# Revolutionizing Retrieval-Augmented Generation with Enhanced PDF Structure Recognition

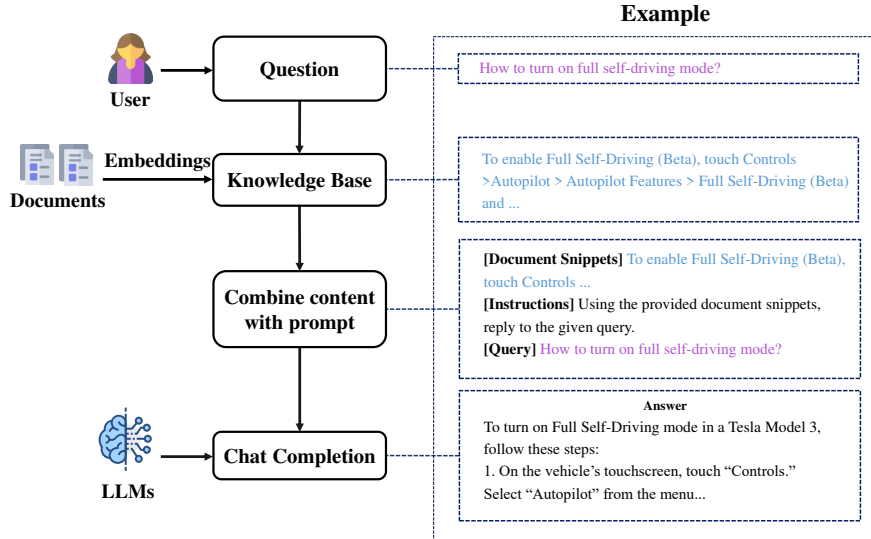**Demiao LIN**
chatdoc.com

## Abstract

With the rapid development of Large Language Models (LLMs), Retrieval-Augmented Generation (RAG) has become a predominant method in the field of professional knowledge-based question answering. Presently, major foundation model companies have opened up Embedding and Chat API interfaces, and frameworks like LangChain have already integrated the RAG process. It appears that the key models and steps in RAG have been resolved, leading to the question: are professional knowledge QA systems now approaching perfection? This article discovers that current primary methods depend on the premise of accessing high-quality text corpora. However, since professional documents are mainly stored in PDFs, the low accuracy of PDF parsing significantly impacts the effectiveness of professional knowledge-based QA. We conducted an empirical RAG experiment across hundreds of questions from the corresponding real-world professional documents. The results show that, ChatDOC (chatdoc.com), a RAG system equipped with a panoptic and pinpoint PDF parser, retrieves more accurate and complete segments, and thus better answers. Empirical experiments show that ChatDOC is superior to baseline on nearly 47% of questions, ties for 38% of cases, and falls short on only 15% of cases. It shows that we may revolutionize RAG with enhanced PDF structure recognition.
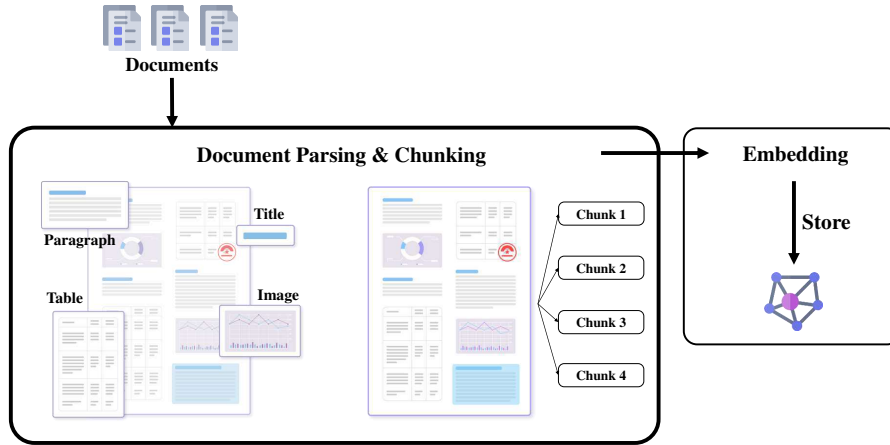
## 1 Introduction

Large language models (LLM) are trained on data that predominantly come from publicly available internet sources, including web pages, books, news, and dialogue texts. It means that LLMs primarily rely on internet sources as their training data, which are vast, diverse, and easily accessible, supporting them to scale up their capabilities. However, in vertical applications, professional tasks require LLMs to utilize domain knowledge, which unfortunately is private, and not part of their pre-training data.

A popular approach to equip LLM with domain knowledge is Retrieval-Augmented Generation (RAG). RAG framework answers a question in four steps: the user proposes a query, the system retrieves the relevant content from private knowledge bases, combines it with the user query as context, and finally asks the LLM to generate an answer. This is illustrated in Figure 1 with a simple example. This process mirrors the typical cognitive process of encountering a problem, including consulting relevant references and subsequently deriving an answer. In this framework, the pivotal component is the accurate retrieval of pertinent information, which is critical for the efficacy of the RAG model.

However, the process of retrieval from PDF files is fraught with challenges. Common issues include inaccuracies in text extraction and disarray in the row-column relationships of tables inside PDF files. Thus, before RAG, we need to convert large documents into retrievable content. The conversion involves several steps, as shown in Figure 2:

**Figure 1.** The workflow of Retrieval-Augmented Generation (RAG).



**Figure 2.** The process of converting PDFs into retrievable contents.

- Document Parsing & Chunking. It involves extracting paragraphs, tables, and other content blocks, then dividing the extracted content into chunks for subsequent retrieval.
- Embedding. It transforms text chunks into real-valued vectors and stores them in a database.
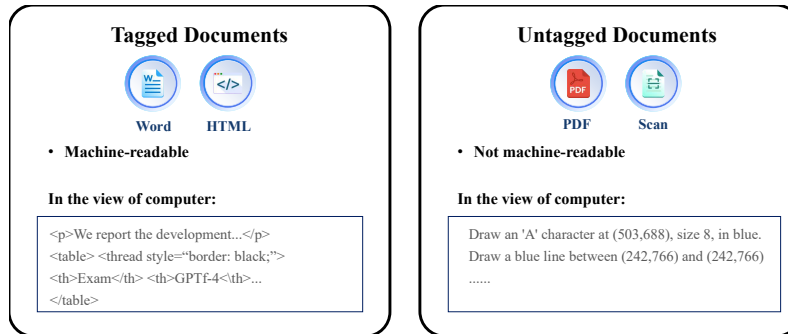
Since each of these steps can lead to information loss, the compounded losses can significantly impact the effectiveness of RAG's responses.

This paper primarily addresses the question of whether the quality of PDF parsing and chunking affects the outcomes of RAG. We will explore the challenges, methodologies, and real-world case studies pertaining to this issue. It will include an examination of two types of methods in this field, namely rule-based and deep learning-based methods, followed by empirical evaluations of their efficacy through practical examples.

## 2  PDF Parsing & Chunking

### 2.1  Challenges and Methods Overview

To humans, the cognitive process of perusing any document page is similar. When we read a page, characters are captured by our retinas. Then, in our brains, these characters are organized into

**Figure 3.** Two types of documents in the view of computers.

paragraphs, tables, and charts, and then understood or memorized. However, computers perceive information as binary codes. From their perspective, as illustrated in Figure 3, documents can be categorized into two distinct types:

- Tagged Documents: Examples include Microsoft Word and HTML documents, which contain special tags like <p>and <table> to organize the text into paragraphs, cells, and tables.
- Untagged Documents: Examples include PDFs, which store instructions on the placement of characters, lines, and other content elements on each document page. They focus on 'drawing' these basic content elements in a way that makes the document legible to human readers. They do not store any structural information of the document, like tables or paragraphs. Thus, untagged documents are only for human e-reading, but are unreadable by machines. This becomes evident when attempting to copy a table from a PDF into MS Word, where the original structure of the table is often completely lost.
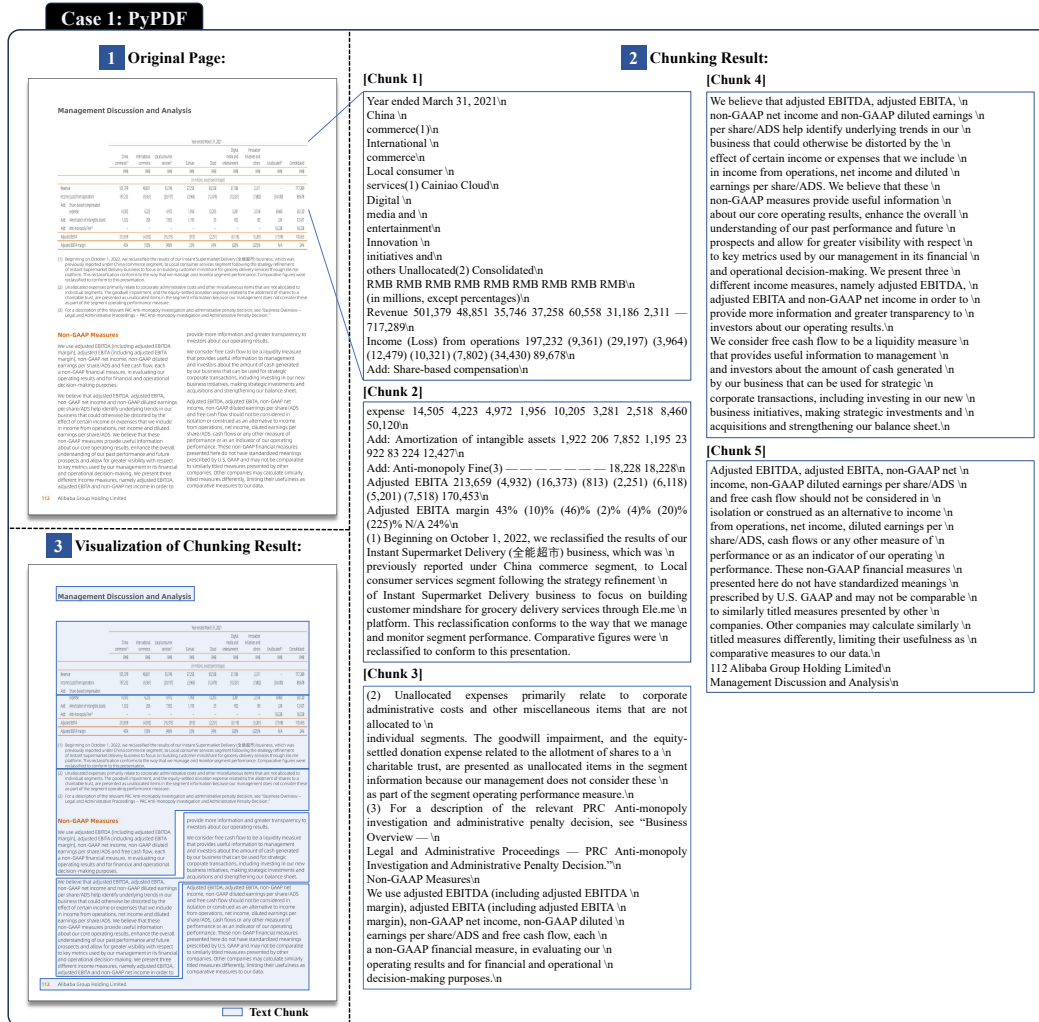
However, Large Language Models (LLMs) exhibit proficiency in processing serialized text. Consequently, to enable LLMs to effectively manage untagged documents, a parser that organizes scattered characters into coherent texts with their structures is necessary. Ideally, a PDF Parser should exhibit the following key features:

- **Document Structure Recognition**: It should adeptly divide pages into different types of content blocks like paragraphs, tables, and charts. This ensures that the divided text blocks are complete and independent semantic units.
- **Robustness in Complex Document Layout**: It should work well even for document pages with complex layouts, such as multi-column pages, border-less tables, and even tables with merged cells.

Currently, there are two main types of methods of PDF Parsing: rule based approaches and deep learning-based approaches. Among them, PyPDF, a widely-used rule-based parser, is a standard method in LangChain for PDF parsing. Conversely, our approach, ChatDOC PDF Parser (`https://pdfparser.io/`), is grounded in the deep learning models. Next, we illustrate the difference between them by introducing the methods and delving into some real-world cases.

## 2.2 Rule-based Method: PyPDF

We first introduce the parsing & chunking workflow based on PyPDF. First, PyPDF serializes characters in a PDF into a long sequence without document structure information. Then, this sequence undergoes segmentation into discrete chunks, utilizing some segmentation rule, such as the "`RecursiveCharacterTextSplitter`" function in LangChain. Specifically, this function divides the document based on a predefined list of separators, such as the newline character "\n". After this initial segmentation, adjacent chunks are merged only if the length of the combined chunks is not bigger than a predetermined limit of $N$ characters. Hereafter, we use "PyPDF" to refer to the method of document parsing and chunking using PyPDF+`RecursiveCharacterTextSplitter`, provided there is no contextual ambiguity. The maximum length of a chunk is set to 300 tokens in the following. Next, we use a case to observe the inherent nature of PyPDF.

**Figure 4.** Parsing and chunking results of PyPDF on Case 1 (original document: [1]). Zoom in to see the details.

**Case 1** in Figure 4 is a page from a document that features a mix of a table and double-column text where their boundaries are difficult to distinguish. Rows in the middle of the table do not have horizontal lines, making it difficult to recognize the rows in the table. And paragraphs have both single-column layout (for notes below the table) and double-columns layout (for paragraphs in the lower part of the page).

The parsing and chunking result of PyPDF is shown in Figure 4. In the "3 Visualization" part, we can see that PyPDF correctly recognizes the one-column and two-column layout parts of the page. But there are three shortcomings of PyPDF:

1. It cannot recognize the boundary of paragraphs and tables. It wrongly splits the table into two parts and merges the second part with the subsequent paragraph as one chunk.

   PyPDF seems to be good at detecting the boundary of a paragraph, as it does not split one paragraph into multiple chunks. But it actually does not parse the boundary of a paragraph. In the "2 Chunking Result" part we can see that each visual text line in the page is parsed as a line ended with "\n" in the result, and there is no special format at the end of a paragraph. It chunks paragraphs correctly because we use a special separator ".\n" that regards a line ending with a period as likely to be the end of a paragraph. However, this heuristic may not hold in many cases.

2. It cannot recognize the structure within a table. In the "2 Chunking Result" part, in chunk1, the upper part of the table is represented as a sequence of short phrases, where a cell may be split

into multiple lines (e.g. the cell "China commerce(1)") and some adjacent cells may be arranged in one line (e.g. the third to the fifth cells in the second line, "services(1) Cainiao Cloud"). So, the structure of the table is completely destroyed. If this chunk is retrieved for RAG, LLM is unable to perceive any meaningful information from it. Similar situation for Chunk 2. Moreover, the headers of the table only exist in Chunk 1, so the lower part of the table in Chunk 2 becomes meaningless.

3. It cannot recognize the reading order of the content. The last line of Chunk 5, "Management Discussion and Analysis" is actually located at the top of the page, but is parsed as the last sentence in the result. This is because PyPDF parses the document by the storage order of the characters, instead of their reading order. This may cause chaotic results when faced with complex layouts.

The result on another case Case 2 features with a complex and cross-page table is shown in Figure 15 in the Appendix.

## 2.3    Deep Learning-based Method: ChatDOC PDF Parser

Next, we turn our attention to the method of deep learning-based parsing, exemplified by our Chat-DOC PDF Parser. The ChatDOC PDF Parser (https://pdfparser.io/) has been trained on a corpus of over ten million document pages. Following the method in [2], it incorporates a sequence of sophisticated steps, including:

1. OCR for text positioning and recognition;
2. Physical document object detection;
3. Cross-column and cross-page trimming;
4. Reading order determination;
5. Table structure recognition;
6. Document logical structure recognition.

Readers might refer to [2] for the details of these steps. After parsing, we use the paragraphs and tables as basic blocks, and merge adjacent blocks until reaching the token limit to form a chunk.

ChatDOC PDF Parser is designed to consistently deliver parsing results in JSON or HTML formats, even for challenging PDF files. It parses a document into content blocks where each block refers to a table, paragraph, chart, or other type. For tables, it outputs the text in each table cell and also tells which cells are merged into a new one. Moreover, for documents with hierarchical headings, it outputs the hierarchical structure of the document. In summary, the parsed result is like a well-organized Word file. Figure 5 shows a scan-copy page and its parsing result. The left side displays the document and the recognized content blocks (with different colored rectangles). The right side shows the parsing result in JSON or HTML format. Readers might refer to [3] for the live demo of this parsing result.

Then, we check the result of ChatDOC PDF Parser on Case 1 in Figure 6. It successfully addresses the three shortcomings of PyPDF.

1. As shown in the "3 Visualization" part, it recognizes the mixed layout and correctly sets the whole table as a separate chunk. For paragraphs, as shown in chunk 2 in the "2 Chunking Result" part, text lines in the same paragraphs are merged together, making it easier to understand.

2. In the "2 Chunking Result" part, in Chunk 1, we can see the table is represented using the `markdown` format, which preserves the table's internal structure. Additionally, ChatDOC PDF Parser can recognize the merged cells inside a table. Since the `markdown` format cannot represent the merged cells, we put the whole text in the merged cell into each original cell in the `markdown` format. As you can see, in Chunk 1 the text "Year ended March 31, 2021" repeats 9 times, which stands for a merged cell with the original 9 ones.

3. Moreover, "Management Discussion and Analysis" and "112 Alibaba Group Holding Limited" is recognized as the page header and footer, and they are placed at the top and bottom of the parsing result which is consistent with reading order.

The result on another case of Case 2 featured with complex and cross-page table is shown in Figure 16 in the Appendix.

**Figure 5.** An example illustrating the results of the ChatDOC PDF Parser. Zoom in to see the details.

# 3  Experiments on the Impact of PDF Recognition on RAG

Back to the main topic of this paper, does the way a document is parsed and chunked affect the quality of answers provided by an RAG system? To answer this, we have carried out a systematic experiment to assess the impacts.

## 3.1  Quantitative Evaluation of RAG Answer Quality

### 3.1.1  Settings

We compared two RAG systems as listed in Table 1:

- ChatDOC: uses ChatDOC PDF Parser to parse the document and leverage the structure information for chunking.
- Baseline: uses PyPDF to parse the document and use `RecursiveCharacterTextSplitter` function for chunking.

Other components, like embedding, retrieval, and QA, are the same for both systems.

### 3.1.2  Data Preparation

For our experiment, we assembled a dataset that closely mirrors real-world conditions, comprising 188 documents from various fields. Specifically, This collection includes 100 academic papers, 28 financial reports, and 60 documents from other categories such as textbooks, courseware, and legislative materials.

We then gathered 800 manually generated questions via crowd-sourcing. After careful screening, we removed low-quality questions and got 302 questions for evaluation. These questions were divided into two categories (as shown in Table 2):

- **Extractive questions** are those that can be answered with direct excerpts from the documents. Usually, they require pinpoint answers because they seek specific information. We found when

**Case 1: ChatDOC PDF Parser**

**1  Original Page:**

Management Discussion and Analysis

**3  Visualization of Chunking Result:**

Management Discussion and Analysis

Paragraph
Table
Header & Footer

**2  Chunking Result:**

**[Chunk 1]**

\<Page Header\> Management Discussion and Analysis\n
| | Year ended March 31, 2021 | Year ended March 31, 2021 | Year ended March 31, 2021 | Year ended March 31, 2021 | Year ended March 31, 2021 | Year ended March 31, 2021 | Year ended March 31, 2021 | Year ended March 31, 2021 | Year ended March 31, 2021 |\n
|--|--|--|--|--|--|--|--|--|\n
| | China  commerce(1) | International  commerce | Local consumer  services(1) | Cainiao | Cloud | Digital  media and entertainment | Innovation  initiatives  and  others | Unallocated(2) | Consolidated |\n
| | RMB | RMB | RMB | RMB | RMB | RMB | RMB | RMB | RMB |\n
| | (in millions, except percentages) | (in millions, except percentages) | (in millions, except percentages) | (in millions, except percentages) | (in millions, except percentages) | (in millions, except percentages) | (in millions, except percentages) | (in millions, except percentages) | (in millions, except percentages) |\n
| Revenue | 501,379 | 48,851 | 35,746 | 37,258 | 60,558 | 31,186 | 2,311 | — | 717,289 |\n
| Income (Loss) from operations | 197,232 | (9,361) | (29,197) | (3,964) | (12,479) | (10,321) | (7,802) | (34,430) | 89,678 |\n
| Add: Share-based compensation  expense | 14,505 | 4,223 | 4,972 | 1,956 | 10,205 | 3,281 | 2,518 | 8,460 | 50,120 |\n
| Add: Amortization of intangible assets | 1,922 | 206 | 7,852 | 1,195 | 23 | 922 | 83 | 224 | 12,427 |\n
| Add: Anti-monopoly Fine(3) | — | — | — | — | — | — | — | 18,228 | 18,228 |\n
| Adjusted EBITA | 213,659 | (4,932) | (16,373) | (813) | (2,251) | (6,118) | (5,201) | (7,518) | 170,453 |\n
| Adjusted EBITA margin | 43% | (10)% | (46)% | (2)% | (4)% | (20)% | (225)% | N/A | 24% |\n

**[Chunk 2]**

(1) Beginning on October 1, 2022, we reclassified the results of our Instant Supermarket Delivery (全能超市) business, which was previously reported under China commerce segment, to Local consumer services segment following the strategy refinement of Instant Supermarket Delivery business to focus on building customer mindshare for grocery delivery services through Ele.me platform. This reclassification conforms to the way that we manage and monitor segment performance. Comparative figures were reclassified to conform to this presentation.\n
(2) Unallocated expenses primarily relate to corporate administrative costs and other miscellaneous items that are not allocated to individual segments. The goodwill impairment, and the equity-settled donation expense related to the allotment of shares to a charitable trust, are presented as unallocated items in the segment information because our management does not consider these as part of the segment operating performance measure.\n
(3) For a description of the relevant PRC Anti-monopoly investigation and administrative penalty decision, see "Business Overview — Legal and Administrative Proceedings — PRC Anti-monopoly Investigation and Administrative Penalty Decision.'"\n
Non-GAAP Measures\n
We use adjusted EBITDA (including adjusted EBITDA margin), adjusted EBITA (including adjusted EBITA margin), non-GAAP net income, non-GAAP diluted earnings per share/ADS and free cash flow, each a non-GAAP financial measure, in evaluating our operating results and for financial and operational decision-making purposes.\n

**[Chunk 3]**

We believe that adjusted EBITDA, adjusted EBITA, non-GAAP net income and non-GAAP diluted earnings per share/ADS help identify underlying trends in our business that could otherwise be distorted by the effect of certain income or expenses that we include in income from operations, net income and diluted earnings per share/ADS. We believe that these non-GAAP measures provide useful information about our core operating results, enhance the overall understanding of our past performance and future prospects and allow for greater visibility with respect to key metrics used by our management in its financial and operational decision-making. We present three different income measures, namely adjusted EBITDA, adjusted EBITA and non-GAAP net income in order to provide more information and greater transparency to investors about our operating results.\n
We consider free cash flow to be a liquidity measure that provides useful information to management and investors about the amount of cash generated by our business that can be used for strategic corporate transactions, including investing in our new business initiatives, making strategic investments and acquisitions and strengthening our balance sheet.\n

**[Chunk 4]**

Adjusted EBITDA, adjusted EBITA, non-GAAP net income, non-GAAP diluted earnings per share/ADS and free cash flow should not be considered in isolation or construed as an alternative to income from operations, net income, diluted earnings per share/ADS, cash flows or any other measure of performance or as an indicator of our operating performance. These non-GAAP financial measures presented here do not have standardized meanings prescribed by U.S. GAAP and may not be comparable to similarly titled measures presented by other companies. Other companies may calculate similarly titled measures differently, limiting their usefulness as comparative measures to our data.\n
\<Page Footer\> 112 Alibaba Group Holding Limited\n

**Figure 6.** Parsing and chunking results of ChatDOC PDF Parser on Case 1 (original document: [4]). Zoom in to see the details.

using LLM for evaluation, it may fail to distinguish the subtle but important differences between answers, so we relied on human assessment. We used a 0-10 scale to rate the results. An annotator is given the retrieved content and answer of both methods and rates the two methods at the same time. We show the retrieved content because it usually cannot evaluate the answer without document content, and show two methods together to promote detailed comparison, especially on partially correct results.

- **Comprehensive analysis questions** necessitate synthesizing information from multiple sources and aspects and making a summary. Since the answer is lengthy and requires a comprehensive understanding of the given document contents, we found it difficult and time-consuming for humans to evaluate. Hence, we used GPT-4 to evaluate the answer quality, scoring from 1-10. We also rate the result of two methods in one request. But we only give the retrieved content without an answer because the answer is lengthy (thus costly) compared with extractive questions and a better retrieved content can imply a better answer (since the used LLM is the same). A pair of results of two methods is scored 4 times to avoid bias [5], and their average value is used. Specifically, for a pair of content $(A, B)$ to be compared for the same question, we feed both

| Steps ↓ | ChatDOC<br>(PDFlux-LLM) | Baseline<br>(PyPDF-LLM) |
|---|---|---|
| PDF Parsing | PDFlux<br>(Deep Learning-based) | PyPDF<br>(Rule-based, default method in LangChain) |
| Chunking | ≈300 tokens per chunk<br>+ chunking via paragraphs, tables etc. | ≈300 tokens per chunk + separator |
| Embedding | text-embedding-ada-002 | |
| Retrieval | ≤3000 tokens | |
| QA | GPT3.5-Turbo | |

**Table 1.** Settings of two RAG systems: ChatDOC and Baseline.

| | Extractive Questions | Comprehensive Analysis Questions |
|---|---|---|
| Number | 86 | 216 |
| Question Examples | *1. Locate the content of section ten, what is the merged operating cost in the income statement?*<br>*2. What is the specific content of table 1.*<br>*3. Extract financial data and profit forecast tables.*<br>*4. Find the long-term loan table.* | *1. Summarize and analyze the profit forecast and valuation in the research report.*<br>*2. Fully report the research approach of this text.*<br>*3. Analyze the long-term debt-paying ability based on this report.*<br>*4. How is the feasibility analysis done in this article?*<br>*5. Give a simple example to explain the encoding steps and algorithm in the paper.* |
| Evaluation | Human Evaluation | GPT 4 evaluation |

**Table 2.** The questions in the dataset are categorized into extractive questions and comprehensive analysis questions.

$A$ and $B$ to GPT-4 to compare and score them twice. We also flip their order, feed $B$ and $A$ to GPT-4, and repeat the request twice.

### 3.1.3 Results

**Results of Extractive Questions**

The results of extractive questions are shown in Table 3. Out of the 86 extractive questions, Chat-DOC performed better than the baseline on 42 cases, tied on 36 cases, and was inferior to Baseline on only 8 cases.

The distribution of rating scores is further detailed in Figure 7. In the distribution table, $T_{ij} = k$ means there are $k$ questions whose answer by ChatDOC is rated as $i$ and the answer by Baseline is rated as $j$. Cases where ChatDOC scores higher than the baseline (ChatDOC wins) are represented in the lower-left half, while cases where the baseline scores higher are in the upper-right. Notably, most samples with a clear winner are in the lower-left half, indicating ChatDOC's superiority. Impressively, ChatDOC achieved full marks (10) in nearly half of these cases, amounting to a total of 40.

**Results of Comprehensive Analysis Questions**

| | Total | ChatDOC wins | Tie | Baseline wins |
|---|---|---|---|---|
| **Extractive Questions** | 86 | **42** **(49%)** | 36 (42%) | 8 (9%) |
| **Comprehensive Questions** | 216 | **101** **(47%)** | 79 (37%) | 36 (17%) |
| **Summary** | 302 | **143** **(47%)** | 115 (38%) | 44 (15%) |

**Table 3.** The comparison result between ChatDOC and Baseline.

**Score of Baseline**

| Score of ChatDOC | 0-2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|
| 0-2 | 5 | 2 | 0 | 1 | 0 | 0 | 1 | 0 | 1 |
| 3 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5 | 3 | 1 | 0 | 3 | 1 | 0 | 0 | 0 | 0 |
| 6 | 1 | 0 | 2 | 1 | 0 | 0 | 0 | 0 | 0 |
| 7 | 2 | 0 | 0 | 3 | 0 | 3 | 0 | 0 | 0 |
| 8 | 1 | 0 | 0 | 3 | 0 | 1 | 2 | 0 | 0 |
| 9 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 4 | 0 |
| 10 | 7 | 3 | 2 | 2 | 0 | 2 | 3 | 1 | 20 |

◿ ChatDOC Wins

**Figure 7.** Distribution of rating scores of extractive questions.

**Score of Baseline**

| Score of ChatDOC | 0-5 | 5-6 | 6-7 | 7-7.5 | 7.5-8 | 8-8.5 | 8.5-9 | 9-9.5 | 9.5-10 |
|---|---|---|---|---|---|---|---|---|---|
| 0-5 | 8 | 2 | 0 | 0 | 1 | 0 | 1 | 2 | 1 |
| 5-6 | 3 | 3 | 1 | 0 | 0 | 2 | 3 | 0 | 0 |
| 6-7 | 1 | 2 | 4 | 2 | 0 | 1 | 1 | 1 | 0 |
| 7-7.5 | 0 | 0 | 3 | 0 | 0 | 3 | 3 | 2 | 0 |
| 7.5-8 | 1 | 1 | 5 | 1 | 5 | 7 | 4 | 0 | 0 |
| 8-8.5 | 0 | 2 | 4 | 4 | 4 | 11 | 5 | 2 | 0 |
| 8.5-9 | 1 | 0 | 1 | 11 | 7 | 21 | 12 | 4 | 1 |
| 9-9.5 | 0 | 0 | 2 | 10 | 8 | 8 | 9 | 7 | 1 |
| 9.5-10 | 0 | 0 | 0 | 0 | 0 | 3 | 2 | 2 | 0 |

◿ ChatDOC Wins

**Figure 8.** Distribution of rating scores of comprehensive analysis questions.

The results of comprehensive analysis questions are shown in Table 3. Out of the 216 comprehensive analysis questions, ChatDOC performed better than the baseline on 101 cases, tied on 79 cases, and was inferior to Baseline on only 36 cases.

In Figure 8, the score distribution chart for these questions shows a denser concentration in the lower-left half. This indicates that ChatDOC frequently outperforms the baseline. Notably, the majority of ChatDOC's retrievals scored between 8.5 and 9.5, demonstrating a high level of retrieval quality.

In summary, ChatDOC significantly outperforms the Baseline, primarily attributed to its superior PDF parsing method.

## 3.2 Case Studies of RAG

To make the comparison more tangible, we show some cases where ChatDOC demonstrates its superiority.

**Case A: ChatDOC**

**Query: How much cargo can I carry at most in terms of size?**

**Retrieved Text Chunk:**
[Chunk 1]

| Cargo Volume | |\n
|-|-|\n
| Front Trunk | 3.1 cu ft (88 L) |\n
| Behind 2nd row | 19.8 cu ft (561 L) |\n
| Maximum total cargo volume with 5 passengers | 22.9 cu ft (649 L) |\n
...

**Answer:**

Based on the document, the maximum cargo capacity in terms of size for the Model 3 is as follows:
- Front Trunk: 3.1 cu ft (88 L)
- Behind 2nd row: 19.8 cu ft (561 L)
- Maximum total cargo volume with 5 passengers: 22.9 cu ft (649 L)

**Original Page Containing the Retrieved Chunk:**

**Retrieved Chunk Highlighted:**

☐ Table

**Figure 9.** ChatDOC's result on finding info in Tesla's Manual (original document: [6]).

**Case A: Baseline (PyPDF + LangChain)**

**Query: How much cargo can I carry at most in terms of size?**

**Retrieved Text Chunk:**
[Chunk 1]

Hip Room Front Rear 53.4 in 52.4 in 1,356 mm 1,331 mm Cargo Volume Front Trunk 3.1 cu ft (88 L) Behind 2nd row 19.8 cu ft (561 L) Maximum total cargo volume with 5 passengers 22.9 cu ft (649 L)
...

**Answer:**

The maximum cargo size capacity for your Model 3 is not explicitly mentioned in the provided document snippets.

**Original Page Containing the Retrieved Chunk:**

**Retrieved Chunk Highlighted:**

☐ Text Chunk

**Figure 10.** Baseline's result on finding info in Tesla's Manual (original document: [6]).

**Case B: ChatDOC**

**Query: What is the dependent variable in table 8?**

**Retrieved Text Chunks:**

[Chunk 1]

TABLE 8ESTIMATES OF RETURNS TO COGNITIVE ABILITY, NONCOGNITIVE ABILITY, AND FACTOR LOADINGS OF UNOBSERVEDHETEROGENEITY IN HIGH SCHOOL AND COLLEGE

[Chunk 2]

|Cognitive ability| Cognitive ability| Cognitive ability| Noncognitive ability| Noncognitive ability| Unobs. heterogeneity| Unobs. heterogeneity|\n

|-|-|-|-|-|-|-|\n

| Year| Estimate| Std. error| Estimate| Std. error| Estimate| Std. error|\n

| High school| High school| High school| High school| High school| High school| High school|\n

| 1968| 0.038| 0.021| 0.064| 0.024| 0.000| 0.030|\n

...

| 2000| 0.047| 0.025| 0.091| 0.028| 0.225| 0.023|\n

| College| College| College| College| College| College| College|\n

| 1968| 0.085| 0.038| 0.082| 0.045| -0.129| 0.054|\n

...

| 2000| 0.151| 0.039| 0.100| 0.037| 0.079| 0.037|\n

...

**Answer:**

The dependent variable in Table 8 is the estimate of returns to cognitive ability, noncognitive ability, and factor loadings of unobserved heterogeneity in high school and college.

**Original Page Containing the Retrieved Chunks (Correct):**

**Retrieved Chunks Highlighted:**

☐ Paragraph    ☐ Table

**Figure 11.** ChatDOC's result on locating a specific table in a research paper (original document: [7]).

### 3.2.1 Case A - Find Specific information in the Tesla manual

Case A involves a query about Tesla's user manual specifically targeting cargo volume information. For this query, ChatDOC and Baseline perform differently as in Figure 9 and Figure 10. The figures show the most relevant chunk(s) retrieved and the LLM's answer. They also show the document page that the relevant chunk(s) are located in, and highlight these chunks. In this case, both models located the table, but the text they feed to LLM is different, so the answer is different. Specifically,

- ChatDOC recognizes the table structure, interpreting the text in the `markdown` format (as shown in the "Retrieved Text Chunks" part), which made it easier for the language model to comprehend.

11

**Figure 12.** Baseline's result in locating a specific table in a research paper (original document: [7])

- Baseline erroneously merges the target table and the table above into one chunk and does not have the table structure. Hence, the text in the chunk is not understandable (as shown in the "Retrieved Text Chunks" part) and the LLM can only answer with "not explicitly mentioned".

This case underscores the effectiveness of ChatDOC's parsing method, particularly in handling tables and presenting them in an LLM-friendly format.

### 3.2.2 Case B - Research paper

In Case B, the user's query is on a specific research paper. It requests the system to identify "Table 8" in the paper and enumerate all the dependent variables it lists. Both the title and the content of the table were necessary for identifying these variables. Figure 11 and Figure 12 show how ChatDOC and Baseline perform in this case.

- ChatDOC effectively retrieves the entire table, encompassing both its title and content. This comprehensive retrieval allows for an accurate response to the query.

**Figure 13.** An example of ChatDOC encountered the ranking and token limit issues.



**Figure 14.** An example that ChatDOC fails to retrieve the relevant table (original document: [8]).

- Baseline does not retrieve true "Table 8", but only a text chunk below "Table 7" (since it contains the text of "Table 8"). Due to the baseline's segmentation strategy, the content of "Table 8" and other content on the same page are combined into a large chunk. This chunk, containing a mix of unrelated content, has a low similarity score and consequently does not show up in the retrieval results.

This case highlights ChatDOC's superior ability to handle complex document structures and its impact on retrieving specific segments for accurate responses.

### 3.3 Discussion on Limitations

While ChatDOC generally performs well, there are instances where its retrieval quality is not as good as Baseline's. We observe two patterns in these cases.

**Ranking and Token Limit Issue.** If ChatDOC retrieves a large, but irrelevant table first, it uses up the context window, preventing access to the relevant information, as the example in Figure 13 shows. This is mainly because the embedding model does not rank the relevant chunk as the top result. This may be addressed by a better embedding model, or a more sophisticated way to handle large tables/paragraphs like only retaining the relevant part of the table for LLM.

**Fine Segmentation Drawback.** Figure 14 shows a case that requires retrieving the whole table with its title. However, ChatDOC wrongly recognizes the title as a regular paragraph, so that the title and the table are stored in different chunks. This led to retrieving only part of the required information, namely the table's title and footnotes, but not the key content within the table. Improving table title recognition could address these issues.

## 4 Applications in ChatDOC

We apply the enhanced PDF structure recognition framework on ChatDOC (chatdoc.com), an AI file-reading assistant that helps to summarize long documents, explain complex concepts, and find key information in seconds.

In terms of reliability and accuracy, it is the top among all ChatPDF products. Here's what makes ChatDOC special:

- Mastery over tables: Simply select any table or text, and dive right into the details.
- Multi-file conversation: Talk about lots of documents at the same time, without worrying about how many pages each one has.
- Citation-backed responses: All answers are supported by direct quotes pulled from the source documents.
- Handle Many File Types: Works seamlessly with scanned files, ePub, HTML, and docx formats.

We are still working on publishing the API of ChatDOC PDF Parser. Please subscribe to the wait list via pdfparser.io.

## 5 Conclusion

Large Language Models (LLMs) are capable of producing more accurate responses when assisted by a PDF parser that effectively extracts and integrates structured information from documents into the prompts. This process enhances the quality and relevance of the data fed into the models, thereby improving their output.

In the future, we will compare more deep learning-based document parsing methods to give a more comprehensive understanding of the relationship between the RAG quality and document parsing quality. Some initial experiments show that some open-sourced PDF parsing methods cannot meet the bar for high-quality RAG.

### References

[1] Alibaba Group Holding Limited. Fiscal year annual report 2023. `https://static.alibabagroup.com/reports/fy2023/ar/ebook/en/index.html`, 2023.

[2] Rongyu Cao, Hongwei Li, Ganbin Zhou, and Ping Luo. Towards document panoptic segmentation with pinpoint accuracy: Method and evaluation. In *16th International Conference on Document Analysis and Recognition*, pages 3–18, 2021.

[3] ChatDOC Team. https://pdfparser.io/.

[4] Daisho Microline Holdings Limited. Fiscal year annual report 2022. https://www1.hkexnews.hk/listedco/listconews/sehk/2022/0626/2022062600094.pdf, 2022.

[5] Peiyi Wang, Lei Li, Liang Chen, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and Zhifang Sui. Large language models are not fair evaluators, 2023.

[6] Tesla Inc. Model 3 owner's manual. https://manual-directory.com/manual/2023-tesla-model-3-owners-manual/, 2023.

[7] Flávio Cunha, Fatih Karahan, and Ilton Soares. Returns to skills and the college premium. *Journal of Money, Credit and Banking*, 43:39–86, 2011. https://sci-hub.hkvisa.net/https://doi.org/10.1111/j.1538-4616.2011.00410.x.

[8] Tom S. Vogl. Height, skills, and labor market outcomes in mexico. *NBER Working Paper Series*, 2012. https://www.nber.org/system/files/working_papers/w18318/w18318.pdf.

# A  More Cases on PDF Parsing & Chunking

**Case 2**    in Figure 15 features a large borderless table that spans two pages. Figure 15 shows the result by PyPDF. A close inspection reveals that tables are represented merely as sequences of text, making them challenging to interpret and understand. And the table is scattered in three chunks. Results on these two cases demonstrate that the rule-based method, like that of PyPDF, tends to dissect a document without a true understanding of its content structure. As a result, tables are often torn apart and paragraphs become jumbled, leading to a disjointed and confusing representation of the original document.

For ChatDOC PDF Parser, shown in Figure 16, the parsing outcome is notably different. It not only preserves the document structure but also effectively segments the document in a way that maintains its inherent meaning. In this case, the table that spans two pages is set into one chunk, with its title at the beginning. So, the information in this chunk is self-contained. If this chunk is retrieved for RAG, the LLM can digest useful information within it.

**Case 2: PyPDF**

**1 Original Pages:**



**2 Chunking Result:**

**[Chunk 1]**

1\n
Hong Kong Exchanges and Clearing Limited and The Stock
Exchange of Hong Kong \n
Limited take no responsibility for the contents of this
announcement, make no \n
representation as to its accuracy or completeness and
expressly disclaim any liability \n
whatsoever for any loss howsoever arising from or in
reliance upon the whole or any \n
part of the contents of this announcement.\n

**[Chunk 2]**

DAISHO MICROLINE HOLDINGS LIMITED\n
(Incorporated in Bermuda with limited liability)\n
(Stock Code: 0567)\n
ANNOUNCEMENT OF ANNUAL RESULTS\n
FOR THE YEAR ENDED 31 MARCH 2022\n
The Board of Directors (the "Board") of Daisho Microline
Holdings Limited (the \n
" Company") announces the preliminary consolidated
results of the Company and its \n
subsidiaries (the "Group") for the year ended 31 March 2022
together with the \n
comparative figures of the previous corresponding year as
follows:\n
CONSOLIDATED STATEMENT OF PROFIT OR LOSS\n
Year ended 31 March 2022\n
2022 2021\n
Note HK$'000 HK$'000\n
Continuing operations\n
Revenue 3 106,471 67,886\n
Cost of sales (98,670) (55,605)\n
Gross profit 7,801 12,281\n
Other income 5 7,341 4,616\n
Selling and distribution expenses (5,083) (3,401)\n
Administrative expenses (31,157) (35,422)\n
Other operating expenses (480) (527)\n
Fair value gain on derivative financial instruments – 101\n
Reversal of (Provision for) impairment loss on\n

**[Chunk 3]**

trade receivables, net 10(b) 1,808 (2,859)\n
Impairment loss on other receivables – (1,780)\n
Impairment loss on property, plant and equipment 15 (5,010)
(2,314)\n
Change in fair value of contingent consideration \n
receivable – 3,311\n
Gain on bargain purchase arising from the \n
acquisition of subsidiaries – 1,197\n
Loss on early redemption of a promissory note – (4,512)\n
Finance costs 6 (2,244) (7,655)\n

**[Chunk 4]**

2\n
2022 2021\n
Note HK$'000 HK$'000\n
Loss before taxation from continuing operations 6 (27,024)
(36,964)\n
Income tax expense 7 (444) (532)\n
Loss for the year from continuing operations (27,468)
(37,496)\n
Discontinued operation\n
Loss for the year from discontinued operation 11 (1,660)
(29,480)\n
Loss for the year (29,128) (66,976)\n
From continuing and discontinued operations\n
Loss per share\n
Basic (Hong Kong cents) 8 (2.80) (10.38)\n
Diluted (Hong Kong cents) 8 (2.80) (10.38)\n
From continuing operations\n
Loss per share\n
Basic (Hong Kong cents) 8 (2.64) (5.81)\n
Diluted (Hong Kong cents) 8 (2.64) (5.81)\n

**3 Visualization of Chunking Result:**



▭ **Text Chunk**

**Figure 15.** Parsing and chunking results of PyPDF on Case 2 (original document: [4]).
Zoom in to see the details.

**Case 2: ChatDOC PDF Parser**

| **1** Original Pages: | **2** Chunking Result: | **3** Visualization of Chunking Result: |

**2** Chunking Result:

[Chunk 1]

Hong Kong Exchanges and Clearing Limited and The Stock Exchange of Hong Kong Limited take no responsibility for the contents of this announcement, make no representation as to its accuracy or completeness and expressly disclaim any liability whatsoever for any loss howsoever arising from or in reliance upon the whole or any part of the contents of this announcement.\n
DAISHO MICROLINE HOLDINGS LIMITED\n
(Incorporated in Bermuda with limited liability)\n
(Stock Code: 0567)\n
ANNOUNCEMENT OF ANNUAL RESULTS\n
FOR THE YEAR ENDED 31 MARCH 2022\n
The Board of Directors (the "Board") of Daisho Microline Holdings Limited (the "Company") announces the preliminary consolidated results of the Company and its subsidiaries (the "Group") for the year ended 31 March 2022 together with the comparative figures of the previous corresponding year as follows:\n

[Chunk 2]

CONSOLIDATED STATEMENT OF PROFIT OR LOSS\n
Year ended 31 March 2022\n
| | Note | 2022 | 2021 |\n
|-|-|-|-|\n
| | Note | HK\$'000 | HK\$'000 |\n
| Continuing operations | | | |\n
| Revenue | 3 | 106,471 | 67,886 |\n
| Cost of sales | | (98,670) | (55,605) |\n
| Gross profit | | 7,801 | 12,281 |\n
| Other income | 5 | 7,341 | 4,616 |\n
| Selling and distribution expenses | | (5,083) | (3,401) |\n
| Administrative expenses | | (31,157) | (35,422) |\n
| Other operating expenses | | (480) | (527) |\n
| Fair value gain on derivative financial instruments | | – | 101 |\n
| Reversal of (Provision for) impairment loss on trade receivables, net | 10(b) | 1,808 | (2,859) |\n
| Impairment loss on other receivables | | – | (1,780) |\n
| Impairment loss on property, plant and equipment | 15 | (5,010) | (2,314) |\n
| Change in fair value of contingent consideration receivable | | – | 3,311 |\n
| Gain on bargain purchase arising from the acquisition of subsidiaries | | – | 1,197 |\n
| Loss on early redemption of a promissory note | | – | (4,512) |\n
| Finance costs | 6 | (2,244) | (7,655) |\n
| Loss before taxation from continuing operations | 6 | (27,024) | (36,964) |\n
| Income tax expense | 7 | (444) | (532) |\n
| Loss for the year from continuing operations | | (27,468) | (37,496) |\n
| Discontinued operation | | | |\n
| Loss for the year from discontinued operation | 11 | (1,660) | (29,480) |\n
| Loss for the year  | | (29,128) | (66,976) |\n
| From continuing and discontinued operations | | | |\n
| Loss per share | | | |\n
| Basic (Hong Kong cents) | 8 | (2.80) | (10.38) |\n
| Diluted (Hong Kong cents) | 8 | (2.80) | (10.38) |\n
| From continuing operations | | | |\n
| Loss per share | | | |\n
| Basic (Hong Kong cents) | 8 | (2.64) | (5.81) |\n
| Diluted (Hong Kong cents) | 8 | (2.64) | (5.81) |\n

**Figure 16.** Parsing and chunking result of ChatDOC PDF Parser on Case 2 (original document: [4]). Zoom in to see the details.