



Information Extraction From Fiscal Documents Using LLMs



Vikram Aggarwal¹, Jay Kulkarni², Aditi Mascarenhas², Aakriti Narang², Siddarth Raman², Ajay Shah², Susan Thomas²

¹Google Inc., ²XKDR Forum



1. Motivation

- Governments periodically release fiscal data that is critical for research
- The data format (PDF) is difficult to analyze
- How can LLMs understand these formats?
- How can we create research-ready datasets at scale?

2. Challenge

Large Language Models (LLMs) are great at language, yet poor at multi-page tabular extraction

- Table structure has multi-level, complex hierarchy
- Documents are too big, larger than context window
- Documents are in Indic languages, where LLMs do worse
- Tables are inconsistently rendered, cells are split or merged

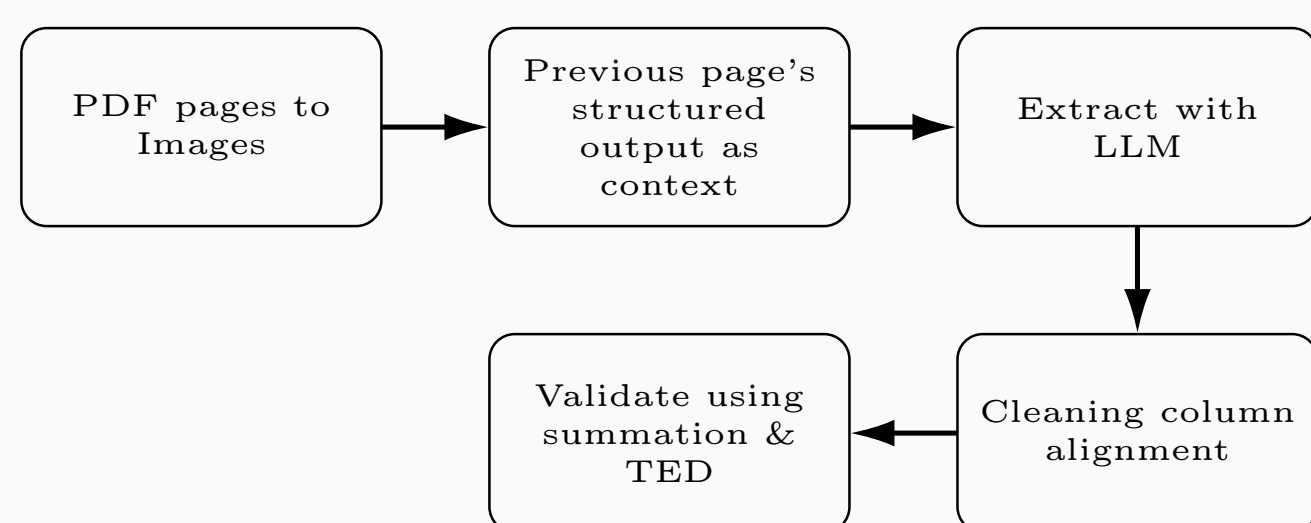
Sources of complexity

- Tables span multiple pages; tables have varying structures
- The files are large: roughly hundreds of pages (84-241 pages in the 7 volumes). Since tables span multiple pages, page-by-page extraction missed vital information
- Regional languages mix with English in inconsistent patterns
- Units and formatting vary across documents

3. Our approach

Our improvements to a naïve LLM document extraction

- **Image-based processing:** Convert PDF pages to high-resolution JPGs (300 DPI). This improves LLM comprehension
- **Sequential context:** Provide previous page's structured output as an aide to understanding the current page.
- **Multi-level validation:** Use column-sums to ensure numerical consistency. Use hierarchy to validate structural consistency
- **Meta-prompting:** Provide domain context, get LLM to write the extraction prompt
- **Intelligent cleaning:** Post-process to improve column alignment



Input PDF to output table

(clockwise) Two pages that capture complexity of the PDF; the output table extracted from the first page

ಸೃಷ್ಟಿ EXPENDITURE 2049 ಬಜ್ಜೆ ಸಂಪನ್ಮೂಲ 2049 Interest Payments									
ಸಂಪನ್ಮೂಲ - 1 ಆದಾಯಾಂಶಗಳ ಸಂಖ್ಯೆ : 29 VOLUME - 1 Demand No : 29 (₹ ಲಕ್ಷ ಸಂಖ್ಯೆ) (₹ in Lakh)									
ಲೆಕ್ಕ ಬಿಡುಗಡೆ, ಆದಾಯಾಂಶಗಳ ಸಂಖ್ಯೆ ಮತ್ತು ಅಂಶ		Heads of Account, Dem No & Dept Code		ಲೆಕ್ಕ Accounts 2018-19	ಬಜೆಟ್ Budget 2019-20	ಸಂಸ್ಕರಣೆ Revised 2019-20	ಬಜೆಟ್ Budget 2020-21		
1		2		3	4	5			
01	ಆಂತರಿಕ ಸಾಲದ ಮೇಲಿನ ಬಜ್ಜೆ	Interest on Internal Debt	C	1235918.78	1541652.00	1499527.55	1829267.00		
03	ಸಣ್ಣ ಉಳಿತಾಯ, ಫಂಡ್ ಮತ್ತು ಫಂಡ್ ಉಳಿತಾಯ ಮೇಲಿನ ಬಜ್ಜೆ	Interest on Small Savings, Provident Funds Etc.	C	235754.02	287656.00	287656.00	320090.00		
04	ಕೇಂದ್ರ ಸರ್ಕಾರದಿಂದ ಪಡೆದ ಸಾಲ ಮತ್ತು ಮುಂಗಡದ ಮೇಲಿನ ಬಜ್ಜೆ	Interest on Loans & Advances from Central Government	C	70610.04	76710.00	76604.00	67400.00		
05	ಮೇಲ್ಕಂಡಿರುವ ಮೇಲಿನ ಬಜ್ಜೆ	Interest on Reserve Funds	C	8.33	12.00	12.00	4882.00		
60	ಇತರ ಹೊಣೆಗಾರಿಕೆಗಳ ಮೇಲಿನ ಬಜ್ಜೆ	Interest on Other Obligations	C	...	1.00	305.83	...		
ಒಟ್ಟು ಬಜ್ಜೆ				Total 2049	C	1542291.17	1906031.00	1864105.38	2221639.00
ಒಟ್ಟು				TOTAL V+C	C	1542291.17	1906031.00	1864105.38	2221639.00
ಒಟ್ಟು ಬಜ್ಜೆ				GRAND TOTAL	C	1542291.17	1906031.00	1864105.38	2221639.00

ಆದಾಯಾಂಶ 2020-21

BUDGET 2020-21

185

ಸೃಷ್ಟಿ EXPENDITURE 2049 ಬಜ್ಜೆ ಸಂಪನ್ಮೂಲ 2049 Interest Payments									
ಸಂಪನ್ಮೂಲ - 1 ಆದಾಯಾಂಶಗಳ ಸಂಖ್ಯೆ : 29 VOLUME - 1 Demand No : 29 (₹ ಲಕ್ಷ ಸಂಖ್ಯೆ) (₹ in Lakh)									
ಲೆಕ್ಕ ಬಿಡುಗಡೆ, ಆದಾಯಾಂಶಗಳ ಸಂಖ್ಯೆ ಮತ್ತು ಅಂಶ		Heads of Account, Dem No & Dept Code		ಲೆಕ್ಕ Accounts 2018-19	ಬಜೆಟ್ Budget 2019-20	ಸಂಸ್ಕರಣೆ Revised 2019-20	ಬಜೆಟ್ Budget 2020-21		
1		2		3	4	5			
01	ಆಂತರಿಕ ಸಾಲದ ಮೇಲಿನ ಬಜ್ಜೆ	Interest on Internal Debt	C	1024916.40	1343614.00	1303940.00	1644445.00		
101	ಮಾರುಕಟ್ಟೆ ಸಾಲಗಳ ಮೇಲಿನ ಬಜ್ಜೆ	Interest on Market Loans	C	5000.00		
102	ಸರ್ಕಾರಿ ಮೇಲಣ್ಣ ಮೇಲಿನ ಬಜ್ಜೆ	Discount on Loans	C	4.55	10.00	
115	ಉಳಿತಾಯ ಉಳಿತಾಯ ಮೇಲಿನ ಬಜ್ಜೆ	Interest on Ways & Means Advances from Reserve Bank of India	C	178961.95	163527.00	163527.00	147784.00		
123	ರಾಜ್ಯ ಸರ್ಕಾರದ ಮೂಲಕ ಕೇಂದ್ರ ಸರ್ಕಾರದ ಬಜ್ಜೆ ಮತ್ತು ಸರ್ಕಾರದ ಯೋಜನೆ ಬಡ್ಡಿ, ಒಳಗಡೆ ಉಳಿತಾಯ ಉಳಿತಾಯ ಮೇಲಿನ ಬಜ್ಜೆ	Interest on Special Securities issued to NSF of the Central Government by State Government	C	28891.49	31050.00	28595.00	27694.00		
200	ಒಳಗಡೆ ಆಂತರಿಕ ಸಾಲಗಳ ಮೇಲಿನ ಬಜ್ಜೆ	Interest on Other Internal Debts	C	3148.94	3461.00	3461.00	4334.00		
305	ಸಾಲದ ನಿರ್ವಹಣೆ	Management of Debt	C	1235918.78	1541652.00	1499527.55	1829267.00		
ಒಟ್ಟು ಬಜ್ಜೆ				Total 01	C	1235918.78	1541652.00	1499527.55	1829267.00
03	ಸಣ್ಣ ಉಳಿತಾಯ, ಫಂಡ್ ಮತ್ತು ಫಂಡ್ ಉಳಿತಾಯ ಮೇಲಿನ ಬಜ್ಜೆ	Interest on Small Savings, Provident Funds Etc.	C	116439.48	139473.00	139473.00	163084.00		
104	ರಾಜ್ಯ ಉಳಿತಾಯ ಮೇಲಿನ ಬಜ್ಜೆ	Interest on State Provident Funds	C	...	1.00	1.00	1.00		
107	ರಾಜ್ಯ ಮತ್ತು ಉಳಿತಾಯ ಮೇಲಿನ ಬಜ್ಜೆ	Interest on Trusts and Endowment	C	119314.54	146881.00	146881.00	155704.00		
108	ಒಮ್ಮೆ ಮತ್ತು ಉಳಿತಾಯ ಮೇಲಿನ ಬಜ್ಜೆ	Interest on Insurance and Pension Fund	C	...	1.00	1.00	1.00		
115	ಒಳಗಡೆ ಉಳಿತಾಯ ಉಳಿತಾಯ ಮೇಲಿನ ಬಜ್ಜೆ	Interest on Other Savings Deposits	C	...	1300.00	1300.00	1300.00		
117	ನಿರ್ದಿಷ್ಟ ಕೆಲಸದ ನಿರ್ದಿಷ್ಟ ಯೋಜನೆ ಮೇಲಿನ ಬಜ್ಜೆ	Interest on Defined Contribution Pension Scheme	C	235754.02	287656.00	287656.00	320090.00		
ಒಟ್ಟು ಬಜ್ಜೆ				Total 03	C	235754.02	287656.00	287656.00	320090.00
04	ಕೇಂದ್ರ ಸರ್ಕಾರದಿಂದ ಪಡೆದ ಸಾಲ ಮತ್ತು ಮುಂಗಡದ ಮೇಲಿನ ಬಜ್ಜೆ	Interest on Loans & Advances from Central Government	C	53690.16	62520.00	62511.00	56064.00		
101	ರಾಜ್ಯ / ಕೇಂದ್ರ ಸರ್ಕಾರದ ಉಳಿತಾಯ ಮೇಲಿನ ಬಜ್ಜೆ	Interest on Loans for State / Union Territory Plan Schemes	C	527.51	483.00	474.00	412.00		
104	ಉಳಿತಾಯ ಮೇಲಿನ ಬಜ್ಜೆ	Interest on Loans for Non-Plan Schemes	C	16364.26	13677.00	13588.00	10900.00		
109	ಉಳಿತಾಯ ಮೇಲಿನ ಬಜ್ಜೆ	Interest on State Plan Loans consolidated in terms of recommendations of	C						

ಆದಾಯಾಂಶ 2020-21

BUDGET 2020-21

186

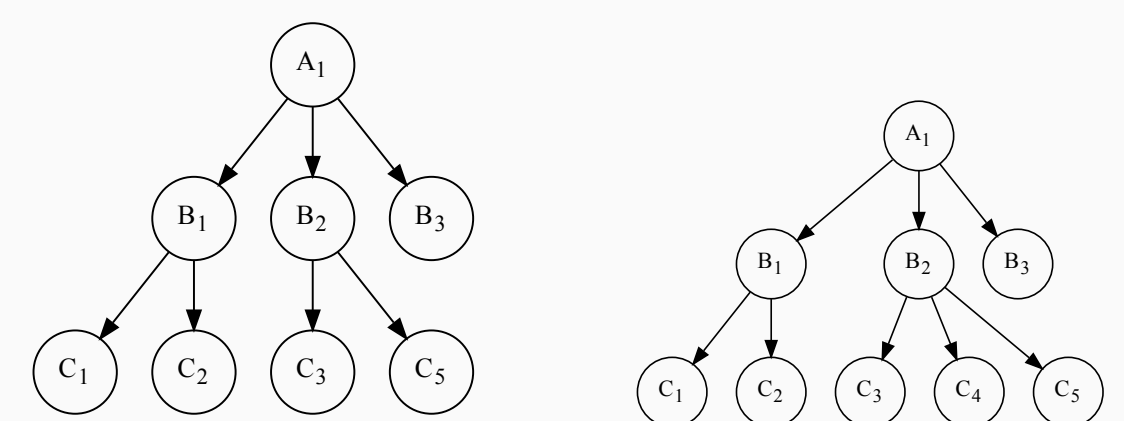
Source Page Num	Vol Num	Demand Number	Major Head Code	Major Head Name	Sub Major Head Code	Sub Major Head Name	Full Account Code	Description	V/C Marker	Row Type	Row Level	Accounts 2018_19	Budget 2019_20	Revised 2019_20	Budget 2020_21
185	1	29	2049	Interest Payments	01	Interest on Internal Debt		Interest on Intern	C	Data	Sub-Major-Head	1235918.78	1541652	1499527.55	1829267
185	1	29	2049	Interest Payments	03	Interest on Small Savings, Provident		Interest on Small	C	Data	Sub-Major-Head	235754.02	287656	287656	320090
185	1	29	2049	Interest Payments	04	Interest on Loans & Advances from		Interest on Loan	C	Data	Sub-Major-Head	70610.04	76710	76604	67400
185	1	29	2049	Interest Payments	05	Interest on Reserve Funds		Interest on Rese	C	Data	Sub-Major-Head	8.33	12	12	4882
185	1	29	2049	Interest Payments	60	Interest on Other Obligations		Interest on Othei	C	Data	Sub-Major-Head			1	305.83
185	1	29	2049	Interest Payments				Total 2049	C	Total	Major-Head	1542291.17	1906031	1864105.38	2221639
185	1	29	2049	Interest Payments				TOTAL V+C	V+C	Total	Major-Head	1542291.17	1906031	1864105.38	2221639
185	1	29	2049	Interest Payments				GRAND TOTAL		Total	Sub-Major-Head	1542291.17	1906031	1864105.38	2221639

4. Validation results

Numerical Consistency: Verify budget head sums within and across schema

Volume	Checks	Passed	Pass Rate %
Volume 1	528	463	88%
Volume 2	463	402	87%
Volume 3	289	233	81%
Volume 4	249	206	83%
Volume 5	390	316	81%
Volume 6	155	134	86%
Volume 7	237	198	84%
All	2311	1952	84%

Structural Consistency: Use Tree Edit Distance Similarity to verify tabular structure



5. Advantages of our method

- Can handle arbitrarily long PDF files. Extract multipage tables from 200+ page PDFs
- Resilient against inconsistent Indic character encoding
- Works in the absence of ground-truth data. Information extraction is at 74%-95% accuracy
- Identifies extraction failure source to page & table location
- Create research-ready datasets for states finances
- Can also create parallel PDF & structured (CSV/JSON) corpus for LLM training