2412.07626>.
- [21] B. Wang, C. Xu, X. Zhao, L. Ouyang, F. Wu, Z. Zhao, R. Xu, K. Liu, Y. Qu, F. Shang, B. Zhang, L. Wei, Z. Sui, W. Li, B. Shi, Y. Qiao, D. Lin, and C. He. *MinerU: An Open-Source Solution for Precise Document Content Extraction*. 2024. arXiv: 2409.18839 [cs.CV]. URL: <https://arxiv.org/abs/2409.18839>.
- [22] D. S. Team. *Docling Technical Report*. Tech. rep. Version 1.0.0. Aug. 2024. DOI: 10.48550/arXiv.2408.09869. eprint: 2408.09869. URL: <https://arxiv.org/abs/2408.09869>.

- [23] H. Xing, F. Gao, Q. Zheng, Z. Zhu, Z. Shao, and M. Yan. “Intelligent Document Parsing: Towards End-to-end Document Parsing via Decoupled Content Parsing and Layout Grounding”. In: *Findings of the Association for Computational Linguistics: EMNLP 2025*. Ed. by C. Christodoulopoulos, T. Chakraborty, C. Rose, and V. Peng. Suzhou, China: Association for Computational Linguistics, Nov. 2025, pp. 19987–19998. ISBN: 979-8-89176-335-7. DOI: 10.18653/v1/2025.findings-emnlp.1088. URL: <https://aclanthology.org/2025.findings-emnlp.1088/>.
- [24] J. Niu, Z. Liu, Z. Gu, B. Wang, L. Ouyang, Z. Zhao, T. Chu, T. He, F. Wu, Q. Zhang, Z. Jin, G. Liang, R. Zhang, W. Zhang, Y. Qu, Z. Ren, Y. Sun, Y. Zheng, D. Ma, Z. Tang, B. Niu, Z. Miao, H. Dong, S. Qian, J. Zhang, J. Chen, F. Wang, X. Zhao, L. Wei, W. Li, S. Wang, R. Xu, Y. Cao, L. Chen, Q. Wu, H. Gu, L. Lu, K. Wang, D. Lin, G. Shen, X. Zhou, L. Zhang, Y. Zang, X. Dong, J. Wang, B. Zhang, L. Bai, P. Chu, W. Li, J. Wu, L. Wu, Z. Li, G. Wang, Z. Tu, C. Xu, K. Chen, Y. Qiao, B. Zhou, D. Lin, W. Zhang, and C. He. *MinerU2.5: A Decoupled Vision-Language Model for Efficient High-Resolution Document Parsing*. 2025. arXiv: 2509.22186 [cs.CV]. URL: <https://arxiv.org/abs/2509.22186>.
- [25] Unstructured.io Team. *Unstructured.io: Open-Source Pre-Processing Tools for Unstructured Data*. URL: <https://unstructured.io> (visited on 01/20/2026).
- [26] Z. Li, Y. Liu, Q. Liu, Z. Ma, Z. Zhang, S. Zhang, Z. Guo, J. Zhang, X. Wang, and X. Bai. *MonkeyOCR: Document Parsing with a Structure-Recognition-Relation Triplet Paradigm*. 2025. arXiv: 2506.05218 [cs.CV]. URL: <https://arxiv.org/abs/2506.05218>.
- [27] N. Livathinos, C. Auer, A. Nassar, R. T. de Lima, M. Lysak, B. Ebouky, C. Berrospi, M. Dolfi, P. Vagenas, M. Omenetti, K. Dinkla, Y. Kim, V. Weber, L. Morin, I. Meijer, V. Kuropiatnyk, T. Strohmeyer, A. S. Gurbuz, and P. W. J. Staar. *Advanced Layout Analysis Models for Docling*. 2025. arXiv: 2509.11720 [cs.CV]. URL: <https://arxiv.org/abs/2509.11720>.
- [28] T. Sun, C. Cui, Y. Du, and Y. Liu. *PP-DocLayout: A Unified Document Layout Detection Model to Accelerate Large-Scale Data Construction*. 2025. arXiv: 2503.17213 [cs.CV]. URL: <https://arxiv.org/abs/2503.17213>.
- [29] J. Redmon, S. K. Divvala, R. B. Girshick, and A. Farhadi. “You Only Look Once: Unified, Real-Time Object Detection”. In: *CoRR abs/1506.02640* (2015). arXiv: 1506.02640. URL: <http://arxiv.org/abs/1506.02640>.
- [30] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko. “End-to-End Object Detection with Transformers”. In: *CoRR abs/2005.12872* (2020). arXiv: 2005.12872. URL: <https://arxiv.org/abs/2005.12872>.
- [31] Y. Huang, T. Lv, L. Cui, Y. Lu, and F. Wei. *LayoutLMv3: Pre-training for Document AI with Unified Text and Image Masking*. 2022. arXiv: 2204.08387 [cs.CL]. URL: <https://arxiv.org/abs/2204.08387>.
- [32] N. Islam, Z. Islam, and N. Noor. “A survey on optical character recognition system”. In: *arXiv preprint arXiv:1710.05703* (2017).

- [33] R. Kittinaradorn and JaidevAI. *EasyOCR*. <https://github.com/JaidevAI/EasyOCR>. 2020.
- [34] R. Smith. “An overview of the Tesseract OCR engine”. In: *Ninth international conference on document analysis and recognition (ICDAR 2007)*. Vol. 2. IEEE. 2007, pp. 629–633.
- [35] Unstructured.io Team. *Unstructured.io: Documentation for the open-source library*. URL: <https://docs.unstructured.io/open-source/introduction/overview> (visited on 01/28/2026).
- [36] H. Wei, C. Liu, J. Chen, J. Wang, L. Kong, Y. Xu, Z. Ge, L. Zhao, J. Sun, Y. Peng, C. Han, and X. Zhang. *General OCR Theory: Towards OCR-2.0 via a Unified End-to-end Model*. 2024. arXiv: 2409.01704 [cs.CV]. URL: <https://arxiv.org/abs/2409.01704>.
- [37] Y. Li, G. Yang, H. Liu, B. Wang, and C. Zhang. *dots.ocr: Multilingual Document Layout Parsing in a Single Vision-Language Model*. 2025. arXiv: 2512.02498 [cs.CV]. URL: <https://arxiv.org/abs/2512.02498>.
- [38] N. Aksoy, Z. A. Güven, and M. O. Ünalır. “Understanding the Impact of Dataset Characteristics on RAG based Multi-hop QA Performance”. In: (2025).
- [39] X. Song, A. Salcianu, Y. Song, D. Dopson, and D. Zhou. “Fast wordpiece tokenization”. In: *Proceedings of the 2021 conference on empirical methods in natural language processing*. 2021, pp. 2089–2103.
- [40] S. J. Mielke, Z. Alyafeai, E. Salesky, C. Raffel, M. Dey, M. Gallé, A. Raja, C. Si, W. Y. Lee, B. Sagot, and S. Tan. *Between words and characters: A Brief History of Open-Vocabulary Modeling and Tokenization in NLP*. 2021. arXiv: 2112.10508 [cs.CL]. URL: <https://arxiv.org/abs/2112.10508>.
- [41] M. Schuster and K. Nakajima. “Japanese and Korean voice search”. In: *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 2012, pp. 5149–5152. doi: 10.1109/ICASSP.2012.6289079.
- [42] T. Kudo and J. Richardson. “SentencePiece: A simple and language independent subword tokenizer and detokenizer for Neural Text Processing”. In: *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Ed. by E. Blanco and W. Lu. Brussels, Belgium: Association for Computational Linguistics, Nov. 2018, pp. 66–71. doi: 10.18653/v1/D18-2012. URL: <https://aclanthology.org/D18-2012/>.
- [43] T. Kudo. “Subword Regularization: Improving Neural Network Translation Models with Multiple Subword Candidates”. In: *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Ed. by I. Gurevych and Y. Miyao. Melbourne, Australia: Association for Computational Linguistics, July 2018, pp. 66–75. doi: 10.18653/v1/P18-1007. URL: <https://aclanthology.org/P18-1007/>.
- [44] J. Ni, C. Qu, J. Lu, Z. Dai, G. H. Abrego, J. Ma, V. Zhao, Y. Luan, K. Hall, M.-W. Chang, et al. “Large dual encoders are generalizable retrievers”. In: *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. 2022, pp. 9844–9855.

- [45] R. Kaur and S. Singh. “A comprehensive review of object detection with deep learning”. In: *Digital Signal Processing* 132 (2023), p. 103812. issn: 1051-2004. doi: <https://doi.org/10.1016/j.dsp.2022.103812>. URL: <https://www.sciencedirect.com/science/article/pii/S1051200422004298>.
- [46] C. He, W. Li, Z. Jin, C. Xu, B. Wang, and D. Lin. “Opendatalab: Empowering general artificial intelligence with open datasets”. In: *arXiv preprint arXiv:2407.13773* (2024).
- [47] Z. Ge, S. Liu, F. Wang, Z. Li, and J. Sun. *YOLOX: Exceeding YOLO Series in 2021*. 2021. arXiv: 2107.08430 [cs.CV]. URL: <https://arxiv.org/abs/2107.08430>.
- [48] A. Nassar, N. Livathinos, M. Lysak, and P. Staar. “TableFormer: Table Structure Understanding With Transformers”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. June 2022, pp. 4614–4623. doi: <https://doi.org/10.1109/CVPR52688.2022.00457>.
- [49] Y. Zhao, W. Lv, S. Xu, J. Wei, G. Wang, Q. Dang, Y. Liu, and J. Chen. *DETRs Beat YOLOs on Real-time Object Detection*. 2024. arXiv: 2304.08069 [cs.CV]. URL: <https://arxiv.org/abs/2304.08069>.
- [50] B. Pfitzmann, C. Auer, M. Dolfi, A. S. Nassar, and P. W. J. Staar. “DocLayNet: A Large Human-Annotated Dataset for Document-Layout Analysis”. In: (2022). doi: 10.1145/3534678.353904. URL: <https://arxiv.org/abs/2206.01062>.
- [51] IBM Research. *Granite Docling Documentation*. IBM. URL: <https://www.ibm.com/granite/docs/models/docling> (visited on 02/03/2026).
- [52] A. Nassar, A. Marafioti, M. Omenetti, M. Lysak, N. Livathinos, C. Auer, L. Morin, R. T. de Lima, Y. Kim, A. S. Gurbuz, M. Dolfi, M. Farré, and P. W. J. Staar. *SmolDocling: An ultra-compact vision-language model for end-to-end multi-modal document conversion*. 2025. arXiv: 2503.11576 [cs.CV]. URL: <https://arxiv.org/abs/2503.11576>.
- [53] J. X. M. Artifex Software Inc. *PyMuPDF*. Version 1.26.7. 2025. URL: <https://github.com/pymupdf/PyMuPDF>.
- [54] Z. Zhao, H. Kang, B. Wang, and C. He. *DocLayout-YOLO: Enhancing Document Layout Analysis through Diverse Synthetic Data and Global-to-Local Adaptive Perception*. 2024. arXiv: 2410.12628 [cs.CV]. URL: <https://arxiv.org/abs/2410.12628>.
- [55] B. Wang, Z. Gu, C. Xu, B. Zhang, B. Shi, and C. He. *UniMERNet: A Universal Network for Real-World Mathematical Expression Recognition*. 2024. arXiv: 2404.15254 [cs.CV].
- [56] G. Comanici, E. Bieber, M. Schaekermann, et al. *Gemini 2.5: Pushing the Frontier with Advanced Reasoning, Multimodality, Long Context, and Next Generation Agentic Capabilities*. 2025. arXiv: 2507.06261 [cs.CL]. URL: <https://arxiv.org/abs/2507.06261>.
- [57] Google. *Image Understanding*. Google AI for Developers. URL: <https://ai.google.dev/gemini-api/docs/image-understanding> (visited on 02/03/2026).

- [58] S. Jaiswal, P. Bisht, K. Kansara, and M. S. Datta. "Comparison of Chunking Techniques Across Diverse Document Types in NLP Retrieval Tasks". In: *2025 International Conference on Responsible, Generative and Explainable AI (ResGenXAI)*. 2025, pp. 1–6. doi: 10.1109/ResgenXAI64788.2025.11344045.
- [59] B. Smith and A. Troynikov. *Evaluating Chunking Strategies for Retrieval*. Tech. rep. Chroma, June 2024. URL: <https://research.trychroma.com/evaluating-chunking>.
- [60] A. Antonacopoulos, S. Pletschacher, D. Bridson, and C. Papadopoulos. "ICDAR 2009 Page Segmentation Competition". In: *2009 10th International Conference on Document Analysis and Recognition*. 2009, pp. 1370–1374. doi: 10.1109/ICDAR.2009.275.
- [61] M. Li, Y. Xu, L. Cui, S. Huang, F. Wei, Z. Li, and M. Zhou. *DocBank: A Benchmark Dataset for Document Layout Analysis*. 2020. arXiv: 2006.01038 [cs.CL]. URL: <https://arxiv.org/abs/2006.01038>.
- [62] X. Zhong, J. Tang, and A. J. Yepes. "PubLayNet: largest dataset ever for document layout analysis". In: *2019 International Conference on Document Analysis and Recognition (ICDAR)*. IEEE. Sept. 2019, pp. 1015–1022. doi: 10.1109/ICDAR.2019.00166.
- [63] K. Benkirane. *publaynet-mini*. <https://huggingface.co/datasets/kenza-ily/publaynet-mini>. 2029.
- [64] Z. Zou, K. Chen, Z. Shi, Y. Guo, and J. Ye. *Object Detection in 20 Years: A Survey*. 2023. arXiv: 1905.05055 [cs.CV]. URL: <https://arxiv.org/abs/1905.05055>.
- [65] A. J. Yepes, X. Zhong, and D. Burdick. *ICDAR 2021 Competition on Scientific Literature Parsing*. 2021. arXiv: 2106.14616 [cs.IR]. URL: <https://arxiv.org/abs/2106.14616>.
- [66] K. Oksuz, B. C. Cam, S. Kalkan, and E. Akbas. *One Metric to Measure them All: Localisation Recall Precision (LRP) for Evaluating Visual Detection Tasks*. 2021. arXiv: 2011.10772 [cs.CV]. URL: <https://arxiv.org/abs/2011.10772>.
- [67] I. Karmanov, A. S. Deshmukh, L. Voegtler, P. Fischer, K. Chumachenko, T. Roman, J. Seppänen, J. Parmar, J. Jennings, A. Tao, et al. "Eclair-Extracting Content and Layout with Integrated Reading Order for Documents". In: *arXiv preprint arXiv:2502.04223* (2025).
- [68] R. Padilla, W. L. Passos, T. L. B. Dias, S. L. Netto, and E. A. B. da Silva. "A Comparative Analysis of Object Detection Metrics with a Companion Open-Source Toolkit". In: *Electronics* 10.3 (2021). issn: 2079-9292. doi: 10.3390/electronics10030279. URL: <https://www.mdpi.com/2079-9292/10/3/279>.
- [69] A. Avetisyan, C. Xie, H. Howard-Jenkins, T.-Y. Yang, S. Aroudj, S. Patra, F. Zhang, D. Frost, L. Holland, C. Orme, et al. "Scenescript: Reconstructing scenes with an autoregressive structured language model". In: *European Conference on Computer Vision*. Springer. 2024, pp. 247–263.
- [70] Z. C. Lipton, C. Elkan, and B. Narayanaswamy. *Thresholding Classifiers to Maximize F1 Score*. 2014. arXiv: 1402.1892 [stat.ML]. URL: <https://arxiv.org/abs/1402.1892>.

- [71] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. "Microsoft coco: Common objects in context". In: *European conference on computer vision*. Springer. 2014, pp. 740–755.
- [72] MiXaiLL76. "Faster-COCO-Eval: Faster and Enhanced COCO Evaluation Library". In: (2024).