

# Final Project Report: NLP-Driven Financial Report Summarization

## Abstract

We developed a Natural Language Processing (NLP) system to automate the summarization and analysis of financial reports. By utilizing state-of-the-art text preprocessing, Named Entity Recognition (NER), and transformer-based summarization models, our system efficiently extracts key financial metrics and provides concise summaries. The goal is to enhance accessibility and support decision-making processes for stakeholders in the financial domain.

## Prior Related Work

1. **NER Models:** Integrated BERT-based NER and custom models fine-tuned for financial text to extract structured financial entities.
2. **Summarization Models:** Employed transformer-based models like BART and T5, which have shown robust performance in summarizing structured and semi-structured data.
3. **Readability Metrics:** Prior research underscored the importance of assessing readability to ensure generated summaries are accessible to non-technical users.

## Methodology

1. **Environment Setup**
  - Configured Python libraries such as transformers, pandas, nltk, and torch.
  - Utilized cloud platforms with GPU acceleration for faster processing.
2. **Text Preprocessing**
  - **Standardization:** Normalized text by removing non-alphanumeric characters and standardizing formats.
  - **Stopword Removal:** Eliminated redundant terms to focus on essential content.
  - **Abbreviation Expansion:** Addressed financial domain-specific abbreviations using a predefined glossary.
3. **NER Integration**

- Utilized a fine-tuned BERT model for financial text to extract entities such as revenue, PAT (Profit After Tax), and EBITDA.
- Ensured accuracy by training the NER model on curated financial datasets.

#### 4. Summarization

- **Model:** Employed the BART-large-CNN model for summarizing financial text.
- **Implementation:** Integrated summarization output into a pipeline for structured reporting.
- **Readability Assessment:** Used readability indices (e.g., Flesch, Gunning Fog) to evaluate output accessibility.

#### 5. Evaluation Metrics

- Assessed summaries on parameters such as readability, conciseness, compression ratio, and semantic similarity.
- 

## Experiments and Results

#### 1. Evaluation Metrics:

- **Flesch Reading Ease:** 44.14 (indicating slightly difficult-to-read text).
- **Gunning Fog Index:** 10.45 (appropriate for educated readers).
- **Smog Index:** 14.4 (suitable for technical audiences).
- **Semantic Similarity:** 0.451 (moderate alignment with input text).
- **Conciseness Ratio:** 0.0137.
- **Compression Ratio:** 72.93% (highly compressed summaries).

#### 2. Extracted Financial Metrics:

- Revenue: .
- Profit After Tax (PAT): N/A
- EBITDA: 25

## ===== Evaluation =====

### Summary Evaluation Metrics:

Flesch Reading Ease: 44.14

Gunning Fog Index: 10.45

Smog Index: 14.4

Automated Readability Index: 12.8

Semantic Similarity: 0.4513867199420929

Conciseness Ratio: 0.01371094537391655

Compression Ratio: 72.93443104969127

Redundancy (Repeated Words): 2295

3. **Generated Summary:** The system generated summaries reflecting key financial insights. Examples include revenue growth trends, equity funding details, and profitability impacts :

Adani family's equity stake in the Adani portfolio companies . NDTV profit industry-leading profitability . EBITDA and PAT of AWL was impacted on account of hedges . This exceptional financial performance drove our PAT to a record high EBITDA in FY 2023-24 of ` 40,129 crore, marking a substantial 70.8% growth . We have continued to deploy latest 33% growth in revenue from power supply to ` 7,735 crore . During the year, we tapped into diversified sources to raise equity and debt equity . We increased the debt funding pool with a clear roadmap aligned with the project cash flows . We received an equity investment of enhanced revenues, margins and at sustained high efficiency . At the heart of our approach is de-risking of projects, ensuring reliable cash flows and accessing low-cost, management framework . This positions us to comfortably manage debt repayment and finance growth surplus cash flows . The high-quality long-term predictable cash flows will further support our run-rate EBITDA of ` 10,462 crore . Gensuite moderate revenue and operating margins could translate into a loss of prospective revenue . The reliance on long-term debt exposes us to the risk of failure in . Adverse currency movements may stretch our debt repayment obligations . Increase in forex debt repayment beyond projected estimates . 44% of our long-term debt was in Indian currency as of March 31, 2024 . Adani Green is committed to balancing profits with ethics and integrity . Cash Profit = PAT + Depreciation + Deferred Tax + Exceptional Items + Distribution to TOTAL (which is part of finance cost as per Total Revenue (A) Revenue from operations 5,133 7,792 9,220 Other Expenses excluding loss on sale of assets and such one-off expenses . Your Company has recorded revenue from operations to the tune of ` 9,220 crore during the financial year 2023-24 Net profit for the FY 20 23-24 is ` 1,260 crore as compared agencies . The equity authorized share capital of your Company is improvement YOY backed by 95.5% plant availability . Issue of equity shares with differential rights as to support and assistance received from the Government of 2.45 crore .

## Discussion

The project successfully demonstrates the use of NLP for financial data summarization. While the system achieves high compression and readability, challenges remain in handling missing or inconsistent data (e.g., Revenue, PAT :N/A). Enhancing semantic similarity and improving entity extraction for ambiguous data points will be prioritized in future iterations.

## Conclusion

This project highlights the potential of NLP in automating financial report analysis. By integrating NER and summarization models, the system provides a scalable solution for generating concise and informative financial summaries. Future improvements include fine-tuning models on larger, domain-specific datasets and improving output interpretability.

## References

1. **BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding**.<https://arxiv.org/abs/1810.04805>
2. **BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension**. <https://arxiv.org/abs/1910.13461>
3. **Abdaljalil, S., & Bouamor, H. (2021).**  
*An Exploration of Automatic Text Summarization of Financial Reports.*  
chrome-extension://efaidnbmnnnibpcajpcglclefindmkaj/[https://aclanthology.org/2021.finnlp-1.1.pdf?utm\\_source=chatgpt.com](https://aclanthology.org/2021.finnlp-1.1.pdf?utm_source=chatgpt.com)
4. **Jin, H., Zhang, Y., Meng, D., Wang, J., & Tan, J. (2023).**  
*Beyond Pure Text: Summarizing Financial Reports Based on Both Textual and Tabular Data.* [Beyond Pure Text: Summarizing Financial Reports Based on Both Textual and Tabular Data](#)
5. **Bozyiğit, F., & Kılınc, D. (2021).** *Practices of Natural Language Processing in the Finance Sector.* [https://link.springer.com/chapter/10.1007/978-981-16-8997-0\\_9](https://link.springer.com/chapter/10.1007/978-981-16-8997-0_9)

## Contributions: Collective-Work

Hriday.M(se22uari056)

Virat.K(se22uari183)

Saahil.M(se22uari145)

Aashrith Reddy(se22uari092)