

PDF to XML Conversión

Eshita Shukla

Jui Pitale

Aditya Tumarada

- VJTI
- VJTI
- IIT Guwahati

Introduction



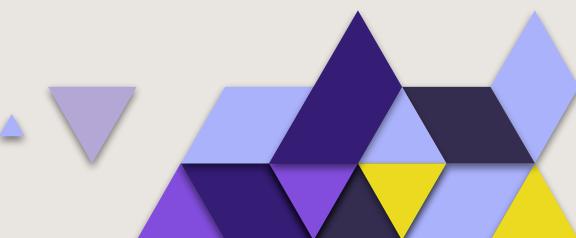
Problem Statement

Write a PDF to XML utility (tool) by leveraging the pdfbox library so that we can use this tool to compare .pdf files to DB tables.

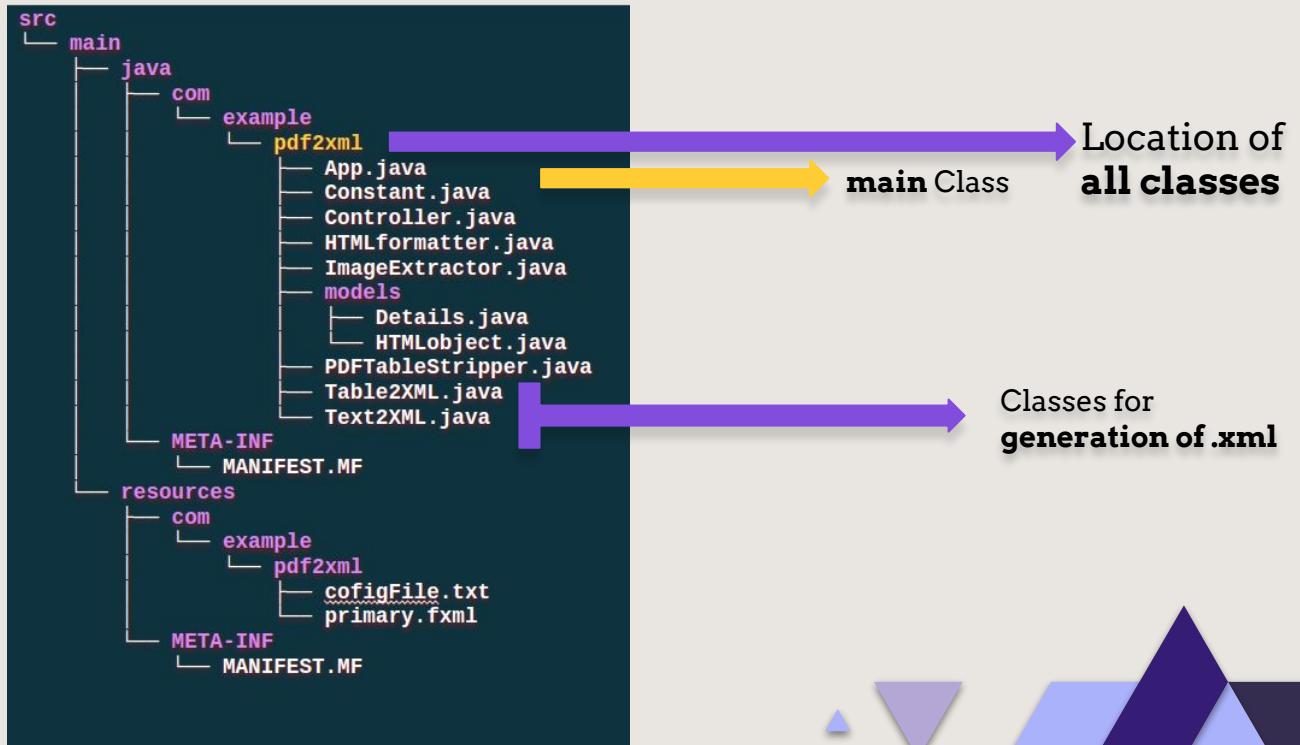
We have a framework, which can compare the XML data with database tables and reports the data differences.

However, the system does not support the data comparison between PDF and table data due to which data is compared manually and leads to error-prone.

So that once we convert the PDF file to the XML file, it can be injected into the existing Framework to compare the data and find the root cause the data discrepancy in the PDF files.



Project Structure



Demo

Approach

Our Target

Implementation
with GUI

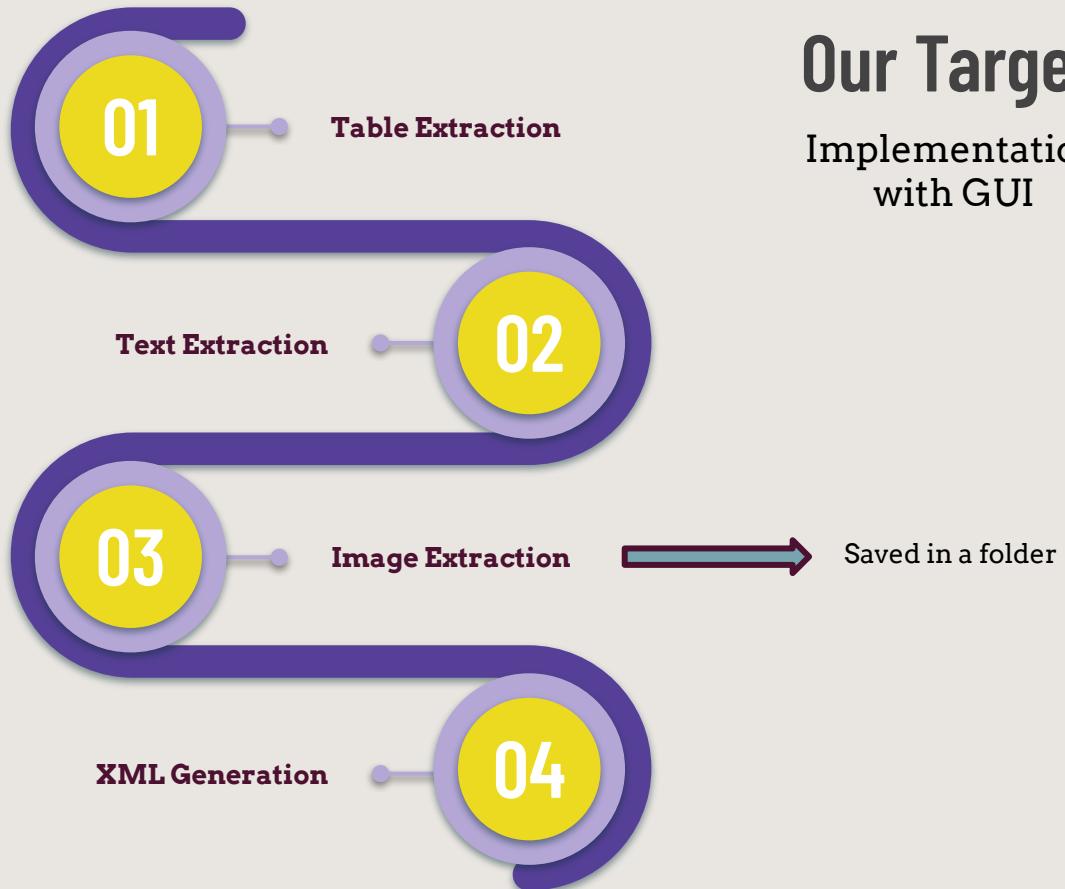


Table Extraction



Group all rows into table
(some in no tables)

purple	purple	purple	purple
purple	purple	purple	purple
purple	purple	purple	purple
white	white	white	white

Group all rows into table
(some in no tables)



Extract rows & their
coordinates

yellow	yellow	yellow	yellow
yellow	yellow	yellow	yellow
yellow	yellow	yellow	yellow
yellow	yellow	yellow	yellow

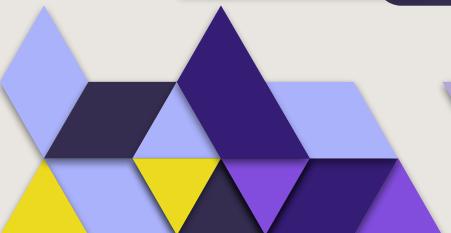
Divide each page into
rectangles

yellow	yellow	yellow	yellow
yellow	yellow	yellow	yellow
yellow	yellow	yellow	yellow
yellow	yellow	yellow	yellow
yellow	yellow	yellow	yellow

Divide each rectangle (row)
into columns

yellow	yellow	yellow

Classify each row ad
"normal" or "heading"



Dividing into rectangles

amazon.in		Tax Invoice/Bill of Supply/Cash Memo (Original for Recipient)						
Sold By : OLD IS GOLD BOOK'S * Gali No-62, H N-629 Santnagar,Near Milan Vihar NEW DELHI, DELHI, 110084 IN		Billing Address : Neelam Shukla 10 B Everest, Anushaktinagar Mumbai 400094 MUMBAI, MAHARASHTRA, 400094 IN						
		State/UT Code: 27						
PAN No: GMFPS6446M GST Registration No: NotApplicable		Shipping Address : Neelam Shukla Neelam Shukla 10 B Everest, Anushaktinagar Mumbai 400094 MUMBAI, MAHARASHTRA, 400094 IN						
		State/UT Code: 27						
Order Number: 403-4367815-0724316 Order Date: 06.03.2020		Invoice Number : IN-9400 Invoice Details : DL-1183551805-1920 Invoice Date : 06.03.2020						
Sl. No	Description	Unit Price	Net Amount	Tax Rate	Tax Type	Tax Amount	Total Amount	
1	Compilers 9332542457 (NX-WAZZ-MW3W)	₹224.00	₹224.00	0%	IGST	₹0.00	₹224.00	
				0%	CGST	₹0.00		
				0%	IGST	₹0.00		
				0%	None	₹0.00		
	Shipping Charges	₹76.00	₹76.00	0%	IGST	₹0.00	₹76.00	
				0%	CGST	₹0.00		
				0%	IGST	₹0.00		
				0%	None	₹0.00		
TOTAL:						₹0.00	₹300.00	
Amount in Words: Three Hundred only								
<p style="text-align: right;">For OLD IS GOLD BOOK'S: <i>Nisha Singh</i> Authorized Signatory</p>								
Whether tax is payable under reverse charge - No								

Classifying into "heading" or not

configFile

Description

DESCRIPTION

No.

Qty

Amount

AMOUNT

Probability = 6/9

Sl. No	Description	Unit Price	Qty	Net Amount	Tax Rate	Tax Type	Tax Amount	Total Amount
1	Compilers 9332542457 (NX-WAZ2-MW3W)	₹224.00	1	₹224.00	0%	IGST	₹0.00	₹224.00
					0%	CGST	₹0.00	
					0%	IGST	₹0.00	
					0%	None	₹0.00	
					0%	IGST	₹0.00	
					0%	CGST	₹0.00	
					0%	IGST	₹0.00	
					0%	None	₹0.00	
TOTAL:							₹0.00	₹300.00

Count	Tax Rate	Tax Type	Tax Amount	Total Amount
1	0%	IGST	₹0.00	₹224.00
	0%	CGST	₹0.00	
	0%	IGST	₹0.00	
	0%	None	₹0.00	
	0%	IGST	₹0.00	
	0%	CGST	₹0.00	
	0%	IGST	₹0.00	
	0%	None	₹0.00	
TOTAL:				₹0.00

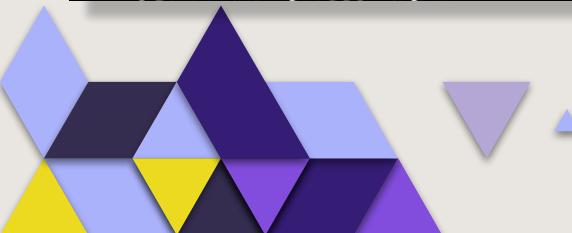
Amount in Words: Three Hundred only	For OLD IS GOLD BOOK'S: Nirbhaya Singh Authorized Signatory
--	---

Whether tax is payable under reverse charge - No

Without probability calculation:
if the keyword
"Description" is present in
the name of the book that
has been bought by the
customer, the algo fails

Probability Calculation for Consecutive rows

Sl. No	Description	Unit Price	Qty	Net Amount	Tax Rate	Tax Type	Tax Amount	Total Amount
1	Compilers 9332542457 (NX-WAZ2-MW3W)	₹224.00	1	₹224.00	0%	IGST	₹0.00	₹224.00
					0%	CGST	₹0.00	
					0%	IGST	₹0.00	
					0%	None	₹0.00	
	Shipping Charges	₹76.00		₹76.00	0%	IGST	₹0.00	₹76.00
					0%	CGST	₹0.00	
					0%	IGST	₹0.00	
					0%	None	₹0.00	
TOTAL:							₹0.00	₹300.00



APPROACH

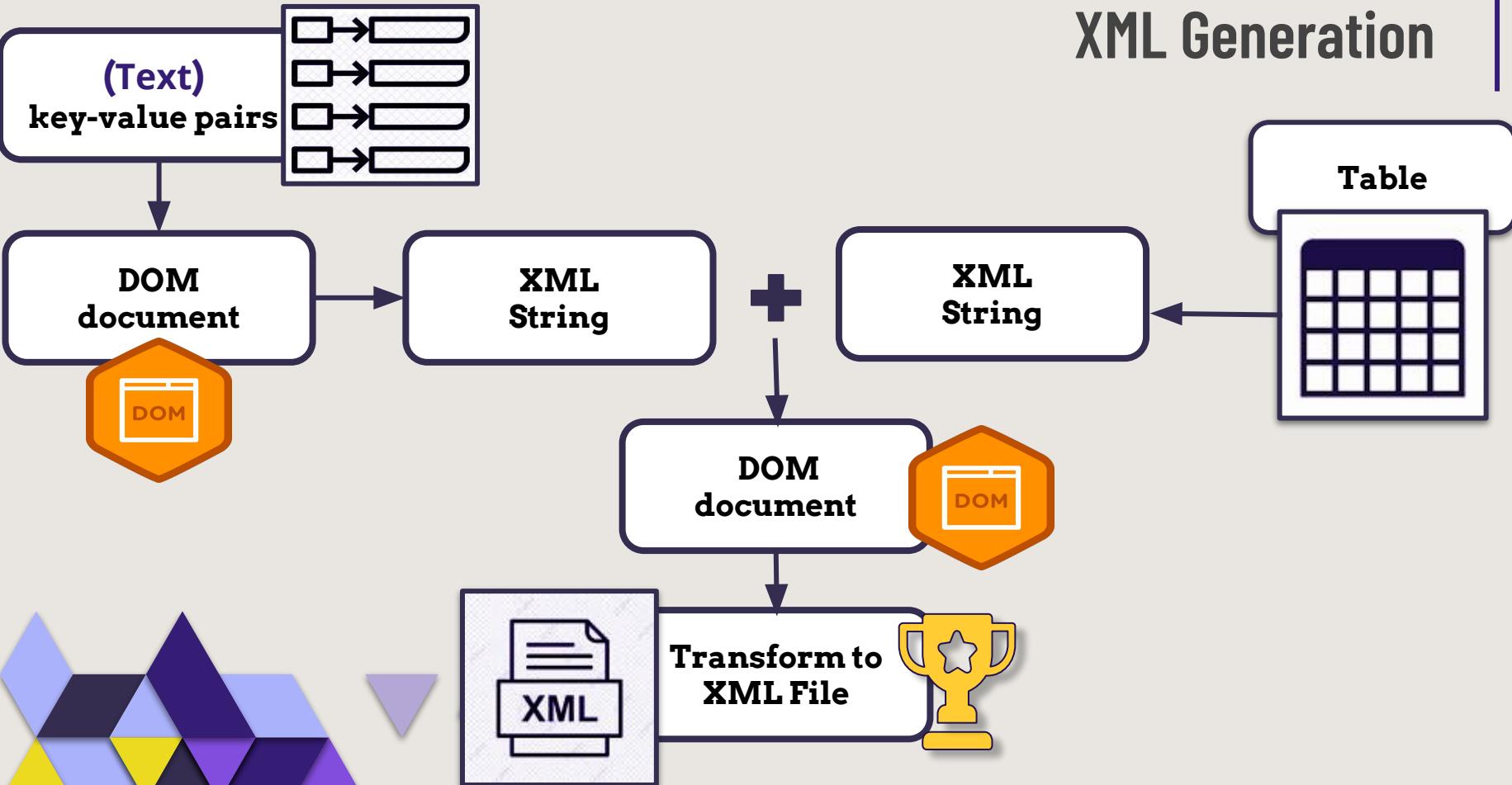
Text Extraction



(page-wise extraction in key value pair form)



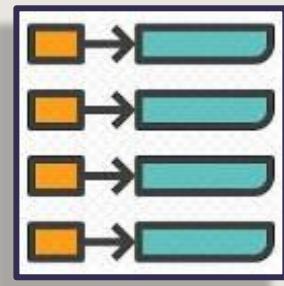
APPROACH



Limitations of the S/W



ReactionGIF.org

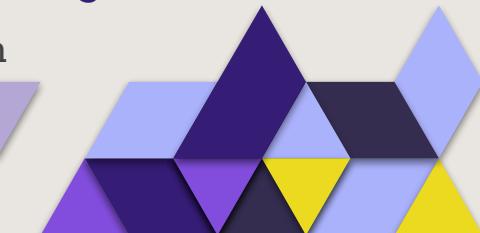


Recognition of the
“key”-“value” pairs

- Change in **text styles**
- **Line-spacing**
- **Could be difficult** to handle



At least some of the
words in the **heading**
should be known





Conclusión

- ❖ Our team has been working on this problem statement for the last **5 weeks**.
 - ❖ We have been able to process **almost all types of PDFs** and successfully extracted all data into an xml, to ensure which we changed our algorithm thrice.
 - ❖ On the **basis of only the alignment**(left/right/center) of the text in each column of the table, we are able to **successfully get the desired XML output**.
 - ❖ The project is made as **flexible** as possible. The user can add more keywords if they want
- 

THANKS!



Questions ??

