# LEVERAGING THE NATURAL LANGUAGE TOOLS AND TECHNIQUES IN ENHANCING THE EFFICACY OF THE TEXT OUTLINE OF FINTECH REQUEST FOR PROPOSALS (RFPs)

**Smriti Narang**
*University of Delhi, New Delhi, India*

## ABSTRACT

*In the present day and age, where tremendous amounts of text-based information are produced consistently, keeping ourselves side by side with new data has become troublesome. Reports in the monetary area recount a quantitative story. On the other hand, the subjective or language content that goes with fiscal reports is a fundamental part of the data set utilized by monetary market members for checking and stewardship. This requires the advancement of proficient innovative strategies for utilizing the presence of these monstrous sums of literary information. Perusing monetary reports like yearly reports is incredibly tedious, and consequently, organizations need to distribute valuable human resources to comprehend, examine, and understand these reports. Accordingly, programmed outline strategies can simplify this assignment by empowering admittance to a more modest yet instructive part of the guaranteed record. In this work, a framework using NLP methods for summing up monetary logs in light of questions given by the client is introduced. This framework conquers the difficulties presented by existing segment-based outline models. As AI strategies have been demonstrated to be compelling on downstream errands, for example, text outline. This work aims to exploit these techniques for clear extractive review and create human-like synopses. The proposed model classifies phrases according to their pertinence by utilizing the strength of solo grouping approaches, and it gets a ROUGE-1 score of 46%.*

## INTRODUCTION

Information is an imperative resource in this century, like what oil was in the last, determined by propels in present day innovation. These days, the assortment and transmission of colossal measures of information parachute the world. With so much information being shared on the web, AI calculations should be fostered that can consequently gather extended texts into brief rundowns and give precise synopses that richly convey the expected items. Archives for example, monetary archives comprise a lot of information and extricating valuable experiences from these records can be a very dreary and tedious interaction.

This data can be consolidated to an edible level utilizing text synopsis. At the point when a client needs to get data from the current exploration works, clients must quickly down to find it. A question-based framework that can look at a report will, in this manner, save time and give the client the data they need. The essential objective of this work is to achieve this objective naturally. Text outline exploration can be separated into nonexclusive and inquiry-based text

synopsis. A nonexclusive text rundown briefly reviews a record that conveys the report's general thought. The interest in consistent inquiry-based outlines is duplicating as various applications create tremendous information measures. This paper presents an inquiry-based rundown of monetary RFPs. The quantity of datasets intended for question-based synopsis is restricted. Furthermore, existing datasets could be more varied in scale and quality. The objective of inquiry-based report rundown is to disengage or create an outline of a record that straightforwardly answers or is connected with the inquiry question. An unaided methodology for text synopsis of RFPs has been utilized.

## DATASET MEASUREMENTS

In this work, we centre around RFPs given by Sarvatra Innovations Private Restricted. An RFP or Solicitation For Proposition is a record used to gather offers for a project or help. The RFP reports were connected with the financial area. Inside the economic area, while the acquisition interaction is at the offering stage, banks issue RFPs to expected merchants. It is often used when picking another monetary specialist organization. Overall, they sort out merchants and solicitation quotes before welcoming offers. All of these financial records were accessible in PDF document design and were written in English. The financial records were significant, with a normal of around 60 pages. Some RFP reports could range from more than 120 pages. A few ordinarily found segments among the RFPs were

- Significant Dates, Qualification Measures, and Annexures. Human-created rundowns or gold synopses were created for evaluation purposes, and the discoveries were contrasted and the created rundown.

## PROPOSED WORK

### A. Technique

An RFP is a task declaration that a bank posts on the web to show how they assess offers from likely assistance providers. For both the bank giving the RFP and the help supplier answering it, the RFP indicates the venture. The dataset that was utilized comprised of different RFPs that were chiefly connected with banking issues.
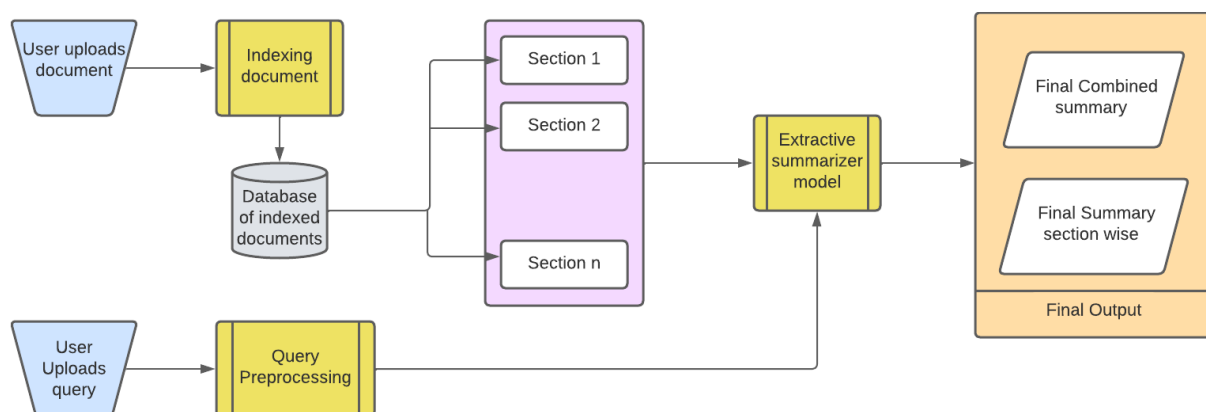


Fig. 1. System architecture diagram for text summarization

Fig. 1. shows the engineering outline that is utilized for text synopsis. In this cycle, preprocessing is the underlying activity. Preprocessing the information is critical because crude details might be arranged conflictingly or not entirely.

Successfully preprocessing raw data can support its rightness, raising project quality and dependability. The initial step was to dispose of the ASCII characters and the list of chapters from the RFP report, as they needed more obvious importance in the report. The accompanying step was to produce text records for each piece of the report and save them in a data set. The whole history was filtered for the different text dimensions. The induction was that the most utilized text dimension is sectioned, given the variety of text dimensions created. That text dimension was given the passage tag. Any text dimensions bigger than that were provided unmistakable header labels. Every assertion in the PDF record was given a title. Afterwards, the views were iterated to lay out the archive's detailed segments.

These parts were then extricated, the substance of which was recorded to a text document and was then put away in the data set. The ensuing stage was to list each section or piece for proficient recovery of parts given the client's question. The ordering activity was done utilizing the Whoosh library. Each part was given legitimate labels and marks for distinguishing proof.

In this work, a question-based extractive synopsis model has been proposed. A custom auto-idea web index for the client input inquiry was made. The search was parsed utilizing Python programming language at the backend. Once the client presents the question, it goes through pre-handling and is shipped off the backend. In the back end, the most relevant recorded segments to the question are separated from the data set. The length of the outline can be concluded by the client utilizing a slider gadget. This upgrades the convenience of the product and gives the client a tweaked insight.

To separate review data, the TextRank calculation was executed. A critical analysis of the formation of the synopsis is word recurrence. If a word or state shows up more now and again, it is vaster and more significant. The relationship between at least two terms in the TextRank calculation is assessed. A square network is made to show how a single word collaborates with every one of the different words in a record. Given its recurrence, a score should be relegated to each expression to create the expression positioning.

The expression positioning chart made for the "Qualification Models." segment of one of the RFP reports is shown in Fig. 2. there are 21 expressions, as should be visible from this figure. The phrases are coordinated in dropping requests of their position, with each file addressing one of the sentences in the segment. The sections with hazier varieties demonstrate that the expressions at that file has a higher recurrence of events in the segment. One method for executing the TextRank calculation is the Python library bundle pyTextRank. It is utilized with the expansion of the spaCy pipeline, which is notable for giving elements like expression acknowledgement. Another strategy is to carry out the TextRank calculation without any preparation. Given your necessities, a technique can be chosen. Finally, a rundown is developed

41

from the selected important areas utilizing the TextRank calculation relying upon the outline length indicated by the client. It is shown to the client on the UI.
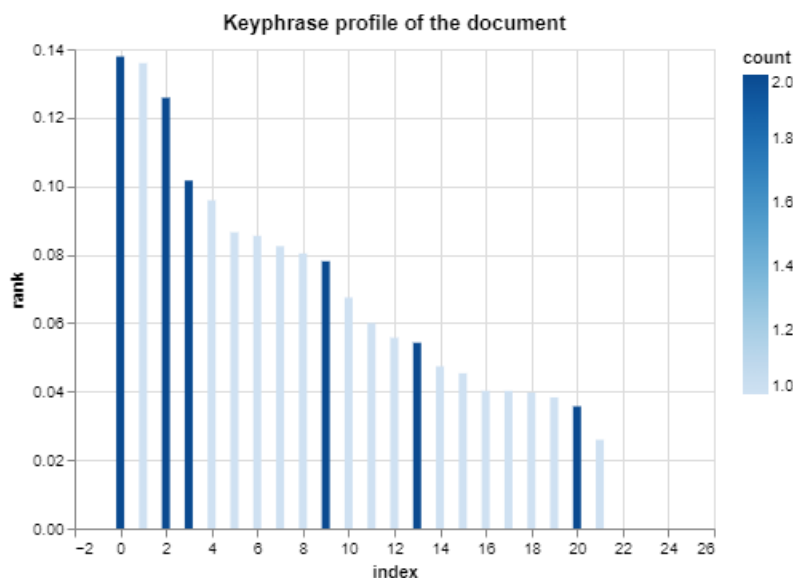


Fig. 2. Graph for key phrase ranking

## B. Calculation

1) Whoosh Ordering: Whoosh is a Python library comprising various capabilities and classes to list your report. Archive ordering is the cycle of allocating names or properties to reports to be effectively looked for and recovered later. Both manual and computerized ordering are accessible. For more modest records, manual ordering is easy to do. However, more extensive records tend to be very drawn-out and tedious. Thus, different programming, APIs, apparatuses, or libraries can assist and make the system simpler.

These files are then used to look through your archive. The library is precious in building custom hunt motors for your application. The whoosh library records each segment with fitting marks and labels in this work. The features most relevant to the client's feedback inquiry are then tracked down utilizing the records.

2) Question Idea: In this work, a question is gathered from the client and used to build the synopsis. The client question is preprocessed before being utilized to waitlist the segments that match the query and give the outline. As the client is composing the question, proposals are given; that is, a rundown of potential questions that the client can choose is shown. These suggestions match the part names present in the RFP archive. This permits clients to view the exact data they seek in the report.

3) TextRank: A calculation in light of PageRank, TextRank is a chart-based model and positioning calculation for text handling. Like PageRank, the sites are positioned by first deciding their weight by building a coordinated chart in which every page is a hub, and the associations between the corners address the sites' interconnections. Thus, while PageRank is

utilized to rank site pages, TextRank is used to type sentences. The central idea is very similar since the text in TextRank is the site page in PageRank.
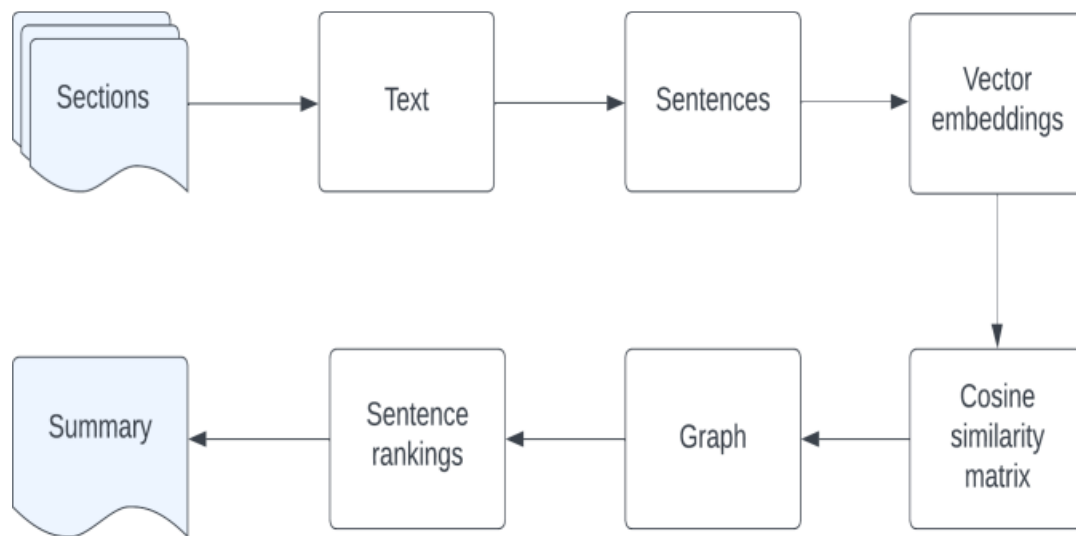


Fig. 3. Flow of TextRank algorithm

The TextRank calculation utilized in this study is delineated in Fig. 3. To start with, the appropriate area is separated into discrete sentences. Using TF-IDF, word embeddings are determined for each sentence. The comparability between the sentences is then determined and kept in a grid. Then, a diagram is made from the framework. The hubs of this diagram address the sentences, and the edge weight addresses the similitude scores. The sentences are positioned in light of the chart, and the highest-level sentences are then linked to make the last synopsis.

## RESULT EXAMINATION

In this work, we utilized the ROUGE score [13] to gauge the nature of the created outlines. ROUGE is an abbreviation that represents Review Situated Student for Listing Assessment. Various measures are utilized for assessing automated text synopsis strategies and machine interpretation. It thinks about an independently produced outline or performance to a set of reference rundowns. These reference outlines are, for the most part, alluded to as gold rundowns. People regularly create these rundowns. There are three assortments of rouge scores: ROUGE-N, ROUGE-S, and ROUGEL.

ROUGE-N searches for unigram, bigram, trigram, and higher request n-gram cross-over, though ROUGE-L searches for the most extended matching arrangement of words utilizing LCS.

ROUGE-1, ROUGE-2, ROUGE-L, and ROUGE-LSum are used to evaluate our work. We present the discoveries from a strategy: The TextRank calculation filled in as a solo learning technique propelled by the PageRank calculation.

| TextRank | Precision | Recall | F1 |
|---|---|---|---|
| ROUGE-1 | 0.580 | 0.386 | 0.464 |
| ROUGE-2 | 0.326 | 0.216 | 0.260 |
| ROUGE-L | 0.360 | 0.240 | 0.288 |
| ROUGE-LSum | 0.360 | 0.240 | 0.288 |

Table 1: ROUGE scores for TextRank summarization model

Compared with related work, the accuracy has improved from 0.37 to 0.58. The review and F measures are genuinely great, too. The ROUGE scores for the top gold outlines had a typical ROUGE-1 F proportion of 0.464. While attempting to deliver brief rundowns, the exactness factor is fundamental. Thus, it is typically desirable to register the accuracy, review, and then the F-Measure. In contrast with the ROUGE-2 score, which demonstrates the cross-over of bigrams between the anticipated outline and the gold synopsis, the ROUGE-1 score, which measures the cross-over of unigrams between the two, is more prominent. This is because of fewer bigram covers as the outlines get progressively longer.

## CONCLUSION

A rundown of monetary records represents a test as the information is addressed in different configurations and is broadly extensive. In this manner, the framework should give a relevant outline for any question the client presents. Therefore, extractive outline techniques can work on this undertaking and give admittance to a more modest yet enlightening piece of a given report. Numerous constraints may be tackled by upgrading the precision of area recognizable proof and extraction. A general language model can be executed to comprehend questions in different dialects. Albeit undesirable relics are taken out during pre-handling, a record with a high level of commotion can diminish the adequacy of the produced synopses. Measurable models don't need high computational ability. Nonetheless, outlines created can feel disengaged and lose indispensable source text data. Another essential component is to embrace

44

an extraneous audit by enrolling the assistance of various money experts to evaluate the nature of the synopses made naturally.

# REFERENCES

[1] Abdaljalil, Samir, and Houda Bouamor. "An exploration of automatic text summarization of financial reports." In Proceedings of the Third Workshop on Financial Technology and Natural Language Processing, pp. 1-7. 2021

[2] El-Haj, M., AbuRa'ed, A., Litvak, M. & Pittaras, N., Giannakopoulos, G. (2020). The Financial Narrative Summarisation Shared Task (FNS 2020).

[3] Van Lierde, H., Chow, T. W. S. (2019). Query-oriented text summarization based on hypergraph transversals. Information Processing and Management, 56(4), 1317–1338.

[4] Khan, R., Qian, Y., Naeem, S. (2019). Extractive-based text summarization using KMEANS and TF-IDF. International Journal of Information Engineering and Electronic Business, 11(3), 33–44.

[5] Araci, D. (2019, August 27). FinBERT: Financial Sentiment Analysis with Pre-trained Language Models.

[6] Laskar, M. T. R. (2020, November 14). Utilizing Bidirectional Encoder Representations from Transformers.

[7] Hernandez-Castaneda, A., Garcia-Hernandez, R. A., Ledeneva, Y. Millan-Hernandez, C. E. (2020). Extractive Automatic Text Summarization Based on Lexical-Semantic Keywords. IEEE Access, 8, 49896–49907.

[8] Yong, S. P., Abidin, A. I. Z., & Chen, Y. Y. (2006). A neural-based text summarization system. Data Mining VII: Data, Text and Web Mining and Their Business Applications.

[9] Zheng, S., Lu, A.,& Cardie, C. (2020). SUMSUM@FNS-2020 Shared Task

[10 ] Litvak, Marina, Natalia Vanetik, and Zvi Puchinsky. "Hierarchical summarization of financial reports with RUNNER." In Proceedings of the 1st Joint Workshop on Financial Narrative Processing and MultiLing Financial Summarisation, pp. 213-225. 2020..

[11 ] Azzi, Abderrahim Ait, and Juyeon Kang. "Extractive summarization system for annual reports." In Proceedings of the 1st Joint Workshop on Financial Narrative Processing and MultiLing Financial Summarisation, pp. 143-147. 2020.

[12 ] Mingjun Zhao, Shengli Yan, Bang Liu, Xinwang Zhong, Qian Hao, Haolan Chen, Di Niu, Bowei Long, Weidong Guo. QBSUM: a Large-Scale Query-Based Document Summarization Dataset from Real-world Applications

[13 ] Lin, Chin-Yew. "Rouge: A package for automatic evaluation of summaries." In Text summarization branches out, pp. 74-81. 2004.

[14 ] Lin, H., Bilmes, J. (2010). Multi-document summarization via budgeted maximization of submodular functions. HLT '10, Human language technologies: The 2010 annual conference of the north American chapter of the association for computational linguistics. St roudsburg, PA, USA: ACL912–920.

[15 ] Teh, Y. W., Jordan, M. I., Beal, M. J., Blei, D. M. (2004). Sharing clusters among related groups: hierarchical Dirichlet processes. Proceedings of the seventeenth international conference on neural information processing systems, NIP S'04. Cambridge, MA, USA: MIT Press1385–1392.