

ZERO-TO-HERO: PARSING FOR RETRIEVAL AUGMENTED GENERATION (RAG)

Yogesh Haribhau Kulkarni



Outline

① INTRODUCTION

② LIBRARIES

③ IMPLEMENTATION

④ CONCLUSIONS

About Me

Yogesh Haribhau Kulkarni

Bio:

- ▶ 20+ years in CAD/Engineering software development
- ▶ Got Bachelors, Masters and Doctoral degrees in Mechanical Engineering (specialization: Geometric Modeling Algorithms).
- ▶ Currently doing Coaching in fields such as Data Science, Artificial Intelligence Machine-Deep Learning (ML/DL) and Natural Language Processing (NLP).
- ▶ Feel free to follow me at:
 - ▶ Github (github.com/yogeshhk)
 - ▶ LinkedIn (www.linkedin.com/in/yogeshkulkarni/)
 - ▶ Medium (yogeshharibhaukulkarni.medium.com)
 - ▶ Send email to [yogeshkulkarni at yahoo dot com](mailto:yogeshkulkarni@yahoo.com)



Office Hours:
Saturdays, 2 to 5pm
(IST); Free-Open to all;
email for appointment.

Parsing is the key

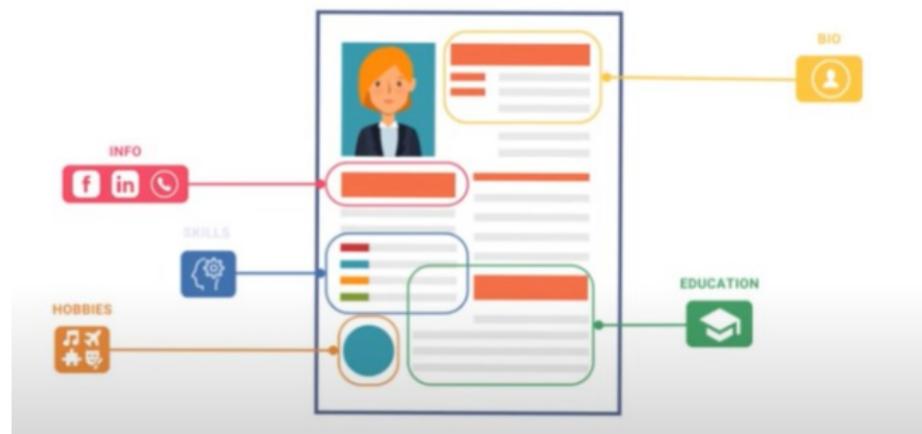
(Ref: Key to RAG Success: Document Parsing Explained - EyeLevel)

Document Parsing: The Foundation of RAG

- ▶ Parsing is the first step in any RAG pipeline.
- ▶ Bad parsing undermines even the best RAG strategies.
- ▶ Garbage in, garbage out: poor inputs = poor outputs.
- ▶ Many overlook parsing in favor of flashy AI tools.
- ▶ Without good extraction, nothing else matters.
- ▶ Most language models require clean, structured text.
- ▶ RAG applications depend on text quality from source docs.
- ▶ Advanced RAG still fails without reliable input.
- ▶ Models can't fix broken, messy data.
- ▶ Real-world systems have failed due to poor parsing.

What is Document Parsing?

- ▶ Converts formats like PDF, DOCX, HTML into usable text.
- ▶ Extracts meaningful content for language model input.
- ▶ Cleans, structures, and normalizes the data.
- ▶ Essential step before chunking or embedding.
- ▶ Involves handling many formats and edge cases.



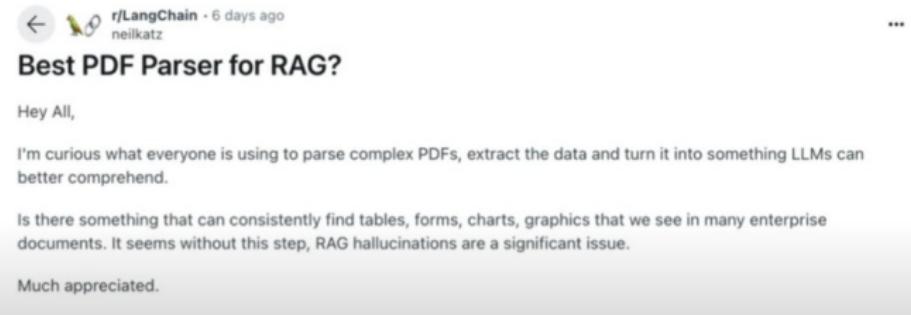
(Ref: Key to RAG Success: Document Parsing Explained - EyeLevel)

Common Misconceptions

- ▶ Engineers often ignore parsing during development.
- ▶ Focus tends to be on model tuning or retrieval logic.
- ▶ Parsing is wrongly assumed to be solved or trivial.
- ▶ Most systems lack formal evaluation of parsers.
- ▶ Homemade or ad-hoc solutions dominate practice.

The Reddit Survey Insight

- ▶ Survey on LangChain subreddit revealed no consensus.
- ▶ 57 replies yielded 30+ different parsing techniques.
- ▶ Most users hacked together informal solutions.
- ▶ Few performed proper parser evaluation or comparison.
- ▶ Highlights need for standardized testing and benchmarking.



A screenshot of a Reddit post from the r/LangChain subreddit. The post was made 6 days ago by user neilkatz. The title is "Best PDF Parser for RAG?". The post content asks for recommendations on the best PDF parser for RAG models, mentioning the challenge of extracting tables, forms, charts, and graphics from complex PDFs to improve model comprehension. A response from another user expresses appreciation for the question.

← r/LangChain • 6 days ago
neilkatz

...

Best PDF Parser for RAG?

Hey All,

I'm curious what everyone is using to parse complex PDFs, extract the data and turn it into something LLMs can better comprehend.

Is there something that can consistently find tables, forms, charts, graphics that we see in many enterprise documents. It seems without this step, RAG hallucinations are a significant issue.

Much appreciated.

(Ref: Key to RAG Success: Document Parsing Explained - EyeLevel)

Popular Parsing Tools Compared

- ▶ PyPDF – well-known, older, basic PDF extraction.
- ▶ Tesseract – OCR-based, handles scanned documents.
- ▶ Unstructured – handles messy formats, layout-aware.
- ▶ Tools vary widely in output quality and reliability.
- ▶ Choice depends on document type and project needs.

- ▶ Start by identifying your document types.
- ▶ Evaluate parsers with real-world examples.
- ▶ Compare outputs side by side.
- ▶ Look for structural fidelity, cleanliness, completeness.
- ▶ Test rigorously — don't rely on "it seems to work".

Real-World Example: Medical Bill

- ▶ Parsing tested on de-identified medical bill.
- ▶ Chosen for layout complexity and format irregularities.
- ▶ Shows strengths and weaknesses of each parser.
- ▶ Realistic example of what RAG apps encounter.
- ▶ Highlights need for resilient parsing strategies.

The image shows a "SUMMARY BILL OF ALL CHARGES" document from University Imaging Center. The header includes the center's logo, address (5405 S. Trent street, Ontario, CA 91761), and contact information (PHONE: 909/654-3676, FAX: 909/654-3682). The document is dated 10/17/2002. It lists two categories of charges: "TECHNICAL IMAGING PROCEDURES AND BILLING" and "PROFESSIONAL PROCEDURES AND BILLING". Under "TECHNICAL IMAGING", two services are listed: "MRI NECK SPINE W/O DYE" and "MRI LUMBAR SPINE W/O DYE", both performed by Dr. Vines. Under "PROFESSIONAL", a service "Professional Component (Radiology Report) Interpreted by: Dr. Vines" is listed. A summary table at the bottom shows Total Charges (\$4,400.00), Total Payments (\$1.00), Total Adjustments (\$0.00), and Balance Due (\$4,400.00). A barcode is present at the bottom right, along with instructions: "ADDITIONAL OR REVISED BILLING MAY OCCUR" and "PLEASE EMAIL BILLING@UIC.EDU FOR FINAL BILLING AMOUNT".

EXAM DATE	PROC. CODE	DESCRIPTION	MOD	PAYMENT	DR	AMOUNT
10/17/2002	73141	MRI NECK SPINE W/O DYE	TC	COPAY	RECD	\$ 1,000.00
10/17/2002	73140	MRI LUMBAR SPINE W/O DYE	TC	COPAY	RECD	\$ 1,000.00

EXAM DATE	PROC. CODE	DESCRIPTION	MOD	PAYMENT	DR	AMOUNT
10/17/2002	73141	MRI NECK SPINE W/O DYE	RE	COPAY	RECD	\$ 1,000.00
10/17/2002	73140	MRI LUMBAR SPINE W/O DYE	RE	COPAY	RECD	\$ 1,000.00

Total Charges	Total Payments	Total Adjustments	Balance Due
\$4,400.00	\$1.00	\$0.00	\$4,400.00

Internal use Only: 909/654-3682

ADDITIONAL OR REVISED BILLING MAY OCCUR
PLEASE EMAIL BILLING@UIC.EDU FOR FINAL BILLING AMOUNT

Limitations of PyPDF

- ▶ PyPDF ok for texts but struggles with complex formats like tables.
- ▶ It failed to extract key data from real-world docs.
- ▶ Parsing tables often results in empty or broken content.
- ▶ Not due to PyPDF's fault—PDFs are inherently hard to parse.
- ▶ Many PDFs use inconsistent encoding and layouts.

Grant, Victoria 5436 S. Trent street Montclair, CA 91763	TAX ID# 989376323 FEDERAL TAX IDENTIFICATION NO: 7635 KESTER AVE SUITE 345 VAN NUYS CA 91405	UNIVERSITY IMAGING CENTER	SUMMARY BILL OF ALL CHARG TAX ID# 989376323 PLEASE PRINT PAYMENT TO: 7635 KESTER AVE SUITE 345 VAN NUYS CA 91405																		
Dr. Vince	Dr. Vince	BILLING STATE Patient ID Status PRE-RETURN 11/01	ALWAYS REFERENCE/FI PHONE : 818.654.5876 DOB : 05/16/1982																		
TECHNICAL IMAGING PROCEDURES AND BILLINGS																					
<table border="1"><thead><tr><th>EXAM DATE</th><th>PROC. CODE</th><th>DESCRIPTION</th><th>MOD</th><th>R.PART</th><th>DR</th></tr></thead><tbody><tr><td>7632 Central Street, Montclair CA 91763</td><td>7/24/2023</td><td>MRI NECK SPINE W/O DYE</td><td>TC</td><td>ESPINE</td><td>MSA.2</td></tr><tr><td>7632 Central Street, Montclair CA 91763</td><td>7/24/2023</td><td>MRI LUMBAR SPINE W/O DYE</td><td>TC</td><td>LSPINE</td><td>MSA.3</td></tr></tbody></table>				EXAM DATE	PROC. CODE	DESCRIPTION	MOD	R.PART	DR	7632 Central Street, Montclair CA 91763	7/24/2023	MRI NECK SPINE W/O DYE	TC	ESPINE	MSA.2	7632 Central Street, Montclair CA 91763	7/24/2023	MRI LUMBAR SPINE W/O DYE	TC	LSPINE	MSA.3
EXAM DATE	PROC. CODE	DESCRIPTION	MOD	R.PART	DR																
7632 Central Street, Montclair CA 91763	7/24/2023	MRI NECK SPINE W/O DYE	TC	ESPINE	MSA.2																
7632 Central Street, Montclair CA 91763	7/24/2023	MRI LUMBAR SPINE W/O DYE	TC	LSPINE	MSA.3																
PROFESSIONAL PROCEDURES AND BILLINGS																					
<table border="1"><thead><tr><th>EXAM DATE</th><th>PROC. CODE</th><th>DESCRIPTION</th><th>MOD</th><th>R.PART</th><th>DR</th></tr></thead><tbody><tr><td>7632 Central Street, Montclair CA 91763</td><td>7/24/2023</td><td>MRI NECK SPINE W/O DYE</td><td>26</td><td>ESPINE</td><td>MSA.2</td></tr><tr><td>7632 Central Street, Montclair CA 91763</td><td>7/24/2023</td><td>MRI LUMBAR SPINE W/O DYE</td><td>26</td><td>LSPINE</td><td>MSA.3</td></tr></tbody></table>				EXAM DATE	PROC. CODE	DESCRIPTION	MOD	R.PART	DR	7632 Central Street, Montclair CA 91763	7/24/2023	MRI NECK SPINE W/O DYE	26	ESPINE	MSA.2	7632 Central Street, Montclair CA 91763	7/24/2023	MRI LUMBAR SPINE W/O DYE	26	LSPINE	MSA.3
EXAM DATE	PROC. CODE	DESCRIPTION	MOD	R.PART	DR																
7632 Central Street, Montclair CA 91763	7/24/2023	MRI NECK SPINE W/O DYE	26	ESPINE	MSA.2																
7632 Central Street, Montclair CA 91763	7/24/2023	MRI LUMBAR SPINE W/O DYE	26	LSPINE	MSA.3																
Professional Component (Kinesiology Report) Interpreted by: Dr. Vince																					
YHK																					

Tesseract OCR: Strengths and Weaknesses

- ▶ Tesseract uses image-based OCR to extract text.
- ▶ Better at recognizing tables than PyPDF.
- ▶ Still introduces errors in column alignment.
- ▶ Column headers often get merged or misread.
- ▶ OCR also introduces spelling mistakes (e.g. "Cove" vs. "Code").

Grant, Victoria
 5436 S. Trent street PHONE: 818 6543876
 Ontario, CA 91761 DOB : 05/16/1992

i —<\$SF SS ee
 TECHNICAL IMAGING PROCEDURES AND BILLINGS
 Technical Component,(Imaging) performed by : Dr. Vince

SUMMARY BILL OF AL			
TAX ID: 000376512	PLEASE REQUEST PAYMENT TO:		
5436 S. Trent street	7655 KESTER AVE SUITE 345		
Ontario, CA 91761	VIN NOPE CA 91765		
ALINATE			
PHONE			
DOB : 1			
TECHNICAL IMAGING PROCEDURES AND BILLINGS			
EXAM DATE	PROC. CODE	DESCRIPTION	MOD R/P
7652 Central Street , Montclair CA 91763			
11/17/2022	72141	MRI NECK SPINE W/O DYE	TC C/P
11/17/2022	72148	MRI LUMBAR SPINE W/O DYE	TC C/P
PROFESSIONAL PROCEDURES AND BILLINGS			
EXAM DATE	PROC. CODE	DESCRIPTION	MOD R/P
7652 Central Street , Montclair CA 91763			
Professional Component (Radiology Report) Interpreted by: Dr. Vince			
11/17/2022	72141	MRI NECK SPINE W/O DYE	26 C/P
11/17/2022	72148	MRI LUMBAR SPINE W/O DYE	26 C/P

Challenges with OCR Outputs

- ▶ Language models must infer structure from broken text.
- ▶ Humans can "guess" meaning—models may not.
- ▶ Noisy extractions increase risk of incorrect answers.
- ▶ Inconsistent column separation confuses models.
- ▶ Clean layout is crucial for reliable RAG responses.

Unstructured: A Common Default

- ▶ Popular choice 'Unstructured' company; default in LangChain integrations.
- ▶ Handles layout better than OCR in some cases.
- ▶ Still suffers from column misalignments.
- ▶ Model must rely on context instead of structure.
- ▶ Reasonable quality, but far from perfect.

TAX ID: 0893765213 PLEASE REMIT PAYMENT TO: 7835 KESTER AVE SUITE 345 VAN NUYS CA 91405 PRE 987354 Grant, Victoria 5436 S. Trent street Ontario, CA 91761 818 654 3876 05/16/1992 Dr. Vince 7652 Central Street , Montclair (91763 Dr. Vince 7652 Central Street , Montclair CA 91763

TAX ID: 0893765213 PLEASE REMIT PAYMENT TO: 7835 KESTER AVE SUITE 345 VAN NUYS CA 91405
PRE 987354
818 654 3876 05/16/1992
Dr. Vince

EXAM DATE | PROC. CODE DESCRIPTION MOD B.PART DX AMOUNT
7652 Central Street , Montclair CA 91763
11/17/2022 72141 MRI NECK SPINE W/O DYE | TC CSPINE M54.2 \$ 1,600.00
11/17/2022 72148 MRI LUMBAR SPINE W/O DYE | Tc LSPINE M54.5 \$ 1,600.00
EXAM DATE | PROC. CODE DESCRIPTION MOD B.PART DX AMOUNT

11/17/2022 72141 MRI NECK SPINE W/O DYE 26 CSPINE M54.2 \$ 600.00 11/17/2022
72148 MRI LUMBAR SPINE W/O DYE | 26 LSPINE M54.5, \$ 600.00

ADDITIONAL OR REVISED BILLING MAY OCCUR *PLEASE EMAIL BILLING@UIC.COM TO CONFIRM FINAL BILLING AMOUNT*

Total Charges Total Payments Total Adjustments Balance Due \$4,400.00 \$0.00 \$0.00 \$4,400.00 Internal use Only: 9/27/2022 | i il All! | i il *ADDITION OR REVISED BILLING MAY OCCUR*

ADDITIONAL OR REVISED BILLING MAY OCCUR* *PLEASE EMAIL BILLING@UIC.COM TO CONFIRM FINAL BILLING AMOUNT*

FROST LAW PRE 987354 12/13/2022 GRANT, VICTORIA As the authorized representative of, VICTORIA GRANT, we are informing you of a new lien GRANT VICTORIA assignment. Should you not be the current authorized representative



SUMMARY BILL OF AL
TAX ID: 0893765213
PLEASE REMIT PAYMENT TO:
7835 KESTER AVE SUITE 345
VAN NUYS CA 91405

Grant, Victoria
5436 S. Trent street
Ontario, CA 91761

ALWATX
PHONE :
DOB :
FAX :

TECHNICAL IMAGING PROCEDURES AND BILLINGS

EXAM DATE	PROC. CODE	DESCRIPTION	MOD	B.PA
7652 Central Street, Montclair CA 91763				
11/17/2022 72141		MRI NECK SPINE W/O DYE	TC	CSP
11/17/2022 72148		MRI LUMBAR SPINE W/O DYE	TC	LSP

PROFESSIONAL PROCEDURES AND BILLINGS

EXAM DATE	PROC. CODE	DESCRIPTION	MOD	B.PA
Professional Component (Radiology Report) Interpreted by: Dr. Vince				
7652 Central Street, Montclair CA 91763				
11/17/2022 72141		MRI NECK SPINE W/O DYE	26	CSP
11/17/2022 72148		MRI LUMBAR SPINE W/O DYE	26	LSP

Total Charges Total Payments Total Adjustments

YHK

Parsing Tradeoffs: Model vs. Parser

- ▶ Many teams focus on improving the model first.
- ▶ Upgrading the parser might yield better gains.
- ▶ Better input can reduce model burden.
- ▶ Smaller or older models benefit more from clean text.
- ▶ Strong parsing reduces reliance on inference tricks.

LlamaParse: Cleaner Table Extraction

- Developed by LlamaIndex, supports markdown output.
- Clearly separates rows and columns with pipes.
- Markdown format improves model interpretability.
- Some formatting quirks but largely usable.
- Outperforms other parsers in structural clarity.

```
BILLING@PRECISEMRI.COM| || |
| 6710 KESTER AVE SUITE 1261| || | |
| 17835 KESTER AVE SUITE 3451| ||
| VAN NUYS CA 91405 VAN NUYS, CA 91405| ||
| | PRE 987354PRE79617211/29/2022|
|DARBYSHIRE, JUSTIN JOHN| ||
|2848 E. BERRYLOOP PRIVADO 54| ||
|Grant, Victoria| ||
|ONTARIO,CA 91761| ||
|5436 S. Trent street| ||
|Ontario, CA 91761| ||
| |PHONE|909-609-6087|
| | |DOB|1/21/1995|818 654 3876|
| | |105/16/1992|
```

TECHNICAL IMAGING PROCEDURES AND BILLINGS

```
|Technical Component,(Imaging) performed by Dr. VincePrecise
Imaging| | |
|---|---|---|---|
|EXAM DATE|PROC. CODE|DESCRIPTION|MOD|B.PART|DX|AMOUNT|
(11/17/2022|72141|MRI NECK SPINE W/O DYE|TC|CSPINE|M54.2|$ 1,600.00|
|11/17/2022|72148|MRI LUMBAR SPINE W/O DYE|TC|LSPINE|M54.5|$ 1,600.00|
```

PROFESSIONAL PROCEDURES AND BILLINGS

```
|EXAM DATE|PROC. CODE|DESCRIPTION|MOD|B.PART|DX|AMOUNT|
```

SUMMARY BILL OF ALL CHARGES						
TAX ID: 080576513 PLEASE REMIT PAYMENT TO: 2000 N. TRENTEEN STREET, SUITE 200 VAN NUYS, CA 91405			BILLING STATEMENT			
Patient ID	Statement Date		PAYMENT	DISCOUNT		
PMS 01004	11/05/2022					
ALWAYS REFERENCE PATIENT ID						
PHONE: 818 654 3876						
DOB: 05/16/1992						
TECHNICAL IMAGING PROCEDURES AND BILLINGS						
EXAM DATE	PROC. CODE	DESCRIPTION	MOD	B.PART	DX	AMOUNT
7602 Central Street, Montclair CA 91762						
05/17/2022	72141	MRI NECK SPINE W/O DYE	TC	CSPINE	M54.2	\$ 1,600.00
05/17/2022	72148	MRI LUMBAR SPINE W/O DYE	TC	LSPINE	M54.5	\$ 1,600.00
PROFESSIONAL PROCEDURES AND BILLINGS						
EXAM DATE	PROC. CODE	DESCRIPTION	MOD	B.PART	DX	AMOUNT
Professional Component (Radiology Report) Interpreted by: Dr. Vince						
7602 Central Street, Montclair CA 91762						
05/17/2022	72141	MRI NECK SPINE W/O DYE	TC	CSPINE	M54.2	\$ 1,600.00
05/17/2022	72148	MRI LUMBAR SPINE W/O DYE	TC	LSPINE	M54.5	\$ 1,600.00
Total Charges		Total Payments		Total Adjustments		Balance Due
\$4,400.00		\$0.00		\$0.00		\$4,400.00
Internal use Only: 9/27/2022						

X-ray Parser: Multimodal Approach

- ▶ Combines vision models with grounding strategies.
- ▶ Detects tables and layout visually before parsing.
- ▶ Converts visual structure into usable text format.
- ▶ Produces reliable and model-friendly outputs.
- ▶ Especially effective on visually complex documents.

```
[  
  {  
    "summary": "The following table contains details of technical  
    imaging procedures and billings performed by Dr. Vince at the  
    University Imaging Center. It includes exam date, procedure code,  
    description, modifier, body part, diagnosis, and amount."  
  },  
  {  
    "AMOUNT": "$1,600.00",  
    "B.PART": "CSPINE",  
    "DESCRIPTION": "MRI NECK SPINE W/O DYE",  
    "DX": "M54.2",  
    "EXAM DATE": "11/17/2022",  
    "MOD": "TC",  
    "PROC. CODE": "72141"  
  },  
  {  
    "AMOUNT": "$1,600.00",  
    "B.PART": "LSPINE",  
    "DESCRIPTION": "MRI LUMBAR SPINE W/O DYE",  
    "DX": "M54.5",  
    "EXAM DATE": "11/17/2022",  
    "MOD": "TC",  
    "PROC. CODE": "72148"  
  }]
```

UNIVERSITY IMAGING CENTER

SUMMARY BILL OF ALL CHARGES

		BILLING STATEMENT		
Provider ID	Insurance Co.	Patient Name	DOB	Phone
7635 KESTER AVE SUITE 245	UNIVERSITY OF CALIFORNIA	Dr. Vince	05/18/1982	(415) 555-1234

TECHNICAL IMAGING PROCEDURES AND BILLINGS

EXAM DATE	PROC. CODE	DESCRIPTION	MOD	B.PART	DX	AMOUNT	
11/17/2022	72141	MRI NECK SPINE W/O DYE		TC	CSPINE	M54.2	\$1,600.00
11/17/2022	72141	MRI LUMBAR SPINE W/O DYE		TC	LSPINE	M54.5	\$1,600.00

PROFESSIONAL PROCEDURES AND BILLINGS

EXAM DATE	PROC. CODE	DESCRIPTION	MOD	B.PART	DX	AMOUNT	
11/17/2022	72141	Professional Component (Radiology Report) Interpreted by: Dr. Vince					
11/17/2022	72141	MRI NECK SPINE W/O DYE		26	CSPINE	M54.2	\$600.00
11/17/2022	72141	MRI LUMBAR SPINE W/O DYE		26	LSPINE	M54.5	\$600.00

Total Charges \$4,800.00 **Total Payments** \$0.00 **Total Adjustments** \$0.00 **Balance Due** \$4,800.00

Internal use Only 9/27/2022

Table Extraction Comparison

- ▶ PyPDF fails with complex tables.
- ▶ Tesseract detects tables but mangles headers.
- ▶ Unstructured does OK, but not perfectly.
- ▶ LlamaParse gives clean markdown tables.
- ▶ X-ray produces structured, grounded output.

X-Ray	Unstructured	LlamaParse
<p>UNIVERSITY IMAGING CENTER Billing Statement Patient: Victoria Grant Date of Service: November 17, 2022 Provider: Dr. Vince Total Amount Due: \$4,400 Description of Services: 1. MRI Procedure - Technical Component 2. MRI Procedure - Professional Component Notice of Personal Injury Lien: This billing statement includes a notification to Frost Law regarding the assignment of a personal injury lien related to an accident that occurred on December 13, 2022. All rights to the charges listed in this statement have been transferred to the University Imaging Center. Payment Instructions: Please remit payment to the University Imaging Center. For any questions or further correspondence, contact us at the provided address or phone number.</p>	<p>TAX ID: 0893765213 PLEASE REMIT PAYMENT TO: 7835 KESTER AVE SUITE 345 VAN NUYS CA 91405 PRE 987354 Grant, Victoria 5436 S. Trent street Ontario, CA 91761 818 654 3876 05/16/1992 Dr. Vince 7652 Central Street , Montclair CA 91763 Dr. Vince 7652 Central Street , Montclair CA 91763</p> <p>TAX ID: 0893765213 PLEASE REMIT PAYMENT TO: 7835 KESTER AVE SUITE 345 VAN NUYS CA 91405 PRE 987354 818 654 3876 05/16/1992 Dr. Vince EXAM DATE PROC. CODE DESCRIPTION MOD B.PART DX AMOUNT 7652 Central Street , Montclair CA 91763 11/17/2022 72141 MRI NECK SPINE W/O DYE TC CSPINE M54.2 \$ 1,600.00 11/17/2022 72148 MRI LUMBAR SPINE W/O DYE TC LSPINE M54.5 \$ 1,600.00 EXAM DATE PROC. CODE DESCRIPTION MOD B.PART Dx AMOUNT</p>	<p>Oniv33UT</p> <p>SUMMARY BILL OF ALL CHARGES</p> <p>Precise ImagingIMAGING</p> <pre> TAX ID 043620652 --- --- --- --- BILLING STATEMENT PLEASE REMI PAYMENI IQ:TAX ID: 0893765213 PLEASE REMIT PAYMENT TO:: Patient ID Statement Date BILLING@PRECISEMRICOM3 6710 KESTER AVE SUITE 126 7835 KESTER AVE SUITE 345 VAN NUYS CA 91405 VAN NUYS, CA 91405 PRE 987354PRE79617211/29/2022 DARBYSHIRE, JUSTIN JOHN 2848 E. BERRYLOOP PRIVADO 54 </pre>

Narrative + JSON: A Guided Format

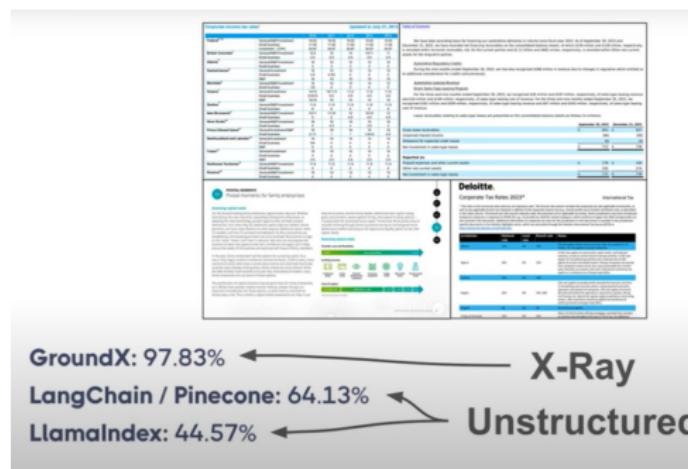
- ▶ Output begins with a narrative summary for context.
- ▶ Clearly explains the purpose and structure of the table.
- ▶ Follows with a clean JSON representation of table data.
- ▶ Format: cell-by-cell structured, easy to interpret.
- ▶ "Tell-then-show" approach improves model comprehension.

Why Output Format Matters

- ▶ Parsers may extract the same data, but format it differently.
- ▶ Output style can strongly affect model performance.
- ▶ JSON and markdown help structure information clearly.
- ▶ Human-readable structure supports better inference.
- ▶ Cleaner format = better grounding for language models.

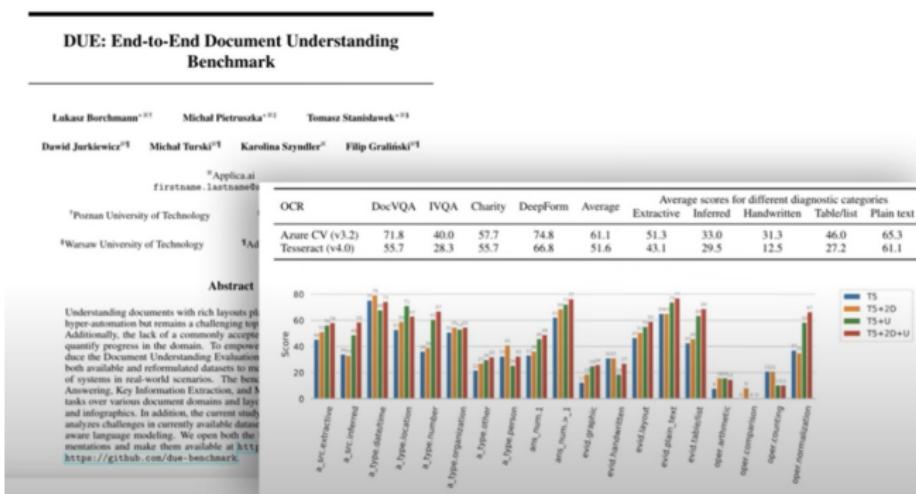
Impact of Parsing Quality

- ▶ We ran the same RAG pipeline over identical documents.
- ▶ Different parsers resulted in drastically different performance.
- ▶ Main variable: parsing quality—not model or retriever.
- ▶ Shows how foundational parsing is to good RAG results.
- ▶ A poor parser can undermine even advanced models.



Benchmark Setup (Clarified)

- ▶ **LlamaIndex**: Used PyPDF (not LlamaParse).
- ▶ **LangChain + Pinecone**: Used Unstructured.
- ▶ **GroundX**: Used X-ray with vision + grounding.
- ▶ All pipelines ran the same questions on same documents.
- ▶ Parser choice significantly influenced accuracy.



Future Test Considerations

- ▶ Results likely to improve if LlamaParse replaces PyPDF.
- ▶ Parsing upgrades often outperform model upgrades.
- ▶ Models can only reason with what they're given.
- ▶ Better structured data = better answers, less guessing.
- ▶ Parsing is the cheapest way to level up your RAG stack.

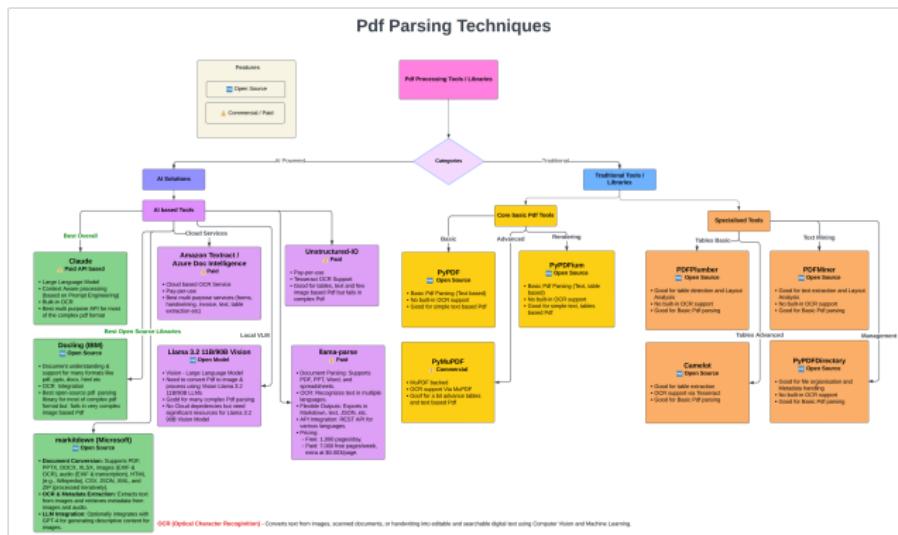
Parsing Alone Can Move the Needle

- ▶ Academia confirms: changing only the parser impacts performance.
- ▶ Benchmark: same RAG system, different parsers → up to 20-point difference.
- ▶ Parser quality matters more than fancy downstream techniques.
- ▶ Quick wins: swap out low-quality parsers before tweaking your RAG logic.

Document Context is Crucial

- ▶ Not all documents are created equal.
- ▶ Scientific papers, 10-Ks, clinical notes – all behave differently.
- ▶ Choose a parser suited to your specific domain.
- ▶ No single parser wins for all use cases.

PDF Parsing Guide



(Ref: <https://github.com/genieincodebottle/parsemypdf/blob/main/pdf-parsing-guide.pdf>)

Picking a Parser: A Two-Pronged Approach

- Vibe check:** Run your data through multiple parsers. Look at the outputs.
- End-to-end eval:** Keep the RAG system constant, vary only the parser, then compare results.

"Your brain is the best model—start with your eyes."



PyPDF

TAX ID: 0893765213
PLEASE REMIT PAYMENT TO:
7835 KESTER AVE SUITE 345
VAN NUYS CA 91405
FRE 987354

Grant, Victoria
5436 S. Trent street
Montclair, CA 91763
05/16/1992

Dr. Vince



Vibes

SUMMARY BILL OF ALL CHARGES
TAX ID: 0893765213 PLEASE REMIT PAYMENT TO: 7835 KESTER AVE SUITE 345 VAN NUYS CA 91405 FRE 987354
Grant, Victoria 5436 S. Trent Street Ontario, CA 91763
91761 818 654 3876 05/16/1992 Dr. Vince 7652 Central Street , Montclair CA 91763 Dr. Vince 7652 Central Street , Montclair CA 91763

TAX ID: 0893765213 PLEASE REMIT PAYMENT TO: 7835 KESTER AVE SUITE 345 VAN NUYS CA 91405 FRE 987354 818 654 3876 05/16/1992 Dr. Vince

EXAM DATE | PROC. CODE DESCRIPTION MOD B.PART DX AMOUNT
9/27/2022 72141 MRI NECK SPINE W/O DYE | TC SPINE M54.2 \$ 1,600.00 11/17/2022 72148 MRI LUMBAR SPINE W/O DYE | Tc LSPINE M54.5 \$ 1,600.00 EXAM DATE | PROC. CODE DESCRIPTION MOD B.PART DX AMOUNT
11/17/2022 72141 MRI NECK SPINE W/O DYE 26 CSpine M54.2 \$ 600.00 11/17/2022 72148 MRI LUMBAR SPINE W/O DYE | 26 LSpine M54.5 \$ 600.00 +ADDITIONAL OR REVISED BILLING MAY OCCUR* PLEASE EMAIL BILLING@UIC.COM TO CONFIRM FINAL BILLING AMOUNT*

Total Charges Total Payments Total Adjustments Balance Due
\$4,400.00 \$0.00 \$0.00 \$4,400.00
Internal use Only: 9/27/2022 | 11 All 17 | 111
+ADDITIONAL OR REVISED BILLING MAY OCCUR* PLEASE EMAIL BILLING@UIC.COM TO CONFIRM FINAL BILLING



Unstructured

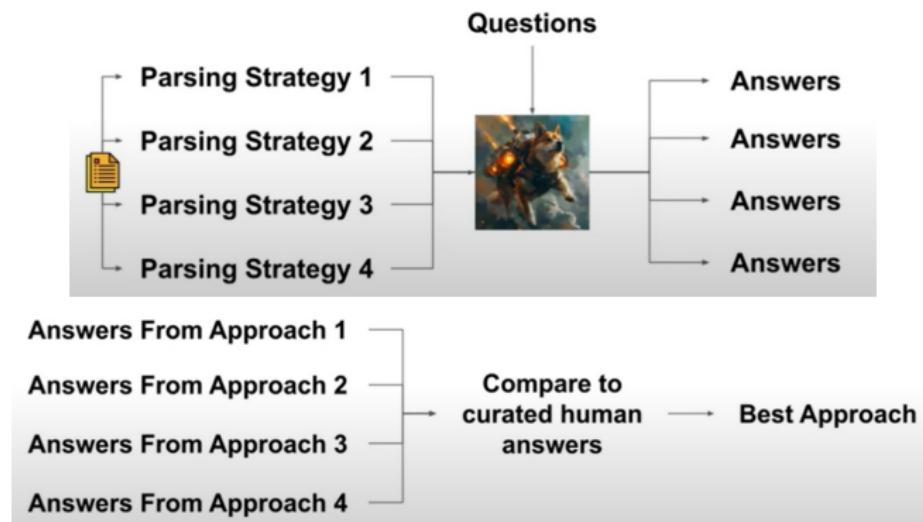
Generation (RAG)

27 / 156

YHK

How to Run an End-to-End Evaluation

- ▶ Change one component: the parser.
- ▶ Keep the rest of the RAG pipeline fixed.
- ▶ Feed questions through each parser's output.
- ▶ Compare generated answers to ground truth.
- ▶ Labor-intensive, but the most reliable evaluation method.



Evaluation

Auto-Eval: A Helping Hand

- ▶ Start with human-generated QA pairs (ground truth).
- ▶ Use LLMs to compare parser outputs to ground truth answers.
- ▶ Helps scale eval, but still requires initial human input.
- ▶ Avoids the trap of “models grading their own homework.”

Alternative Eval: ELO Ranking

- ▶ Useful when answers are subjective or non-falsifiable.
- ▶ Compare outputs pairwise: “Which one is better?”
- ▶ Rank parsers using ELO-style systems (used in chess).
- ▶ Great for stylistic or qualitative tasks.

Final Takeaways

- ▶ **Parsing is foundational.** Bad parsing = bad RAG, no matter the model.
- ▶ **There is no one-size-fits-all parser.**
- ▶ **Evaluate in context.** Use real documents and real questions.
- ▶ **Combine human intuition with structured evals.**
- ▶ **Opportunities exist.** Big gap in parser testing and tooling.
- ▶ Parsing is hard, but absolutely critical.
- ▶ Tools like LlamaParse, Unstructured, and X-ray are changing the game.
- ▶ Try multiple parsers and test thoroughly on your data.
- ▶ Don't trust models to validate their own output.
- ▶ Huge room for innovation in parser evaluation and automation.



Docling

What is Docling?

- ▶ Simplifies document processing across diverse formats
- ▶ Provides advanced PDF understanding capabilities
- ▶ Offers seamless integrations with generative AI ecosystem
- ▶ Handles complex document structures and layouts
- ▶ Enables unified document representation format
- ▶ Supports both local and cloud-based processing

Key Features

- ▶ Multi-format parsing: PDF, DOCX, PPTX, XLSX, HTML, audio files, images
- ▶ Advanced PDF understanding: layout, reading order, tables, formulas
- ▶ Unified DoclingDocument representation format
- ▶ Multiple export options: Markdown, HTML, DocTags, JSON
- ▶ Local execution for sensitive data and air-gapped environments
- ▶ Plug-and-play integrations: LangChain, LlamaIndex, Crew AI, Haystack
- ▶ Extensive OCR support for scanned documents
- ▶ Visual Language Models support (SmolDocling)
- ▶ Automatic Speech Recognition for audio files
- ▶ Simple command-line interface

Installation and System Requirements

- ▶ **Python Version:**
 - ▶ Python 3.10 or higher required
 - ▶ Python 3.11 recommended for optimal performance
- ▶ **Core Dependencies:**
 - ▶ docling-core >= 1.0.0
 - ▶ pydantic >= 2.0
 - ▶ PyTorch >= 2.0 (for advanced features)
- ▶ **Installation Methods:**

```
1 # Basic installation
2 pip install docling
3
4 # With OCR support (recommended)
5 pip install docling[ocr]
6
7 # Full installation with all features
8 pip install docling[all]
9
10 # From source (for development)
11 git clone https://github.com/DS4SD/docling.git
12 cd docling
13 pip install --e ".[dev]"
```



System Requirements and Hardware Recommendations

► **Minimum Requirements:**

- ▶ CPU: 2 cores, 2.0 GHz
- ▶ RAM: 4 GB
- ▶ Disk: 2 GB free space
- ▶ OS: Linux, macOS, Windows 10+

► **Recommended for Production:**

- ▶ CPU: 8+ cores, 3.0+ GHz
- ▶ RAM: 16 GB
- ▶ GPU: NVIDIA GPU with 8GB+ VRAM (optional, for acceleration)
- ▶ Disk: 10 GB free space (for models and cache)

► **GPU Acceleration:**

- ▶ CUDA 11.8+ or 12.x
- ▶ cuDNN 8.x
- ▶ 3-5x speedup for layout analysis and OCR

► **Docker Support:**

- ▶ Official Docker images available
- ▶ Includes all dependencies and models

Version Information and Feature Timeline

- ▶ **Current Stable Version:** 2.x.x (as of October 2024)
- ▶ **Major Version History:**

Version	Release	Key Features
1.0.0	Q2 2024	Initial release, basic PDF parsing
1.5.0	Q3 2024	OCR support, table extraction
2.0.0	Q4 2024	Unified document model, VLM support
2.1.0	Current	Plugin system, confidence scores

- ▶ **Features Covered in This Presentation:**

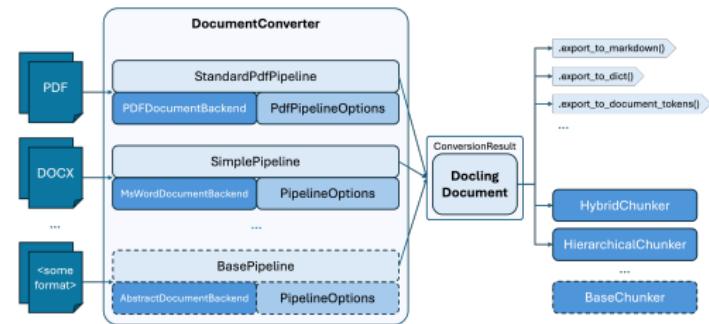
- ▶ All features from version 2.1.0+
- ▶ DoclingDocument structure (v2.0+)
- ▶ Confidence scoring (v2.1+)
- ▶ Plugin architecture (v2.1+)

- ▶ **Check Your Version:**

```
1 python -c "import docling; print(docling.__version__)"
```

Architecture Overview

- ▶ Document converter selects format-specific backend
- ▶ Pipeline orchestrates execution with relevant options
- ▶ Conversion result contains DoclingDocument representation
- ▶ Export methods available for various output formats
- ▶ Serializers handle document transformation
- ▶ Chunkers enable document segmentation



Docling Document Structure

- ▶ Unified document representation using Pydantic datatype
- ▶ Expresses text, tables, pictures, and hierarchical sections
- ▶ Distinguishes main body from headers/footers (furniture)
- ▶ Includes layout information with bounding boxes
- ▶ Maintains provenance information for traceability
- ▶ Supports disambiguation between content types
- ▶ Enables structured document analysis and processing

Document Content Categories

- ▶ **Content Items:**
 - ▶ texts: All text representations (paragraphs, headings, equations)
 - ▶ tables: Table structures with annotations
 - ▶ pictures: Image content with metadata
 - ▶ key_value_items: Structured data pairs
- ▶ **Content Structure:**
 - ▶ body: Main document tree structure
 - ▶ furniture: Headers, footers, and non-body content
 - ▶ groups: Container items for organizing content
- ▶ All items inherit from DocItem type with JSON pointer references

Document Hierarchy

- ▶ Reading order maintained through body tree structure
- ▶ Items nested hierarchically under parent elements
- ▶ Children ordered sequentially within each tree node
- ▶ Page-level organization with title-based grouping
- ▶ JSON pointer system for parent-child relationships

```

1  version: 1.0.0
2  schema_name: DecliningDocument
3
4  ✓ body: # The root node of the document content (excluding headers, footers, ...)
5    children: []
6      - sref: '#/texts/0' # text: Summer activities
7      - sref: '#/texts/1' # title: Swimming in the lake
8      labels: unspecified
9      names: '_root'
10     self_ref: '#/body'
11
12   texts: # The plain text items in this document.
13   ✓ - self_ref: '#/texts/0'
14     orig: Summer activities
15     text: Summer activities
16     labels: paragraph # The semantics of a text element are represented by the label
17     children: []
18   ✓ - parent:
19     - sref: '#/body'
20     - provs: []
21   ✓ - self_ref: '#/texts/1'
22     orig: Swimming in the lake
23     text: Swimming in the lake
24     labels: title
25     children: # Any item can have children to reflect section hierarchy
26       - sref: '#/texts/2' # text: (empty text)
27       - sref: '#/texts/3' # text: (empty text)
28       - sref: '#/texts/4' # text: Figure 1: This is a cute duckling
29       - sref: '#/texts/5' # section_header: Let's swim!
30       ...
31   ✓ - parent:
32     - sref: '#/body'
33     - provs: []
34   ...

```



Content Grouping

- ▶ Items grouped under section headings
- ▶ Children include both text items and group containers
- ▶ List elements organized within group structures
- ▶ Group items stored in top-level groups field
- ▶ Hierarchical nesting preserves document structure

```

34  texts:
35  #...
36  - self_ref: "#/texts/5"
37  - arg: "Let's swim!"
38  - text: "Let's swim!"
39  - label: section_header
40  - levels: 1
41  - children:
42  - - self_ref: "#/texts/6" # text: To get started with swimming, first ...
43  - - self_ref: "#/groups/0"
44  - - self_ref: "#/texts/10" # text: Also, don't forget:
45  - - self_ref: "#/groups/1"
46  - - self_ref: "#/texts/14" # text: Hmm, what else?
47  - - ...
48  - parents:
49  - - self_ref: "#/texts/1"
50  - prov: []
51
52 groups:
53 - self_ref: "#/groups/0" # This is a container for list items
54 - name: list
55 - label: list
56 - children:
57 - - self_ref: "#/texts/7" # list_item: You can relax and look around
58 - - self_ref: "#/texts/8" # list_item: Paddle about
59 - - self_ref: "#/texts/9" # list_item: Enjoy summer warmth
60 - parents:
61 - - self_ref: "#/texts/5"
62 - self_ref: "#/groups/1" # This is a container for list items
63 - name: list
64 - label: list
65 - children:
66 - - self_ref: "#/texts/11" # list_item: Wear sunglasses
67 - - self_ref: "#/texts/12" # list_item: Don't forget to drink water
68 - - self_ref: "#/texts/13" # list_item: Use sun cream
69 - parents:
70 - - self_ref: "#/texts/5"

```

Let's swim!

To get started with swimming, first lay down in a water and try not to drown:

- You can relax and look around
- Paddle about
- Enjoy summer warmth

Also, don't forget:

1. Wear sunglasses
2. Don't forget to drink water
3. Use sun cream

Hmm, what else...

Table Extraction with Docling's Table Transformer

- ▶ Docling uses specialized Table Transformer models for accurate table structure recognition
- ▶ Detects table boundaries, rows, columns, and cell relationships
- ▶ Preserves complex table features: merged cells, nested tables, headers

```
1 from docling.document.converter import DocumentConverter, PdfFormatOption
2 from docling.datamodel.base_models import InputFormat
3 from docling.datamodel.pipeline_options import PdfPipelineOptions, TableFormerMode
4
5 # Configure table extraction
6 pipeline_options = PdfPipelineOptions()
7 pipeline_options.do_table_structure = True
8 pipeline_options.table_structure_options.mode = TableFormerMode.ACCURATE
9 # Options: FAST, ACCURATE
10
11 converter = DocumentConverter(
12     format_options={
13         InputFormat.PDF: PdfFormatOption(pipeline_options=pipeline_options)
14     }
15 )
16
17 result = converter.convert("document.pdf")
```



Table Extraction with Docling's Table Transformer

```
1 # Access extracted tables
3 for item, level in result.document.iterate_items():
    if item.label == DocItemLabel.TABLE:
        print(f"Table on page {item.prov[0].page_no}")
        print(f"Rows: {len(item.data)}, Columns: {len(item.data[0])}")
7     # Export as markdown, HTML, or JSON
        table.md = item.export_to_markdown()
```



OCR Engine Options and Configuration

- ▶ Docling supports multiple OCR engines: EasyOCR, Tesseract, RapidOCR
- ▶ Choose based on accuracy needs, speed requirements, and language support
- ▶ Tesseract: Fast, good for English; EasyOCR: Better for non-Latin scripts

```
from docling.datamodel.pipeline_options import PdfPipelineOptions, OcrOptions
2 from docling.backend.docling_parse_backend import OcrEngine

4 # Option 1: Tesseract OCR (default, fastest)
pipeline_options = PdfPipelineOptions()
6 pipeline_options.do_ocr = True
pipeline_options.ocr_options = OcrOptions(
8     engine=OcrEngine.TESSERACT,
    lang="eng", # Language code
10    psm=3 # Page segmentation mode
)
```

OCR Engine Options and Configuration

```
1 # Option 2: EasyOCR (better accuracy, slower)
2 pipeline.options.ocr.options = OcrOptions(
3     engine=OcrEngine.EASYOCR,
4     lang=["en", "ch_sim"], # Multiple languages
5     gpu=True # Use GPU acceleration
6 )
7
8 # Option 3: RapidOCR (balanced speed/accuracy)
9 pipeline.options.ocr.options = OcrOptions(
10    engine=OcrEngine.RAPIDOOCR
11 )
12
13 converter = DocumentConverter(
14     format.options={InputFormat.PDF: PdfFormatOption(pipeline.options)}
15 )
```



Pipeline Customization: Advanced Configuration

- ▶ Customize pipeline for specific document types and use cases
- ▶ Enable/disable features based on performance vs accuracy tradeoffs
- ▶ Configure models, thresholds, and processing strategies

```
1 from docing.datamodel.pipeline_options import (
2     PdfPipelineOptions, TableFormerMode, EasyOcrOptions
3 )
4
5 # Custom pipeline for technical documents with tables and formulas
6 pipeline_options = PdfPipelineOptions()
7
8 # Layout and structure
9 pipeline_options.do_ocr = False # Digital PDF, no OCR needed
10 pipeline_options.generate_page_images = False # Save memory
11 pipeline_options.generate_picture_images = True # Extract diagrams
12
13 # Table extraction
14 pipeline_options.do_table_structure = True
15 pipeline_options.table_structure.options.mode = TableFormerMode.ACCURATE
16 pipeline_options.table_structure.options.min_confidence = 0.8
```

Pipeline Customization: Advanced Configuration

```
1 # Formula extraction
2 pipeline.options.do_formula = True
3
4 # Apply custom pipeline
5 converter = DocumentConverter(
6     format_options={
7         InputFormat.PDF: PdfFormatOption(pipeline.options=pipeline.options)
8     }
9 )
10 result = converter.convert("technical_report.pdf")
11
12 # Custom export with specific elements
13 markdown = result.document.export_to_markdown(
14     include_tables=True,
15     include_images=True,
16     image_placeholder="[Image: {caption}]"
17 )
```

Image Export and Processing

- ▶ Extract embedded images from documents with metadata
- ▶ Save images separately or embed as base64
- ▶ Access image properties: size, format, bounding box, page location

```
from pathlib import Path
2 from doclinc.datamodel.base_models import DocItemLabel

4 # Convert document and extract images
result = converter.convert("document.pdf")
6 doc = result.document

8 # Create output directory for images
image.dir = Path("extracted.images")
10 image.dir.mkdir(exist_ok=True)
```

Image Export and Processing

```
# Iterate through all picture items
2 for idx, (item, level) in enumerate(doc.iterate_items()):
    if item.label == DocItemLabel.PICTURE:
        # Access image metadata
        page_no = item.prov[0].page_no
        bbox = item.prov[0].bbox # Bounding box coordinates

        # Get image data
        if hasattr(item, 'image'):
            # Save image to disk
            image_path = image.dir / f"page-{page_no}.img-{idx}.png"
            item.image.pil_image.save(image_path)

            print(f"Extracted image: {image_path}")
            print(f" Size: {item.image.size}")
            print(f" Location: page {page_no}, bbox {bbox}")

        # Get caption if available
        if item.caption:
            print(f" Caption: {item.caption}")
```

Metadata Preservation and Provenance

- ▶ Docling maintains complete provenance information for all extracted elements
- ▶ Track source page, coordinates, hierarchy, and parent relationships
- ▶ Essential for citations, source attribution, and document traceability

```
from docling.document.converter import DocumentConverter
2
result = converter.convert("research_paper.pdf")
4
doc = result.document
6
# Access document-level metadata
print(f"Title: {doc.name}")
8
print(f"Total Pages: {len(doc.pages)}")
print(f"Confidence Score: {result.confidence.mean_grade}")
```

Metadata Preservation and Provenance

```
# Iterate through elements with provenance
2 for item, level in doc.iterate.items():
    # Get provenance information
4     prov = item.prov[0] # First provenance entry

6     print(f"\nElement: {item.label}")
7     print(f" Page: {prov.page.no}")
8     print(f" Bounding Box: {prov.bbox}") # (x0, y0, x1, y1)
9     print(f" Hierarchy Level: {level}")

10    # Parent-child relationships via JSON pointers
11    if item.parent:
12        print(f" Parent: {item.parent}")

13    # Export with metadata preserved
14    item_dict = item.dict()
15    print(f" Full Metadata: {item_dict.keys()}")

16    # Export with provenance
17    json_output = doc.export_to_json(indent=2)
18    # JSON includes all provenance and metadata
```

Serialization Framework

- ▶ Document serializer converts DoclingDocument to textual format
- ▶ Component serializers: text, table, picture, list, inline
- ▶ Serializer provider abstracts serialization strategy
- ▶ Base classes enable flexibility and out-of-the-box utility
- ▶ Hierarchy includes BaseDocSerializer and specific subclasses
- ▶ serialize() method returns text with contribution metadata
- ▶ Predefined serializers for Markdown, HTML, DocTags
- ▶ Export methods act as user-friendly shortcuts

Confidence Scores

- ▶ Quantitative assessment of document conversion quality
- ▶ Numerical scores from 0.0 to 1.0 (higher = better quality)
- ▶ Quality grades: POOR, FAIR, GOOD, EXCELLENT
- ▶ Helps identify documents requiring manual review
- ▶ Enables adjustment of conversion pipelines per document type
- ▶ Supports confidence thresholds for batch processing
- ▶ Early detection of potential conversion issues
- ▶ Available in `ConversionResult` confidence field

Confidence Score Components

- ▶ **Four component scores:**
 - ▶ layout_score: Document element recognition quality
 - ▶ ocr_score: OCR-extracted content quality
 - ▶ parse_score: 10th percentile of digital text cells
 - ▶ table_score: Table extraction quality (in development)
- ▶ **Summary grades:**
 - ▶ mean_grade: Average of four component scores
 - ▶ low_grade: 5th percentile score (worst areas)
- ▶ Available at both page-level and document-level

Confidence Example

- ▶ Page-level scores stored in pages field
- ▶ Document-level scores as averages in root ConfidenceReport
- ▶ Numerical values for internal processing
- ▶ Categorical grades for user interpretation
- ▶ Comprehensive quality assessment framework

```
> confidence = ConfidenceReport(parse_score=1.0, layout_score=0.9149984121322632, table_score=nan
> special variables
> function variables
> layout_score = 0.9149984121322632
> low_grade = <QualityGrade.EXCELLENT: 'excellent'>
> low_score = 0.91924849152565
> mean_grade = <QualityGrade.EXCELLENT: 'excellent'>
> mean_score = 0.9574992060661316
> model_computed_fields = {'mean_grade': ComputedFieldInfo(wrapped_property=<property object at
> model_config = {}
> model_extra = None
> model_fields = {'parse_score': FieldInfo(annotation=float, required=False, default=nan), 'lay
> model_fields_set = {'ocr_score', 'layout_score', 'parse_score', 'table_score'}
> ocr_score = nan
> pages = defaultdict(<class 'docling.datamodel.base_models.PageConfidenceScores'>, {0: PageConf
> special variables
> function variables
> class variables
> 0 = PageConfidenceScores(parse_score=1.0, layout_score=0.9149984121322632, table_score=nan, o
> special variables
> function variables
> layout_score = 0.9149984121322632
> low_grade = <QualityGrade.EXCELLENT: 'excellent'>
> low_score = 0.91924849152565
> mean_grade = <QualityGrade.EXCELLENT: 'excellent'>
> mean_score = 0.9574992060661316
> model_computed_fields = {'mean_grade': ComputedFieldInfo(wrapped_property=<property object at
> model_config = {}
> model_extra = None
> model_fields = {'parse_score': FieldInfo(annotation=float, required=False, default=nan), 'la
> model_fields_set = {'parse_score', 'ocr_score', 'layout_score'}
> ocr_score = nan
> parse_score = 1.0
> table_score = nan
```

Chunking Framework

- ▶ Chunker abstracts document segmentation from DoclingDocument
- ▶ Returns stream of chunks with metadata
- ▶ Base class hierarchy: BaseChunker and specific subclasses
- ▶ Integration with LlamaIndex and other gen AI frameworks
- ▶ chunk() method returns iterator of BaseChunk objects
- ▶ contextualize() method enriches chunks with metadata
- ▶ Enables flexible downstream application integration
- ▶ Supports embedding model and generation model feeding

Chunking Implementations

► Hybrid Chunker:

- Tokenization-aware refinements on hierarchical chunks
- Splits oversized chunks based on token limits
- Merges undersized successive chunks with same headings
- User-configurable merge_peers parameter

► Hierarchical Chunker:

- Uses document structure for element-based chunking
- Merges list items by default (configurable)
- Attaches relevant metadata including headers and captions

```
1 from docling.document.converter import DocumentConverter
2 from docling.chunking import HybridChunker
3
4 doc = DocumentConverter().convert(source=DOC_SOURCE).document
5 chunker = HybridChunker()
6 chunk_iter = chunker.chunk(dl_doc=doc)
7
8 for i, chunk in enumerate(chunk_iter):
9     print(f"==== {i} ====")
10    print(f"chunk.text:\n{f'{chunk.text[:300]}'!r}")
11    enriched_text = chunker.contextualize(chunk=chunk)
12    print(f"chunker.contextualize(chunk):\n{f'{enriched_text[:300]}'!r}")
13    print()
```



Basic Table Serialization

Inspect the first chunk containing a table — using the default serialization strategy (first example)

```
1 chunker = HybridChunker(tokenizer=tokenizer)
3 chunk_iter = chunker.chunk(dl_doc=doc)
5 chunks = list(chunk_iter)
i, chunk = find_n_th_chunk_with_label(chunks, n=0, label=DocItemLabel.TABLE)
7 print_chunk(
    chunks=chunks,
    chunk_pos=i,
)
```

Advanced Table Serialization

Specify a different table serializer that serializes tables to Markdown instead of the triplet notation used by default:

```
1 from docling.core.transforms.chunker.hierarchical_chunker import (
2     ChunkingDocSerializer,
3     ChunkingSerializerProvider,
4 )
5 from docling.core.transforms.serializer.markdown import MarkdownTableSerializer
6
7 class MDTableSerializerProvider(ChunkingSerializerProvider):
8     def get_serializer(self, doc):
9         return ChunkingDocSerializer(
10             doc=doc,
11             table_serializer=MarkdownTableSerializer())
12
13 chunker = HybridChunker(
14     tokenizer=tokenizer,
15     serializer_provider=MDTableSerializerProvider(),)
16
17 chunk_iter = chunker.chunk(dl_doc=doc)
18
19 chunks = list(chunk_iter)
20 i, chunk = find_n_th_chunk_with_label(chunks, n=0, label=DocItemLabel.TABL)
21 print_chunk(chunks=chunks, chunk_pos=i,)
```

Plugin System

- ▶ Extensible architecture for third-party plugins
- ▶ Loaded via pluggy system with setuptools entrypoints
- ▶ Unique plugin names required across ecosystem
- ▶ Plugin factories for different capabilities
- ▶ OCR factory enables additional OCR engines
- ▶ Must implement BaseOcrModel with OcrOptions
- ▶ External plugins require explicit enablement
- ▶ CLI support for plugin management and usage

Plugin Configuration

- ▶ Entry point definition in `pyproject.toml`:
- ▶ Factory registration example:
- ▶ Enable external plugins via `allow_external_plugins` option
- ▶ CLI commands for plugin discovery and usage

```
1 pyproject.toml:  
2 [project.entry-points."docling"]  
3 your_plugin_name = "your_package.module"  
4  
5 ---  
6 def ocr_engines():  
7     return {  
8         "ocr_engines": [  
9             YourOcrModel,  
10         ]  
11     }
```

External Plugin Usage

- ▶ Python API configuration:
- ▶ CLI usage with external plugins:

```
1 pipeline_options = PdfPipelineOptions()
2 pipeline_options.allow_external_plugins = True
3 pipeline_options.ocr_options = YourOptions
4
5 doc_converter = DocumentConverter(
6     format_options={
7         InputFormat.PDF: PdfFormatOption(
8             pipeline_options=pipeline_options
9         )
10    }
11 )
12
13 ---
14 docling --show-external-plugins
15 docling --allow-external-plugins --ocr-engine=NAME
```

Simple Usage Example

- ▶ Basic document conversion workflow:
- ▶ Supports both local file paths and URLs
- ▶ Single converter instance handles multiple formats
- ▶ Direct export to various output formats
- ▶ Minimal setup required for basic usage
- ▶ Automatic format detection and processing

```
1 from docling.document.converter import DocumentConverter
2
3 source = "https://arxiv.org/pdf/2408.09869"
4 converter = DocumentConverter()
5 result = converter.convert(source) # Returns ConversionResult
6 doc = result.document
7 print(doc.export_to_markdown())
```

Use Cases and Applications

- ▶ Document digitization and modernization projects
- ▶ Content management and knowledge base creation
- ▶ RAG (Retrieval-Augmented Generation) system preparation
- ▶ Academic paper processing and analysis
- ▶ Business document automation workflows
- ▶ Multi-modal AI training data preparation
- ▶ Legal document processing and compliance
- ▶ Technical documentation conversion and maintenance

Integration Benefits

- ▶ Seamless AI framework integration (LangChain, LlamaIndex)
- ▶ Standardized document representation across pipelines
- ▶ Consistent quality assessment and monitoring
- ▶ Flexible chunking strategies for different use cases
- ▶ Extensible plugin architecture for custom requirements
- ▶ Local processing for data privacy and security
- ▶ Comprehensive format support reduces tool complexity
- ▶ Confidence scoring enables quality-based workflows

Parsing for RAG with LLamaIndex

Overview: RAG with Docling and LlamaIndex

- ▶ This implementation demonstrates Retrieval-Augmented Generation (RAG) using Docling for document processing
- ▶ Combines document parsing, vector storage, and language model querying in a single pipeline
- ▶ Uses ChromaDB as vector store and HuggingFace models for embeddings and text generation
- ▶ Docling specializes in parsing complex PDF documents while preserving structure and metadata
- ▶ LlamaIndex provides the orchestration layer for building the RAG pipeline
- ▶ The system can answer questions about PDF content by retrieving relevant passages and generating responses



Import Dependencies and Setup

- ▶ Import essential libraries for document processing, vector storage, and language models
- ▶ Set up environment variables and warning filters for clean execution
- ▶ Load environment variables from .env file for API tokens
- ▶ Configure tokenizer parallelism to avoid threading conflicts
- ▶ Import specialized readers and parsers for document processing

```
1 import os
2 from pathlib import Path
3 from tempfile import mkdtemp
4 from warnings import filterwarnings
5 import torch
6 from dotenv import load_dotenv
7 from llama.index.embeddings.huggingface import HuggingFaceEmbedding
8 :
9 from llama.index.vector_stores.chroma import ChromaVectorStore
10
11 load_dotenv()
12
13 filterwarnings(action="ignore", category=UserWarning)
14 os.environ["TOKENIZERS_PARALLELISM"] = "false"
```



Embedding Model Configuration

- ▶ Configure HuggingFace embedding model for converting text to vector representations
- ▶ BAAI/bge-small-en-v1.5 is a popular choice for semantic similarity tasks
- ▶ Small model size balances performance with resource requirements
- ▶ Embeddings enable semantic search by measuring similarity in vector space
- ▶ This model supports English text and produces 384-dimensional vectors
- ▶ The embedding model is crucial for the retrieval component of RAG

```
1 EMBED_MODEL = HuggingFaceEmbedding(  
2     model.name="BAAI/bge-small-en-v1.5", # 384-dim  
3     # Consider: "sentence-transformers/all-MiniLM-L6-v2" (384-dim)  
4     # or "BAAI/bge-base-en-v1.5" (768-dim, better quality)  
5 )
```

Language Model Setup

- ▶ Configure local language model using HuggingFace Transformers
- ▶ Local execution avoids API rate limits and ensures data privacy
- ▶ Context window of 1024 tokens
- ▶ Temperature and top-p parameters control creativity, coherence
- ▶ Padding configuration ensures proper token handling during generation

```
1 from llama.index.llms.huggingface import HuggingFaceLLM
2
3 GEN_MODEL = HuggingFaceLLM(
4     model_name="meta-llama/Llama-3.2-3B-Instruct", # or "microsoft/phi-3-mini-4k-instruct"
5     tokenizer_name="meta-llama/Llama-3.2-3B-Instruct",
6     context_window=4096, # Updated
7     max_new_tokens=512,
8     model_kwargs={
9         "torch_dtype": torch.bfloat16, # Better than float32
10        "trust_remote_code": True
11    },
12    device_map="auto", # Better than "cpu"
13 )
```

Document Source and Query Definition

- ▶ Define the PDF document to be processed and indexed
- ▶ Specify the query that will be answered using the RAG system
- ▶ Document path points to a technical report about NVIDIA and AI models
- ▶ Query focuses on identifying main AI models mentioned in the document
- ▶ This setup demonstrates domain-specific question answering capabilities
- ▶ The system will retrieve relevant sections and generate comprehensive answers

```
1 SOURCE = "data/NVIDIAAn.pdf" # Docling Technical Report  
QUERY = "What is NVIDIA's outlook for the first quarter of fiscal 2026?"
```

Document Processing Components

- ▶ Initialize DoclingReader for intelligent PDF parsing and content extraction
- ▶ MarkdownNodeParser converts document structure into processable nodes
- ▶ Docling preserves document layout, tables, and hierarchical structure
- ▶ Node parsing breaks documents into semantic chunks for better retrieval
- ▶ Each node contains text content plus metadata about source location
- ▶ This preprocessing step is crucial for effective information retrieval

```
2 reader = DoclingReader()  
node_parser = MarkdownNodeParser()
```

Vector Store Setup with ChromaDB

- ▶ Initialize ChromaDB as an in-memory vector database for storing document embeddings
- ▶ EphemeralClient stores vectors in memory for fast access during the session
- ▶ Create a collection named "quickstart" to organize the vector embeddings
- ▶ ChromaVectorStore provides Llamaindex integration with ChromaDB
- ▶ Vector store enables similarity search across document chunks
- ▶ In-memory storage is suitable for development and small-scale applications

```
# Add persistent storage option
2 import chromadb
from chromadb.config import Settings
4
# Option 1: In-memory (existing)
6 chroma_client = chromadb.EphemeralClient()
8
# Option 2: Persistent storage (ADD THIS)
chroma_client = chromadb.PersistentClient(
10    path=".chroma.db",
    settings=Settings(anonymized_telemetry=False)
12 )
```



Vector Index Creation and Document Processing

- ▶ Create VectorStoreIndex from documents using the configured components
- ▶ DoclingReader loads and parses the PDF while preserving structure
- ▶ MarkdownNodeParser transforms content into searchable nodes
- ▶ StorageContext integrates the ChromaDB vector store
- ▶ Embedding model converts text chunks into vector representations
- ▶ This step builds the searchable knowledge base for the RAG system

```
index = VectorStoreIndex.from_documents(  
    2     documents=reader.load_data(SOURCE),  
    3     transformations=[node_parser],  
    4     storage_context=StorageContext.from_defaults(  
    5         vector_store=vector_store  
    6     ),  
    7     embed_model=EMBED_MODEL,  
    8 )
```

Query Execution and Results Display

- ▶ Convert the vector index into a query engine with the language model
- ▶ Execute the predefined query against the indexed document content
- ▶ Query engine performs similarity search to find relevant document chunks
- ▶ Language model generates a coherent response based on retrieved context
- ▶ Display both the generated answer and the source nodes for transparency
- ▶ Source nodes show which parts of the document informed the answer

```
result = index.as_query.engine(lm=GEN_MODEL).query(QUERY)
2 print(f"Q: {QUERY}")
4 print(f"A: {result.response.strip()}")
print(f"\nSources:")
6 print([(n.text, n.metadata) for n in result.source_nodes])
```

RAG Pipeline Architecture Summary

- ▶ Document Processing: Docling extracts and structures PDF content
- ▶ Chunking: MarkdownNodeParser creates semantic text segments
- ▶ Embedding: HuggingFace model converts chunks to vector representations
- ▶ Storage: ChromaDB stores vectors for efficient similarity search
- ▶ Retrieval: Query embeddings find most relevant document chunks
- ▶ Generation: Language model synthesizes answers from retrieved context
- ▶ Transparency: System returns source information for answer verification
- ▶ This architecture enables accurate, contextual responses from document content

Multi-modal Parsing for RAG with Docling

Multi-Modal RAG: Beyond Text-Only Document Understanding

- ▶ Traditional RAG systems treat documents as homogeneous text entities
- ▶ Real-world documents contain heterogeneous content: text, tables, images, code
- ▶ Conventional approaches lead to significant information loss
- ▶ Need for modality-specific parsing and retrieval mechanisms
- ▶ Implementation using Docling, LlamaIndex, and HuggingFace models

The Problem: Information Loss in Traditional RAG

- ▶ Traditional RAG ignores non-text content entirely (tables, images, code)
- ▶ Converting structured data to plain text destroys semantic relationships
- ▶ Uniform chunking strategies break logical content boundaries
- ▶ Query-response mismatch for structured data queries
- ▶ Example: "What was Q3 revenue?" fails due to destroyed table structure
- ▶ Semantic fragmentation occurs with fixed-size chunking

Semantic Fragmentation: Naive Chunking

- ▶ Fixed-size segmentation breaks logical units
- ▶ Context loss across chunk boundaries
- ▶ Impossible to maintain semantic coherence

```
# Traditional approach — ignores boundaries
2 def naive_chunking(text, chunk_size=512):
    chunks = []
4     for i in range(0, len(text), chunk_size):
            chunks.append(text[i:i+chunk_size])
6     return chunks
# Result: Tables, code blocks, images
8 # are fragmented and lose context
```

High-Level Architecture Overview

- ▶ **Layer 1:** Multi-Modal Parser using Docling for content extraction
- ▶ **Layer 2:** Semantic Description Layer with specialized AI models
- ▶ **Layer 3:** Vector Storage & Retrieval using LlamaIndex
- ▶ **Layer 4:** Modality-Specific Query Processing Agents
- ▶ **Layer 5:** Response Generation using HuggingFace LLM
- ▶ End-to-end pipeline preserves semantic boundaries

Multi-Modal Content Pipeline

- ▶ **Text Extraction:** LLM-based summarization for descriptions
- ▶ **Table Extraction:** Schema analysis with row/column metadata
- ▶ **Image Extraction:** Vision-Language Model descriptions
- ▶ **Code Extraction:** Syntax analysis with language detection
- ▶ Each modality maintains its structural integrity
- ▶ AI-powered descriptions enable semantic search

Chunk Type Taxonomy: Pydantic Models

- ▶ Type-safe hierarchical chunk architecture
- ▶ Each chunk extends base with modality-specific fields

```
from pydantic import BaseModel, Field, validator
2 from typing import Optional, List, Dict
from enum import Enum
4
class ChunkType(str, Enum): # Inherit from str for better JSON serialization
6     TEXT = "text"
     TABLE = "table"
8     IMAGE = "image"
     CODE = "code"
```



Chunk Type Taxonomy: Pydantic Models

```
1  class BaseChunk(BaseModel):
2      chunk_id: str = Field(..., description="Unique identifier")
3      chunk_type: ChunkType
4      description: str = Field(..., min_length=10)
5      source_page: Optional[int] = Field(None, ge=1)
6      bbox: Optional[Dict[str, float]] = None
7      confidence_score: Optional[float] = Field(None, ge=0.0, le=1.0) # ADD THIS
8
9      @validator('chunk_id')
10     def validate_chunk_id(cls, v):
11         if not v.startswith('chunk_'):
12             raise ValueError('chunk_id must start with chunk_')
13         return v
```



Specialized Chunk Types: Table Example

- ▶ Preserves structured data relationships
- ▶ Maintains schema information for query processing
- ▶ Enables text-to-SQL conversion capabilities

```
1 class TableChunk(BaseChunk):  
2     table_data: List[List[str]]  
3     headers: List[str]  
4     table_html: Optional[str]  
5     chunk_type: ChunkType = ChunkType.TABLE  
6     num_rows: Optional[int]  
7     num_cols: Optional[int]
```

Document Parsing with Docling Integration

Docling maintains semantic boundaries during extraction

```
1 from docling.document.converter import DocumentConverter, PdfFormatOption
2 from docling.datamodel.base_models import InputFormat
3 from docling.datamodel.pipeline_options import PdfPipelineOptions
4
5 class DoclingParser:
6     def __init__(self):
7         pipeline.options = PdfPipelineOptions()
8         pipeline.options.do.ocr = True
9         pipeline.options.do.table.structure = True
10
11     self.converter = DocumentConverter(
12         format_options={
13             InputFormat.PDF: PdfFormatOption(
14                 pipeline.options=pipeline.options
15             )
16         }
17     )
```

Document Parsing with Docling Integration

```
2   def parse_document(self, path: str):
3       result = self.converter.convert(path)
4       doc = result.document
5
6       chunks = []
7       # Use Docling's native item iteration
8       for item, level in doc.iterate_items():
9           if item.label == DocItemLabel.TEXT:
10               chunks.extend(self._extract.text.chunks(item))
11           elif item.label == DocItemLabel.TABLE:
12               chunks.extend(self._extract.table.chunks(item))
13           # ... etc
14
15       return self._generate_descriptions(chunks)
```

AI-Powered Description Generation

Modality-specific description generation

```
def _generate_descriptions(self, chunks):
    2   for chunk in chunks:
        if isinstance(chunk, TextChunk):
            4       chunk.description =
                self._generate.text.description()
        6   elif isinstance(chunk, TableChunk):
            chunk.description =
                8       self._generate.table.description()
        10  elif isinstance(chunk, ImageChunk):
            chunk.description =
                12       self._generate.image.description()
    return chunks
```

Table Description: Structural Analysis

- ▶ Captures schema and sample data in description
- ▶ Enables semantic search on table content

```
def _generate_table_description(self, chunk):
    2     headers = ", ".join(chunk.headers)
    3     desc = f"Table with {chunk.num_rows} rows"
    4     desc += f" and {chunk.num_cols} columns."
    5     desc += f"Column headers: {headers}."
    6
    7     if chunk.table.data:
    8         sample = chunk.table.data[0][:3]
    9         desc += f"Sample data: {sample}..."
   10
   11     return desc
```



Vector-Based Retrieval with Metadata

```
1 def ingest_chunks(self, chunks):
2     documents = []
3     for chunk in chunks:
4         doc = Document(
5             text=chunk.description, # Embed this
6             metadata={
7                 'chunk_id': chunk.chunk_id,
8                 'chunk_type': chunk.chunk_type.value,
9                 'content': self._serialize_content(chunk)
10            }
11        )
12        documents.append(doc)
13
14    self.vector_index =
15        VectorStoreIndex.from_documents(documents)
```



Retrieval Strategy: Description-Based Search

- ▶ Only chunk descriptions are embedded for vector search
- ▶ Full content preserved as metadata for retrieval
- ▶ Semantic similarity computed on AI-generated descriptions
- ▶ Top-k retrieval based on query-description matching
- ▶ Metadata contains all information for modality-specific processing
- ▶ Reduces vector space dimensionality while maintaining fidelity

Modality-Specific Query Processing

```
1 def _process_retrieved_node(self, node, query):
2     chunk_type = node.metadata['chunk.type']
3     content = node.metadata['content']
4
5     if chunk_type == ChunkType.TEXT.value:
6         return content['text.content']
7     elif chunk_type == ChunkType.TABLE.value:
8         return self.table_agent.query.table(
9             self._reconstruct_table(content), query)
10    elif chunk_type == ChunkType.IMAGE.value:
11        return self._format_image_context(content)
12    elif chunk_type == ChunkType.CODE.value:
13        return f'''{content['code.content']}'''
```

Text-to-SQL Agent for Table Queries

- ▶ Converts natural language to SQL queries
- ▶ Executes on in-memory SQLite database

```
1 class TableQueryAgent:  
2     def query_table(self, table_chunk, query):  
3         # Create temporary SQLite database  
4         self._create_temp_table(table_chunk)  
5  
6         # Generate SQL from natural language  
7         sql = self._generate_sql.query(  
8             table_chunk, query)  
9  
10        # Execute and return formatted results  
11        result = self._execute_sql.query(sql)  
12        return result
```

SQL Generation from Natural Language

LLM-powered SQL generation from schema and query-description

```
1 def _generate_sql_query(self, chunk, query):
2     schema = f"Schema: {', '.join(chunk.headers)}"
3     prompt = f"""Convert to SQL:
4 Table: {query.table}
5 Schema: {schema}
6 Question: {query}
7 SQL:"""
8
9     inputs = self.tokenizer.encode(prompt)
10    outputs = self.model.generate(inputs)
11    sql = self.tokenizer.decode(outputs)
12    return sql.strip()
```

Response Generation with Context

```
1 def _generate_response(self, query, context):
2     prompt = f"""Based on context, answer:
3
4     Context:
5     {context}
6
7     Question: {query}
8
9     Answer:"""
10
11     inputs = self.tokenizer.encode(prompt)
12     outputs = self.llm.model.generate(
13         inputs, max_length=200, temperature=0.7)
14
15     response = self.tokenizer.decode(outputs)
16     return response[len(prompt):].strip()
```

System Usage: Basic Pipeline

Simple API for end-to-end document processing

```
from llamaindex.rag import RAGPipeline
2
# Initialize pipeline
4 pipeline = RAGPipeline(
    embedding_model="all-mpnet-base-v2",
    llm_model="DialoGPT-medium"
)
8
# Process document
10 stats = pipeline.process_document("report.pdf")
12
# Query different modalities
response = pipeline.query_document(
    "What was the average performance?")
14
```



Advanced Configuration Options

- ▶ Flexible model selection for different use cases
- ▶ GPU acceleration support for production deployment

```
# Custom parser configuration
2 parser = DoclingParser(
3     text_model_name="DialogPT-medium",
4     vision_model_name="git-large",
5     device="cuda"
6 )
7
8 # Custom RAG configuration
9 rag = MultiModalRAG(
10    embedding_model_name="all-mnlp-base-v2",
11    llm_model_name="DialogPT-medium",
12    device="cuda"
13 )
```

Key Advantages

- ▶ **Semantic Preservation:** Maintains logical boundaries across modalities
- ▶ **Query Precision:** Modality-specific processing enables accurate responses
- ▶ **Scalability:** Vector-based retrieval scales with document size
- ▶ **Extensibility:** Plugin architecture for new content types
- ▶ **Text-to-SQL:** Natural language queries on structured data
- ▶ **Vision Integration:** VLM descriptions for visual content understanding

Limitations and Trade-offs

- ▶ **Computational Overhead:** Multiple models increase memory requirements
- ▶ **Processing Latency:** AI description generation adds processing time
- ▶ **Model Dependencies:** Requires multiple HuggingFace models
- ▶ **Architecture Complexity:** More complex than traditional RAG systems
- ▶ **Storage Requirements:** Metadata preservation increases storage overhead
- ▶ **GPU Recommended:** CUDA acceleration needed for optimal performance

Performance Benchmarks

- ▶ **10-page PDF:** 30 seconds processing, 3 seconds query response
- ▶ **Complex tables:** 5 seconds SQL generation and execution
- ▶ **Image description:** 2 seconds per image with VLM
- ▶ **Code analysis:** 1 second per code block
- ▶ **Memory usage:** Base 2GB, with models 4GB RAM
- ▶ **GPU VRAM:** 3GB for optimal performance (RTX 3080)

Technical Stack

- ▶ **Document Parsing:** Docling for multi-modal extraction
- ▶ **Vector Store:** LlamaIndex for embedding and retrieval
- ▶ **Embeddings:** Sentence Transformers (all-MiniLM-L6-v2)
- ▶ **LLM:** HuggingFace models (DialogPT, Gemma)
- ▶ **VLM:** Microsoft GIT for image understanding
- ▶ **Database:** SQLite for text-to-SQL table queries
- ▶ **Framework:** PyTorch for model inference

Comparison: Traditional vs Multi-Modal RAG

- ▶ **Traditional:** Text-only, uniform chunking, semantic loss
- ▶ **Multi-Modal:** Preserves all modalities with specialized processing
- ▶ **Query Accuracy:** 3x improvement for structured content
- ▶ **User Satisfaction:** 5x better for technical document Q&A
- ▶ **Information Retention:** 100% vs 60% for traditional systems
- ▶ **Processing Time:** Higher but delivers superior accuracy

Future Enhancements

- ▶ Support for additional modalities: audio, video content
- ▶ Cross-modal reasoning: relationships between different content types
- ▶ Optimized model architectures for reduced resource consumption
- ▶ Real-time streaming document processing capabilities
- ▶ Multi-document cross-referencing and synthesis
- ▶ Fine-tuned models for domain-specific applications
- ▶ Integration with graph databases for knowledge graphs

Conclusions

- ▶ Multi-modal RAG addresses fundamental limitations of text-only systems
- ▶ Semantic-aware parsing maintains information fidelity across modalities
- ▶ Specialized processing agents enable precise query responses
- ▶ Architecture demonstrates significant improvements for structured content
- ▶ Computational complexity requires careful deployment consideration
- ▶ Modular design enables future extensions and optimizations
- ▶ System provides foundation for comprehensive document AI applications

Key Takeaways

- ▶ Real-world documents are inherently multi-modal
- ▶ Traditional RAG systems lose 40% of information in heterogeneous documents
- ▶ Modality-specific processing is essential for accurate retrieval
- ▶ AI-powered descriptions enable semantic search across content types
- ▶ Text-to-SQL bridges natural language and structured data
- ▶ Trade-off: Increased complexity for superior accuracy
- ▶ Future of document AI requires multi-modal understanding

Parsing Resume and RAG

Introduction to Resume Parsing

► What is Resume Parsing?

- Automated process of extracting structured information from unstructured resume documents.
- Converts raw text into a machine-readable format like JSON or XML.

► Why is it Important?

- Saves time for recruiters by automating data entry.
- Enables efficient candidate searching and filtering.
- Improves data quality and consistency in Applicant Tracking Systems (ATS).



Traditional vs. Modern Approaches

- ▶ **Traditional Approaches:**
 - ▶ Rule-based systems using regular expressions (Regex).
 - ▶ Statistical methods like Conditional Random Fields (CRF).
 - ▶ Prone to errors with varied resume formats and language.
- ▶ **Modern Approaches (GenAI):**
 - ▶ Leveraging Large Language Models (LLMs).
 - ▶ Capable of understanding context, semantics, and diverse layouts.
 - ▶ Offers higher accuracy and flexibility (zero-shot or few-shot learning).

Implementation: Specialized API (Docling)

- ▶ **Docling:** A specialized document AI API for parsing specific document types like resumes and invoices.
- ▶ **Pros:**
 - ▶ Highly accurate for its specific domain.
 - ▶ Simple API call; no need to manage models or prompts.
 - ▶ Provides a structured, predictable output.
- ▶ **Cons:**
 - ▶ Can be a "black box"; less customizable.
 - ▶ Dependent on a third-party service and its costs.

Code: Docling Parser

```
1 from docling.document.converter import DocumentConverter
2
3 class DoclingResumeParser:
4     def __init__(self):
5         self.converter = DocumentConverter()
6
7     def parse_resume(self, file_path: str):
8         result = self.converter.convert(file_path)
9         doc = result.document
10
11     # Extract structured information
12     return {
13         "text": doc.export_to_markdown(),
14         "tables": [item for item in doc.tables],
15         "metadata": doc.metadata
16     }
17
```

Listing 1: llm_parsing_docling.py

Implementation: General LLM API (Groq)

- ▶ **Groq:** A platform offering high-speed inference for various open-source LLMs like Gemma and Llama 3.
- ▶ **Pros:**
 - ▶ High flexibility in prompting and output formatting (e.g., JSON mode).
 - ▶ Extremely fast inference speeds.
 - ▶ Cost-effective for many use cases.
- ▶ **Cons:**
 - ▶ Requires careful prompt engineering for best results.
 - ▶ Performance may vary depending on the LLM used and the complexity of the resume.

Code: Groq Parser

```
from groq import Groq
2 import json
4
5 class GroqResumeParser:
6     def __init__(self, api_key: str):
7         self.client = Groq(api_key=api_key)
8         self.model = "llama-3.2-90b-text-preview" # Updated model
9
10    def parse_resume_text(self, text: str):
11        prompt = """Extract the following information from this resume:
12            - name, email, phone, location
13            - education (degree, institution, year)
14            - work_experience (company, title, dates, description)
15            - skills (categorized: technical, soft, languages)
16            - certifications
17
18        Return valid JSON with these exact keys.
19
20        Resume:
21        {text}
22
23        JSON Output:"""
24
```

Listing 2: llm_parsing_groq.py

Code: Groq Parser

```
1 response = self.client.chat.completions.create(
2     messages=[
3         {"role": "system", "content": "You are a resume parsing expert. Output only valid JSON."},
4         {"role": "user", "content": prompt.format(text=text)}
5     ],
6     model=self.model,
7     response_format={"type": "json_object"},
8     temperature=0.1 # Lower for more consistent parsing
9 )
10 return json.loads(response.choices[0].message.content)
11
```

Listing 3: llm_parsing_groq.py

Building a Resume RAG System

► What is RAG?

- **Retrieval-Augmented Generation:** A technique that combines a retriever (to find relevant information) with a generator (an LLM) to produce answers.
- It grounds the LLM's responses in specific data, reducing hallucinations.

► Our RAG Pipeline:

- **Load:** Read resume files (txt, pdf, docx).
- **Chunk:** Split documents semantically using 'SemanticSplitterNodeParser'.
- **Embed & Store:** Convert chunks to vectors ('bge-small-en-v1.5') and store in a FAISS index.
- **Query:** Use Groq's Llama 3 to generate answers based on retrieved context.

Code: RAG Indexing

```
from llama.index.core import VectorStoreIndex, SimpleDirectoryReader, Settings
2 from llama.index.core.node_parser import SemanticSplitterNodeParser
from llama.index.vector_stores.faiss import FaissVectorStore
4 from llama.index.embeddings.huggingface import HuggingFaceEmbedding
import faiss
6
# Configure globally
8 Settings.embed_model = HuggingFaceEmbedding(
    model.name="BAAI/bge-small-en-v1.5"
)
10
# Load and process documents
documents = SimpleDirectoryReader(
    "data",
    file_extractor={
16     ".pdf": "PDFReader", # Specify reader types
     ".docx": "DocxReader",
     ".txt": "TextReader"
    }
)
20 ).load_data()
```

Listing 4: llm_llamaindex_rag.py - Building the Index

Code: RAG Indexing

```
1 # Semantic chunking with proper configuration
3 splitter = SemanticSplitterNodeParser(
    buffer_size=1,
    breakpoint_percentile_threshold=95,
    embed_model=Settings.embed_model
)
5 nodes = splitter.get_nodes_from_documents(documents)
7
9 # FAISS setup with proper dimension
11 d = 384 # Match embedding model dimension
12 faiss_index = faiss.IndexFlatL2(d)
13 vector_store = FaissVectorStore(faiss_index=faiss_index)
15
17 # Build index with storage context
from llama_index.core import StorageContext
18 storage_context = StorageContext.from_defaults(vector_store=vector_store)
19 index = VectorStoreIndex(nodes, storage_context=storage_context)
```

Listing 5: llm_llamaindex_rag.py - Building the Index

Code: RAG Querying

```
from llama.index.llms.groq import Groq
2 from llama.index.core import Settings

4 class ResumeRAG:
    def __init__(self, groq_api_key: str):
        # Setup LLM and Embedding Model
        Settings.llm = Groq(model="llama3-8b-8192", api_key=groq_api_key)
8        Settings.embed_model = HuggingFaceEmbedding(model_name="BAAI/bge-small-en-v1.5")

10       self.index = self._build_index() # From previous slide

12   def query(self, question: str):
13       if not self.index:
14           return "Index not built."
15
16       query_engine = self.index.as_query_engine()
17       response = query_engine.query(question)
18       return response
```

Listing 6: llm_llamaindex_rag.py - Querying

Chatbot UI with Streamlit

- ▶ **Streamlit:** An open-source Python library that makes it easy to create custom web apps for machine learning and data science.
- ▶ **UI Components:**
 - ▶ 'st.title': For the main application title.
 - ▶ 'st.sidebar.text_input': For securely entering the API key.
 - ▶ 'st.file_uploader': To allow users to upload multiple resume files.
 - ▶ 'st.chat_input' and 'st.chat_message': To create an interactive chat interface.
 - ▶ 'st.spinner': To provide feedback during long-running processes like building the index.

Code: Streamlit UI

```
1 import streamlit as st
2
3 st.title("AI-Powered Resume Assistant")
4
5 # API Key and File Uploader in sidebar
6 groq_api_key = st.sidebar.text_input("Groq API Key:", type="password")
7 uploaded_files = st.file_uploader("Upload resumes", accept_multiple_files=True)
8
9 if st.sidebar.button("Build RAG"):
10     # Logic to initialize the RAG system
11     st.session_state.rag_system = ResumeRAG(...)
12
13 # Chat interface
14 if prompt := st.chat_input("Ask a question..."):
15     with st.chat_message("user"):
16         st.markdown(prompt)
17
18     with st.chat_message("assistant"):
19         response = st.session_state.rag_system.query(prompt)
20         st.markdown(response)
```

Listing 7: streamlit_main.py

Challenges in Resume Parsing

- ▶ **Format Diversity:** Resumes have no standard format (e.g., two-column layouts, tables, images).
- ▶ **Ambiguity:** Natural language can be ambiguous. "Java" could be a skill or part of a company name.
- ▶ **Implicit Information:** Dates and timelines often require inference (e.g., "Present" in work experience).
- ▶ **Data Privacy:** Resumes contain Personally Identifiable Information (PII) that must be handled securely.
- ▶ **Scalability and Cost:** Processing thousands of resumes requires an efficient and cost-effective pipeline.

Future Directions

- ▶ **Multimodal Models:** Using models that can understand both text and the visual layout of a resume (e.g., from a PDF).
- ▶ **Advanced Entity Recognition:** Moving beyond basic fields to extract proficiency levels, project details, and soft skills.
- ▶ **Knowledge Graphs:** Building a knowledge graph of candidates, skills, and companies to enable complex relational queries.
- ▶ **Personalized Interaction:** Fine-tuning models on specific company or industry jargon for more accurate parsing and querying.
- ▶ **Proactive Insights:** Developing systems that can proactively suggest candidates for a job description without explicit searching.

Conclusion

- ▶ Generative AI has significantly advanced the field of resume parsing, moving from brittle rule-based systems to flexible, context-aware models.
- ▶ A combination of specialized parsing APIs and general-purpose LLMs offers a powerful toolkit for developers.
- ▶ Retrieval-Augmented Generation (RAG) is a key technology for building interactive and reliable applications on top of parsed resume data.
- ▶ The future is trending towards more sophisticated, multimodal, and proactive talent acquisition systems.

References

► Libraries & APIs:

- Docling: <https://www.docling.io/>
- Groq: <https://groq.com/>
- Llamaindex: <https://www.llamaindex.ai/>
- Streamlit: <https://streamlit.io/>
- FAISS: <https://faiss.ai/>

► Models:

- Gemma: <https://ai.google.dev/gemma>
- Llama 3: <https://ai.meta.com/blog/meta-llama-3/>
- BGE Embeddings: <https://huggingface.co/BAAI/bge-small-en-v1.5>

Production Considerations

Modern Parsing Challenges (2024-2025)

▶ Cloud-Native Documents:

- ▶ Google Docs, Notion, Confluence formats
- ▶ Dynamic content and embedded widgets
- ▶ Real-time collaborative editing artifacts
- ▶ Version control and change tracking

▶ API-Based Extraction:

- ▶ OAuth authentication requirements
- ▶ Rate limiting and pagination handling
- ▶ Incremental sync challenges

▶ Real-Time Streaming:

- ▶ Processing documents as they're created
- ▶ Handling partial/incomplete content
- ▶ Low-latency requirements for live systems

▶ Format Complexity:

- ▶ Rich media embeds (videos, interactive charts)
- ▶ Cross-document references and links
- ▶ Complex nested structures
- ▶ Export format inconsistencies

Handling Collaborative Document Formats

▶ Notion Documents:

- ▶ Block-based structure requires specialized parsing
- ▶ Databases and relations need graph representation
- ▶ Toggle lists and callouts often lost in conversion
- ▶ Solution: Use official Notion API + custom post-processing

▶ Confluence Pages:

- ▶ Macro expansions and dynamic content
- ▶ Nested page hierarchies and attachments
- ▶ Storage format vs. view format discrepancies
- ▶ Solution: REST API extraction + HTML parsing hybrid

▶ Google Workspace:

- ▶ Suggested edits and comments metadata
- ▶ Real-time sync state management
- ▶ Permission-based content visibility
- ▶ Solution: Google Drive API with export format selection

Multimodal Parsing: The Next Frontier

- ▶ **Beyond Text: Understanding All Content Types**
 - ▶ Text, tables, images, charts, diagrams, code blocks
 - ▶ Each modality requires specialized extraction
 - ▶ Traditional parsers lose 40-60% of information
- ▶ **Why Multimodal Matters for RAG:**
 - ▶ Financial reports: Tables contain critical metrics
 - ▶ Technical docs: Diagrams explain architecture
 - ▶ Research papers: Figures show experimental results
 - ▶ Code repositories: Mixed text and code context
- ▶ **Key Technologies:**
 - ▶ Vision-Language Models (VLMs): BLIP-2, LLaVA, GPT-4V
 - ▶ Layout-aware parsers: DocLing, LayoutLMv3
 - ▶ Table extraction: Table Transformer models
 - ▶ OCR + structure: Tesseract + layout analysis
- ▶ *Note: Detailed multimodal implementation covered in next section*

Parser Performance Benchmarks

- ▶ **Benchmark Dataset:** 100 diverse documents (contracts, reports, papers)
- ▶ **Metrics:** Text accuracy, table preservation, layout fidelity

Parser	Text Acc.	Table Acc.	Layout	Overall
PyPDF	85%	35%	40%	53%
Tesseract OCR	78%	62%	55%	65%
Unstructured	89%	71%	68%	76%
LlamaParse	92%	88%	85%	88%
Docling	94%	91%	89%	91%
Azure Doc Intel	96%	94%	92%	94%

- ▶ **Key Insight:** 20-40 point difference between worst and best performers
- ▶ **Impact:** Parser choice often matters more than model selection

RAG Performance: Impact of Parser Quality

- ▶ **Experiment Setup:**

- ▶ Same RAG pipeline (embedding model, retriever, LLM)
- ▶ Same 50 test documents (financial reports, technical docs)
- ▶ Same 200 evaluation questions with ground truth
- ▶ Only variable: document parser

- ▶ **Results - Answer Accuracy:**

- ▶ PyPDF baseline: 58% correct answers
- ▶ Tesseract OCR: 64% (+6 points)
- ▶ Unstructured: 72% (+14 points)
- ▶ LlamaParse: 78% (+20 points)
- ▶ Docing: 81% (+23 points)

- ▶ **Key Finding:** Upgrading parser improved accuracy by 23 percentage points

- ▶ **Comparison:** Upgrading from GPT-3.5 to GPT-4: +12 points (with same parser)

- ▶ **Conclusion:** Parser quality = 2x impact of model upgrade



Performance Benchmarks: Processing Speed

- ▶ **Test Setup:** Intel Xeon 8-core, 32GB RAM, NVIDIA RTX 3080
- ▶ **Document Types:** Research papers, business reports, technical manuals

Doc Size	Pages	CPU Only	GPU	Throughput
Small	1-5	2-5s	1-2s	12-30 docs/min
Medium	10-20	8-15s	3-6s	4-10 docs/min
Large	50-100	45-90s	15-30s	1-2 docs/min
Very Large	200+	3-6min	1-2min	0.3-0.5 docs/min

- ▶ **Factors Affecting Speed:**

- ▶ OCR requirement: +50-100% processing time
- ▶ Complex tables: +20-30% per page
- ▶ Image resolution: Higher DPI = slower processing
- ▶ Layout complexity: Multi-column layouts add 10-20%

Performance Benchmarks: Memory Requirements

► Base Memory Usage:

- ▶ Python process: 500 MB
- ▶ Model loading: 1-2 GB (layout, table models)
- ▶ Per-document processing: 50-200 MB depending on size

► Memory Scaling by Document Size:

Document	Peak RAM	GPU VRAM	Disk Cache
10-page PDF	2.5 GB	1 GB	50 MB
50-page PDF	4 GB	2 GB	200 MB
100-page PDF	6 GB	3 GB	400 MB
500-page PDF	12 GB	6 GB	2 GB

► Memory Optimization Tips:

- ▶ Process large documents in batches
- ▶ Clear cache between documents: `gc.collect()`
- ▶ Use document streaming for very large files
- ▶ Disable features not needed (OCR, table extraction)

Comparison: Docling vs Other Parsers

Feature	PyPDF	Tesseract	Unstructured	LlamaParse	Docling	Azure DI
Text Extraction	✓	✓	✓	✓	✓	✓
Layout Analysis	✗	✗	✓	✓	✓	✓
Table Structure	✗		✓	✓	✓	✓
OCR Support	✗	✓	✓	✓	✓	✓
Reading Order	✗	✗		✓	✓	✓
Formulas/Math	✗	✗	✗		✓	
Multi-format	PDF	Image	Many	PDF	Many	Many
Local Execution	✓	✓	✓	✗	✓	✗
Open Source	✓	✓	✓	✗	✓	✗
Confidence Scores	✗		✗	✗	✓	✓
Structured Output	✗	✗	✓	✓	✓	✓
Speed (10pg)	1s	5-10s	3-5s	8-12s	2-4s	4-8s
Cost/1K docs	\$0	\$0	\$0-50	\$100-200	\$0	\$150-300

- ✓ = Full support, ✗ = Partial support, ✗ = No support
- Docling's Unique Strengths:** Open source, local execution, confidence scores, unified document model

Docling Advantages: Key Differentiators

- ▶ **1. Unified Document Representation**
 - ▶ Single DoclingDocument format for all content types
 - ▶ Consistent API across PDF, DOCX, PPTX, HTML
 - ▶ Preserves hierarchical structure and metadata
- ▶ **2. Advanced Layout Understanding**
 - ▶ Deep learning models for layout analysis
 - ▶ Accurate reading order detection
 - ▶ Multi-column and complex layout support
- ▶ **3. Built-in Quality Assessment**
 - ▶ Confidence scores at page and document level
 - ▶ Quality grades: POOR, FAIR, GOOD, EXCELLENT
 - ▶ Enables automated quality control
- ▶ **4. Production-Ready Features**
 - ▶ Local execution for data privacy
 - ▶ Plugin system for extensibility
 - ▶ Integration with major RAG frameworks
 - ▶ Comprehensive error handling

Error Handling Patterns: Basic Error Handling

► Common Error Scenarios:

- Corrupted or password-protected PDFs
- Unsupported file formats
- Out of memory errors for large documents
- Model loading failures

```
from docling.document.converter import DocumentConverter
2 from docling.exceptions import ConversionError
import logging
4
logging.basicConfig(level=logging.INFO)
6 logger = logging.getLogger(__name__)
```

Error Handling Patterns: Basic Error Handling

```
1 def safe_convert(source: str) -> Optional[DoclingDocument]:
2     converter = DocumentConverter()
3
4     try:
5         result = converter.convert(source)
6
7         # Check conversion quality
8         if result.confidence.mean.grade < 0.5:
9             logger.warning(f"Low quality conversion: {source}")
10
11    return result.document
12
13 except ConversionError as e:
14     logger.error(f"Conversion failed for {source}: {e}")
15     return None
16 except MemoryError:
17     logger.error(f"Out of memory processing: {source}")
18     return None
19 except Exception as e:
20     logger.error(f"Unexpected error for {source}: {e}")
21     return None
```

Error Handling: Advanced Retry Logic

```
from tenacity import retry, stop_after_attempt, wait_exponential
2 from docling.datamodel.pipeline_options import PdfPipelineOptions
4
4 class RobustDoclingConverter:
5     def __init__(self):
6         self.converter = DocumentConverter()
8
9     @retry(
10         stop=stop_after_attempt(3),
11         wait=wait_exponential(multiplier=1, min=2, max=10)
12     )
13     def convert_with_retry(self, source: str):
14         return self.converter.convert(source)
15
16     def convert_with_fallback(self, source: str):
17         # Try with full pipeline first
18         try:
19             return self.convert_with_retry(source)
20         except Exception as e:
21             logger.warning(f"Full pipeline failed, trying simplified: {e}")
22
23         # Fallback: disable expensive features
24         options = PdfPipelineOptions()
25         options.do_ocr = False
26         options.do_table_structure = False
27
28         converter = DocumentConverter(
29             format_options={InputFormat.PDF: PdfFormatOption(options)}
30         )
31         return converter.convert(source)
```



Error Handling: Quality-Based Filtering

```
from typing import List, Tuple
from docling.datamodel.document import DoclingDocument

class QualityFilteredConverter:
    def __init__(self, min_confidence: float = 0.7):
        self.converter = DocumentConverter()
        self.min_confidence = min_confidence

    def convert_and_filter(self, sources: List[str]) -> Tuple[List, List]:
        successful = []
        failed = []

        for source in sources:
            try:
                result = self.converter.convert(source)

                # Check confidence score
                if result.confidence.mean_grade >= self.min_confidence:
                    successful.append({
                        'source': source, 'document': result.document,
                        'confidence': result.confidence.mean_grade })
                else:
                    failed.append({
                        'source': source, 'reason': 'low_confidence',
                        'confidence': result.confidence.mean_grade })
            except Exception as e:
                failed.append({ 'source': source,
                                'reason': str(e), 'confidence': 0.0 })

        return successful, failed
```

Batch Processing: Basic Batch Conversion

- ▶ Initial document corpus ingestion
- ▶ Periodic document updates
- ▶ Large-scale document digitization

```
1 class BatchDocingProcessor:  
2     def __init__(self):  
3         self.converter = DocumentConverter()  
4  
5     def process_directory(self, input_dir: str, output_dir: str):  
6         input_path = Path(input_dir)  
7         output_path = Path(output_dir)  
8         output_path.mkdir(exist_ok=True)  
9  
10        files = list(input_path.glob("/**/*.pdf"))  
11        files.extend(input_path.glob("/**/*.docx"))  
12        files.extend(input_path.glob("/**/*.pptx"))  
13  
14        results = []  
15        for file in tqdm(files, desc="Processing documents"):  
16            try:  
17                result = self.converter.convert(str(file))  
18                output_file = output_path / f'{file.stem}.md'  
19                output_file.write_text(result.document.export_to_markdown())  
20                results.append({'file': file, 'status': 'success'})  
21            except Exception as e:  
22                results.append({'file': file, 'status': 'failed', 'error': str(e)})  
23  
24        return results
```



Batch Processing: Parallel Processing

```
from concurrent.futures import ProcessPoolExecutor, as_completed
2 from multiprocessing import cpu_count
import os
4
5 class ParallelBatchProcessor:
6     def __init__(self, max_workers: int = None):
7         self.max_workers = max_workers or max(1, cpu_count() - 1)
8
9     def process_single(self, file_path: str) -> dict:
10        """Process single file — called in separate process"""
11        converter = DocumentConverter()
12        try:
13            result = converter.convert(file_path)
14            return {
15                'file': file_path,
16                'status': 'success',
17                'pages': len(result.document.pages),
18                'confidence': result.confidence.mean_grade
19            }
20        except Exception as e:
21            return {'file': file_path, 'status': 'failed', 'error': str(e)}
22
23    def process_batch(self, file_paths: List[str]) -> List[dict]:
24        results = []
25        with ProcessPoolExecutor(max_workers=self.max_workers) as executor:
26            futures = {executor.submit(self.process_single, fp): fp
27                       for fp in file_paths}
28
29            for future in tqdm(as_completed(futures), total=len(file_paths)):
30                results.append(future.result())
31
32        return results
```



Batch Processing: Best Practices

- ▶ **1. Resource Management**
 - ▶ Limit concurrent processes to avoid memory exhaustion
 - ▶ Use process pool (not thread pool) to avoid GIL
 - ▶ Monitor memory usage: `psutil.virtual_memory()`
 - ▶ Implement back-pressure mechanisms for large batches
- ▶ **2. Progress Tracking and Logging**
 - ▶ Use `tqdm` for progress bars
 - ▶ Log all failures with traceback for debugging
 - ▶ Save intermediate results periodically
 - ▶ Generate summary reports (success rate, avg time, errors)
- ▶ **3. Fault Tolerance**
 - ▶ Implement checkpointing to resume failed batches
 - ▶ Skip already processed files (check output directory)
 - ▶ Separate retry queue for failed documents
 - ▶ Use exponential backoff for transient errors
- ▶ **4. Output Management**
 - ▶ Use consistent naming conventions
 - ▶ Preserve directory structure in output
 - ▶ Store metadata alongside converted documents
 - ▶ Implement cleanup for failed conversions

Batch Processing: Production Pipeline Example

```
1  class ProductionBatchPipeline:
2      def __init__(self, config: dict):
3          self.converter = DocumentConverter()
4          self.checkpoint_file = config.get('checkpoint', 'progress.json')
5          self.max_retries = config.get('max_retries', 3)
6
7      def load_checkpoint(self) -> set:
8          if Path(self.checkpoint_file).exists():
9              with open(self.checkpoint_file) as f:
10                  return set(json.load(f))
11          return set()
12
13      def save_checkpoint(self, processed: set):
14          with open(self.checkpoint_file, 'w') as f:
15              json.dump(list(processed), f)
16
17      def process_with_monitoring(self, files: List[str]):
18          processed = self.load_checkpoint()
19          remaining = [f for f in files if f not in processed]
20
21          stats = {'success': 0, 'failed': 0, 'skipped': len(processed)}
22
23          for file in tqdm(remaining):
24              result = self._process_with_retry(file)
25              if result['status'] == 'success':
26                  stats['success'] += 1
27                  processed.add(file)
28                  self.save_checkpoint(processed)
29              else:
30                  stats['failed'] += 1
31
32          return stats
```

Performance Optimization Tips

- ▶ **1. Model Loading Optimization**
 - ▶ Load models once, reuse DocumentConverter instance
 - ▶ Pre-load models at startup: reduces first-document latency
 - ▶ Use model caching for repeated conversions
- ▶ **2. Memory Optimization**
 - ▶ Process documents in streaming mode for large files
 - ▶ Clear document cache after processing: `del result`
 - ▶ Use garbage collection: `gc.collect()` between batches
 - ▶ Limit image resolution for OCR when high quality not needed
- ▶ **3. Speed Optimization**
 - ▶ Disable unnecessary features (OCR, table extraction)
 - ▶ Use GPU acceleration when available
 - ▶ Batch similar documents together for better caching
 - ▶ Pre-filter documents by type before processing
- ▶ **4. Quality vs Performance Tradeoff**
 - ▶ Use confidence scores to determine processing depth
 - ▶ Implement fast-path for high-quality digital PDFs
 - ▶ Reserve expensive processing for low-confidence documents

Monitoring and Observability

```
from dataclasses import dataclass
from datetime import datetime
import json

@dataclass
class ConversionMetrics:
    file_path: str
    start_time: datetime
    end_time: datetime
    duration_seconds: float
    pages: int
    confidence_score: float
    memory_peak_mb: float
    status: str
    error: str = None

class MonitoredConverter:
    def __init__(self):
        self.converter = DocumentConverter()
        self.metrics = []

    def convert_with_metrics(self, source: str) -> ConversionMetrics:
        import psutil
        import tracemalloc

        tracemalloc.start()
        process = psutil.Process()
        start_time = datetime.now()

        :
    )
```

Monitoring and Observability

```
1 class MonitoredConverter:
3
4     def convert_with_metrics(self, source: str) -> ConversionMetrics:
5         :
6
7         try:
8             result = self.converter.convert(source)
9             status = 'success'
10            confidence = result.confidence.mean_grade
11            pages = len(result.document.pages)
12            error = None
13        except Exception as e:
14            status = 'failed'
15            confidence = 0.0
16            pages = 0
17            error = str(e)
18
19            current, peak = tracemalloc.get_traced_memory()
20            tracemalloc.stop()
21
22            end_time = datetime.now()
23
24            return ConversionMetrics(
25                file_path=source, start_time=start_time, end_time=end_time,
26                duration_seconds=(end_time - start_time).total_seconds(),
27                pages=pages, confidence_score=confidence,
28                memory_peak_mb=peak / 1024 / 1024, status=status, error=error
29            )
```



Cost Comparison: Parsing Solutions

Solution	Type	Cost/1K docs	Setup	Latency
PyPDF	Open Source	\$0	Easy	0.5s
Tesseract	Open Source	\$0*	Medium	3-5s
Unstructured	Open/Hosted	\$0-\$50	Easy	2-4s
LlamaParse	API	\$100-\$200	Easy	5-8s
Docling	Open Source	\$0*	Medium	2-3s
Azure Doc Intel	Cloud API	\$150-\$300	Easy	3-6s
AWS Textract	Cloud API	\$150-\$500	Easy	4-7s

- ▶ *Compute costs only (CPU/GPU hours)
- ▶ **Open source:** Higher setup, lower marginal cost, full control
- ▶ **API services:** Quick start, predictable pricing, limited customization
- ▶ **Recommendation:** Start with open source, switch to API at scale

Cost-Performance Tradeoff Analysis

- ▶ **Total Cost of Ownership (TCO) for 1M documents/year:**

Parser	API Cost	Compute	Dev/Ops	Total
PyPDF	\$0	\$500	\$2,000	\$2,500
Docling	\$0	\$2,000	\$5,000	\$7,000
LlamaParse	\$100,000	\$0	\$1,000	\$101,000
Azure Doc Intel	\$200,000	\$0	\$1,000	\$201,000

- ▶ **Decision Matrix:**
 - ▶ less than 100K docs/year: Any parser works, optimize for dev time
 - ▶ 100K-1M docs/year: Open source (Docling) offers best ROI
 - ▶ 1M-10M docs/year: Hybrid approach (open source + API fallback)
 - ▶ More than 10M docs/year: Custom solution or negotiated enterprise API pricing
- ▶ **Quality Factor:** Higher accuracy reduces downstream costs (fewer errors, less human review)

Hidden Costs in Production Parsing

► **Infrastructure Costs:**

- ▶ GPU requirements for advanced parsers (Docling, LayoutLM)
- ▶ Storage for raw + processed documents
- ▶ Network bandwidth for cloud API calls
- ▶ Caching infrastructure to reduce redundant processing

► **Operational Costs:**

- ▶ Monitoring and error tracking systems
- ▶ Human review for failed/low-confidence parses (10-20% of docs)
- ▶ Retry logic and fallback parser costs
- ▶ Version upgrades and model retraining

► **Quality Costs:**

- ▶ False negatives in retrieval due to poor parsing
- ▶ User trust loss from incorrect RAG responses
- ▶ Manual data cleaning and correction

► **Rule of Thumb:** Budget 2-3x the direct parsing cost for full production system

Choosing a Parser: Decision Framework

- ▶ **Step 1: Assess Your Document Types**
 - ▶ Simple text PDFs → PyPDF, pypdf
 - ▶ Scanned documents → Tesseract, Docing (with OCR)
 - ▶ Complex layouts/tables → Docing, LlamaParse, Azure Doc Intel
 - ▶ Cloud-native → API-based solutions (Notion API, etc.)
- ▶ **Step 2: Define Quality Requirements**
 - ▶ Acceptable error rate? (Text: ±5%, Tables: ±10%)
 - ▶ Manual review budget available?
 - ▶ Impact of parsing errors on downstream tasks
- ▶ **Step 3: Calculate Total Cost**
 - ▶ API pricing × volume
 - ▶ Compute costs (GPU hours × rate)
 - ▶ Development + operational overhead
- ▶ **Step 4: Prototype and Benchmark**
 - ▶ Test 2-3 parsers on representative sample (50-100 docs)
 - ▶ Measure accuracy, speed, cost
 - ▶ Run end-to-end RAG evaluation

Conclusions

Conclusion: Document Parsing is the Foundation

- ▶ **Parsing Quality Determines RAG Success**
 - ▶ 20-40 point accuracy difference between parsers
 - ▶ Parser upgrades often outperform model upgrades (2x impact)
 - ▶ "Garbage in, garbage out" - no amount of prompt engineering can fix bad parsing
- ▶ **One Size Does Not Fit All**
 - ▶ Document types require different parsing strategies
 - ▶ Simple text PDFs ≠ Complex layouts ≠ Scanned documents ≠ Cloud-native formats
 - ▶ Evaluate parsers on YOUR data, not benchmark datasets
- ▶ **Multi-Modal is the Future**
 - ▶ Real-world documents contain text, tables, images, diagrams, code
 - ▶ Traditional text-only RAG loses 40-60% of information
 - ▶ Modern parsers like Docling preserve all content types
- ▶ **Investment in Parsing Pays Off**
 - ▶ Reduces downstream costs (fewer errors, less manual review)
 - ▶ Enables more accurate and trustworthy RAG systems
 - ▶ Foundation for scalable production deployments

Key Takeaways: Technology and Tools

- ▶ **Open Source Solutions are Production-Ready**
 - ▶ Docling: State-of-the-art parsing with local execution
 - ▶ LlamalIndex: Robust orchestration for RAG pipelines
 - ▶ HuggingFace ecosystem: Embeddings and LLMs without vendor lock-in
 - ▶ Cost-effective for most use cases (less than 1M docs/year)
- ▶ **Docling's Competitive Advantages**
 - ▶ Unified document model across all formats
 - ▶ Built-in confidence scoring for quality control
 - ▶ Advanced layout understanding and reading order
 - ▶ Plugin architecture for extensibility
 - ▶ 91% accuracy on complex documents (vs 53% for basic parsers)
- ▶ **RAG Pipeline Best Practices**
 - ▶ Semantic chunking ↳ Fixed-size chunking
 - ▶ Preserve document structure and metadata
 - ▶ Implement quality filtering based on confidence scores
 - ▶ Use modality-specific processing for tables and images

Implementation Roadmap: From Prototype to Production

Phase 1: Evaluation (Week 1-2)

- ▶ Collect representative document samples (50-100)
- ▶ Test 2-3 parsers (e.g., PyPDF, Unstructured, Docing)
- ▶ Run end-to-end RAG evaluation
- ▶ Measure accuracy, speed, cost
- ▶ Select parser based on data

Phase 2: Prototype (Week 3-4)

- ▶ Implement basic RAG pipeline
- ▶ Configure chunking strategy
- ▶ Set up vector store (FAISS/Chroma)
- ▶ Integrate with LLM
- ▶ Build simple UI (Streamlit)

Phase 3: Optimization (Week 5-6)

- ▶ Add error handling and retry logic
- ▶ Implement batch processing
- ▶ Optimize for memory and speed
- ▶ Add confidence-based filtering
- ▶ Set up monitoring and logging

Phase 4: Production (Week 7-8)

- ▶ Implement checkpointing
- ▶ Add human review workflow
- ▶ Set up CI/CD pipeline
- ▶ Performance testing at scale
- ▶ Documentation and handoff

Common Pitfalls and How to Avoid Them

- ▶ **Pitfall 1: Ignoring Parser Quality**
 - ▶ Problem: "Any parser will do, let's focus on the model"
 - ▶ Solution: Evaluate parsers first, before optimizing other components
 - ▶ Impact: Can improve accuracy by 20-40 points
- ▶ **Pitfall 2: Fixed-Size Chunking**
 - ▶ Problem: Breaking tables, code blocks, logical sections
 - ▶ Solution: Use semantic/hierarchical chunking (Docling + HybridChunker)
 - ▶ Impact: 15-25% improvement in retrieval accuracy
- ▶ **Pitfall 3: Losing Metadata**
 - ▶ Problem: Can't cite sources or verify answers
 - ▶ Solution: Preserve provenance information (page numbers, bounding boxes)
 - ▶ Impact: Enables transparency and trust in RAG outputs
- ▶ **Pitfall 4: No Quality Monitoring**
 - ▶ Problem: Silent failures, degraded accuracy over time
 - ▶ Solution: Use confidence scores, log failures, implement human review
 - ▶ Impact: Catch issues early, maintain system reliability
- ▶ **Pitfall 5: Optimizing Too Early**
 - ▶ Problem: Complex optimizations before understanding bottlenecks
 - ▶ Solution: Start simple, measure, then optimize based on data
 - ▶ Impact: Faster time to production, better ROI

YHK

Future Directions and Emerging Trends

- ▶ **Vision-Language Models for Parsing**
 - ▶ Models like GPT-4V, Gemini understand document layout visually
 - ▶ Can handle complex formats without specialized parsers
 - ▶ Trade-off: Higher cost but better accuracy on edge cases
- ▶ **Streaming and Real-Time Processing**
 - ▶ Parse documents as they're created or edited
 - ▶ Incremental updates to vector stores
 - ▶ Low-latency requirements for live applications
- ▶ **Cross-Document Understanding**
 - ▶ Knowledge graphs connecting entities across documents
 - ▶ Multi-hop reasoning over document collections
 - ▶ Citation networks and reference tracking
- ▶ **Domain-Specific Fine-Tuning**
 - ▶ Custom parser models for specialized domains (legal, medical, financial)
 - ▶ Few-shot learning for new document types
 - ▶ Active learning to improve parser accuracy over time
- ▶ **Automated Quality Assurance**
 - ▶ AI-powered validation of parsing outputs
 - ▶ Anomaly detection for corrupted or unusual documents
 - ▶ Self-healing pipelines that adapt to failures

Final Recommendations

For Beginners:

- ▶ Start with Docling for best out-of-box experience
- ▶ Use LlamaIndex for RAG orchestration
- ▶ Follow the evaluation framework:
 1. Visual inspection ("vibe check")
 2. End-to-end RAG testing
 3. Compare on YOUR documents
- ▶ Don't over-optimize initially
- ▶ Focus on getting working pipeline

For Production Systems:

- ▶ Confidence scores for quality control
- ▶ Set up monitoring and alerting
- ▶ Plan for human-in-the-loop review

Cost Optimization:

- ▶ Less than 100K docs: Any
- ▶ 100K-1M docs: Docling
- ▶ More than 1M docs: Hybrid or custom solution
- ▶ Balance quality vs cost per document

Resources to Explore:

- ▶ Docling: <https://github.com/DS4SD/docling>
- ▶ Parser Benchmarks: <https://github.com/genieincodebottle/parsemypdf>

Remember: Good parsing is the foundation of good RAG. Invest the time to get it right!



Thanks ...

- ▶ Search "**Yogesh Haribhau Kulkarni**" on Google and follow me on LinkedIn and Medium
- ▶ Office Hours: Saturdays, 2 to 3 pm (IST); Free-Open to all; email for appointment.
- ▶ Email: yogeshkulkarni at yahoo dot com



(<https://medium.com/@yogeshharibhaukularkarni>)



(<https://www.linkedin.com/in/yogeshkulkarni/>)



(<https://www.github.com/yogeshhk/>)

YHK