# HYBRID DOCUMENT SUMMERIZATION

Submitted by,

Naveen Pandurangi(1VK15CS034)

Nikhil Chandran(1VK15CS038)

Srivatsa V Jamadagni(1VK15CS054)

Mohammed Abu Talha Ahmed(1VK15CS030)

Under The Guidance Of

Mrs. Chandramma B.E,M.Tech,Ph.D

Assoc.Professor ,Dept Of CSE,VKIT

# TABLE OF CONTENTS

- **Introduction**
- **Problem statement**
- **Existing system and its disadvantages**
- **Proposed system**
- **Significance of the project**
- **Methodology**
- **Expected outcomes**
- **Conclusion**
- References

# HYBRID DOCUMENT SUMMERIZATION

## INTRODUCTION :

- Document Summarization is the technique by which the huge parts of content are retrieved.

- The Document Summarization plays out the summarization task by unsupervised learning system.

- The significance of a sentence in info content is assessed by the assistance of Simplified Lesk calculation.

- Document Summarization assesses the weights most sentences utilizing the Simplified Lesk calculation and arranges them in decreasing order of their weights.

# PROBLEM STATEMENT

- Reading the whole document, dismembering it.

- Isolating the critical thoughts from the crude content require some serious energy and exertion.

- Existing system increase the human effort while creating a synopsis.

- A few vital products compress records as well as website pages.

- The persons cannot quickly determine which points are imported for reading.

# EXISTING SYSTEM AND ITS DISADVANTAGES

- Usually summarization is done by hand by the editor or user.

- Most of the methods currently in use do summarization for text .

- Summarization of document is hardly even possible

- Most methods utilize supervised learning.

- This method is tedious, time consuming and very inefficient.

# PROPOSED SYSTEM- Objectives

• A solitary or single input content is going to be outlined utilizing unsupervised learning.

• Sentences with induced weights are composed in sliding solicitation concerning their weights.

• Certain quantities of sentences are chosen as an outline.

• The proposed computations, abridges solitary or single report content utilizing unsupervised learning approach.

• Heaviness of every sentence in a substance is resolved using streamlined Lesk's computation and wordnet.

• Summarization procedure is performed as indicated by the given rate of synopsis.

# SIGNIFICANCE OF THE PROJECT

- It requires some serious energy and exertion.

- Summarization of a 600 word doc takes atleast 10 minutes.

- Programming summarizes 5000 words doc in a brief instant.

- Requires less information to get the most essential data.

- It reduces the human effort while creating a synopsis.

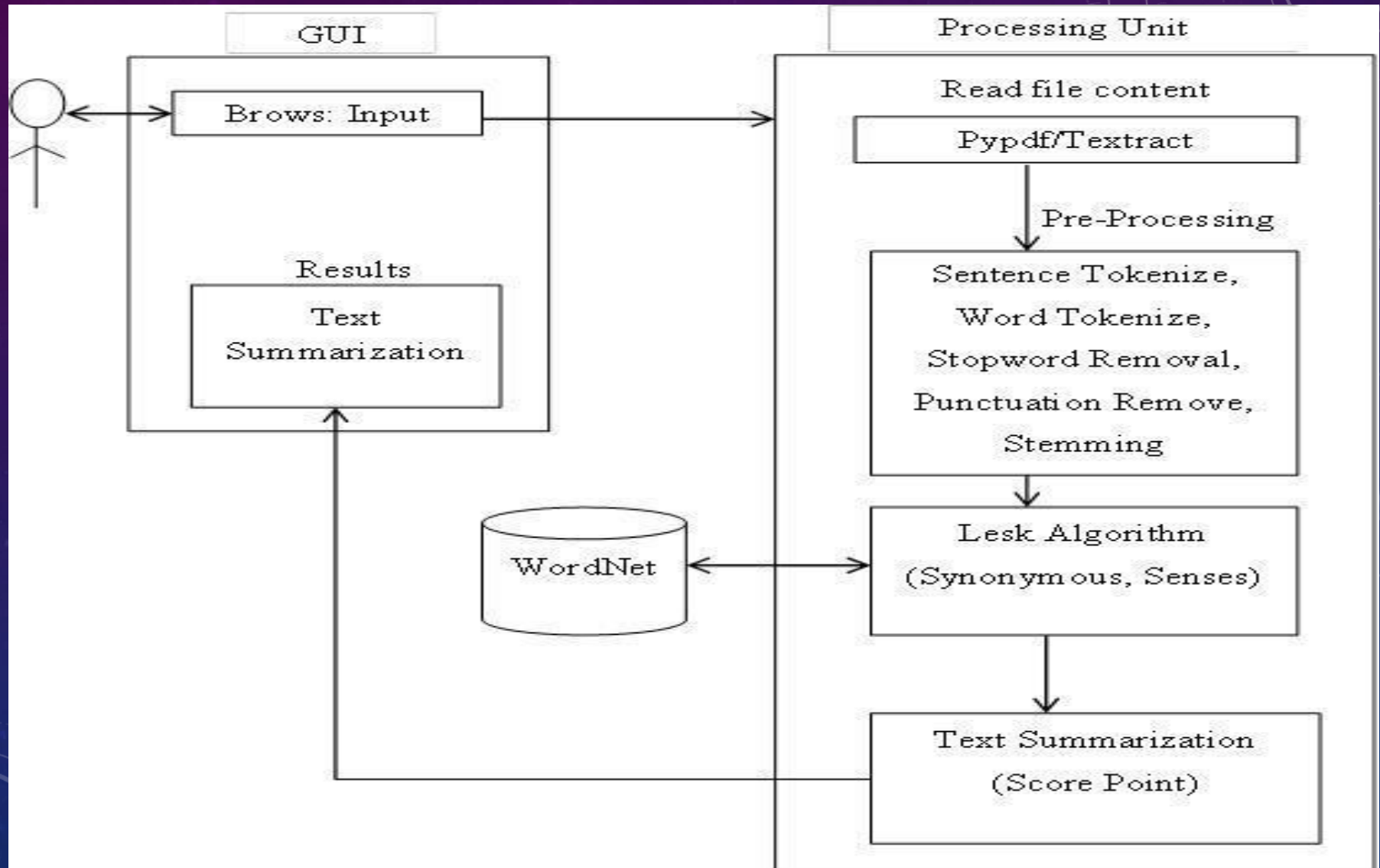- A few vital products compress records as well as website pages.

# METHODOLOGY-NLTK TOOL KIT

- The **Natural Language Toolkit** is a suite of libraries and programs for symbolic and statistical natural language processing (NLP).

- It was developed by Steven Bird and Edward Loper in the Department of Computer and Information Science at the University of Pennsylvania.

- NLTK includes graphical demonstrations and sample data.

- It is accompanied by a book that explains the underlying concepts behind the language processing tasks supported by the toolkit.

# METHODOLOGY-WORDNET

• **WordNet** is a lexical database for the English language.

• It groups English words into sets of synonyms called *synsets*.

• WordNet can is a combination of dictionary, and thesaurus.

• Primarily used  in automatic text analysis and artificial intelligence  applications.

• The database and software tools are freely available for download from the WordNet website.

• Both the lexicographic data (*lexicographer files*) and the compiler (called *grind*) for producing the distributed database are available.

# METHODOLOGY-
# DATA FLOW MODEL

# METHODOLOGY-LESK ALGORITHM

• The Lesk algorithm is a classical algorithm for word sense disembogues introduced by Michael. E.Lesk in 1986.

• The Lesk algorithm is based on the assumption that words in a given "neighbourhood" will tend to share a common topic.

• A simplified version of the Lesk algorithm is to check the definition of a word with the terms contained in its neighbourhood.

• Versions have been adapted to use wordnet An implementation might look like this:
Calculation 1: This calculation compresses a single report content utilizing unsupervised learning approach. In This approach , the heaviness of each sentence in a content is determined utilizing Improved Lesk calculation and WordNet

# THE ALGORITHM

**Step 1:** The list of distinct sentences of the content is prepared.

**Step 2:** Repeat steps 3 to 7 for each of the sentences.

**Step 3:** A sentence is gotten from the list.

**Step 4:** Stop words are expelled from the sentence as they don't take an interest straightforwardly in sense assessment system.

**Step 5**: Glosses(dictionary definitions) of all the important words are extricated utilizing the WordNet.

**Step 6:** Intersection is performed between the sparkles and the information content itself.

**Step 7:** Summation of all the crossing point comes about speaks to the heaviness of the sentence.

**Step 8:** Weight appointed sentences are arranged in descending request concerning their weights.

**Step 9:** Desired number of sentences are chosen by the level of summarization.

**Step 10:** Selected sentences are re-orchestrated by their real sequence in the info content.

**Step 11:** Stop.

# METHODOLOGY-SINGLE DOCUMENT SUMMARIZATION

**Input:**

- Text Data for which Summary is required.

- Value of N – for generating top N frequent Terms.

**Output:**

- Summary for the Original Text Data.

- Compression Ratio.

- Retention Ratio.

**Steps:**

- Data Preprocessing Phase Retrieve data

- Eliminate Stop Word

- For the entire text content

- Get the N frequent Terms

- Generate Term-Frequency List

- For all N-Frequent Terms

- Generate Sentences from the Original Data

- If the sentence consists of a term that is present in frequent-terms-list then

- add sentence to summary-sentence-list.

- Calculate Compression Ratio and Retention Ratio

# EXPECTED OUTCOMES

SAMPLE PARAGRAPH:

Thank you all so very much. Thank you to the Academy. Thank you to all of you in this room. I have to congratulate the other incredible nominees this year. *The Revenant* was the product of the tireless efforts of an unbelievable cast and crew.

# TOKENIZATION:

1. Thank you all so verymuch.

2. Thank you to theAcademy.

3. Thank you to all of you in this room.

4. I have to congratulate the other incredible nominees this year.

5. *The Revenant* was the product of the tireless efforts of an
   - unbelievable cast andcrew.

# PRE-PROCESS:

1. thank youall so very much

2. thank you to theacademy

3. thank you to all of you in this room

4. i have to congratulate the other incredible nominees this year

5. *the revenant* was the product of the tireless efforts of an unbelievable cast and crew

# HISTOGRAM:

| Word | Count | Word | Count | Word | Count |
|------|-------|------|-------|------|-------|
| thank | 3 | in | 1 | revenant | 1 |
| you | 4 | this | 2 | was | 1 |
| all | 2 | room | 1 | product | 1 |
| so | 1 | i | 1 | tireless | 1 |
| very | 1 | have | 1 | efforts | 1 |
| much | 1 | congratulate | 1 | an | 1 |
| to | 3 | other | 1 | unbelievable | 1 |
| the | 5 | incredible | 1 | cast | 1 |
| academy | 1 | nominees | 1 | and | 1 |
| of | 3 | year | 1 | crew | 3 |

# WEIGHTED HISTOGRAM:

| Word | Weight | Word | Weight | Word | Weight |
|------|--------|------|--------|------|--------|
| thank | 3/5 | in | 1/5 | revenant | 1/5 |
| you | 4/5 | this | 2/5 | was | 1/5 |
| all | 2/5 | room | 1/5 | product | 1/5 |
| so | 1/5 | i | 1/5 | tireless | 1/5 |
| very | 1/5 | have | 1/5 | efforts | 1/5 |
| much | 1/5 | congratulate | 1/5 | an | 1/5 |
| to | 3/5 | other | 1/5 | unbelievable | 1/5 |
| the | 5/5 | incredible | 1/5 | cast | 1/5 |
| academy | 1/5 | nominees | 1/5 | and | 1/5 |
| of | 3/5 | year | 1/5 | crew | 1/5 |

**MAXIMUM IS 5**

# SENTENCE SCORES:

| | |
|---|---|
| thank | 0.5 |
| you | 0.8 |
| All | 0.4 |
| So | 0.2 |
| Very | 0.2 |
| Much | 0.2 |
| | 2.3 |

# SENTENCE SCORES BY ORDER:

| Sentence | Score |
|---|---|
| *the revenant* was the product of the tireless efforts of an unbelievable cast and crew | 6.2 |
| thank you to all of you in this room | 4.3 |
| i have to congratulate the other incredible nominees this year | 3.4 |
| thank you to the academy | 3.1 |
| thank you all so very much | 2.3 |

# PICK n LARGEST SENTENCES:

| Sentence | Score |
|---|---|
| *the revenant* was the product of the tireless efforts of an unbelievable cast and crew | 6.2 |
| thank you to all of you in this room | 4.3 |

HERE n = 2

# EXPECTED OUTCOME

*The Revenant* was the product of the tireless efforts of an unbelievable cast and crew.

Thank you to all of you in this room.

# CONCLUSION

• **Automatic Text Summarization approach depends on upon the semantic data of the concentration in a substance.**

• **So this way, gathered parameters like approaches, spots of different substances are not considered.**

• **In this recommendation, Lesk mean for word sense disambiguation by utilizing the vocabulary definitions to the electronic dictionary information base on utilizing wordnet.**

• **This goal is clear from covering sentence, couple of fusing words that give the setting of the word, in this not utilizing the late using the definitional shines of those words, other than those of words related to them through with the unmistakable relations portrayed in wordnet.**

# REFERENCES

[1] Alexander M., Rush Sumit Chopra, Jason Weston- A Neural Attention Model for Abstractive Sentence Summarization arXiv:1509.00685v2 [cs.CL] 3 Sep 2015

[2] Çaglar˘Gulçehre˙ Bing Xiang, Ramesh Nallapati, Bowen Zhou, Cicerodos Santos - Abstractive Text Summarization using Sequence-to-sequence RNNs and Beyond, arXiv:1602.06023v5 [cs.CL] 26 Aug 2016

[3] Santosh Kumar Bharti, Korra Sathya Babu, Sanjay Kumar Jena - Automatic Keyword Extraction for Text Summarization: A Survey, National Institute of Technology, Rourkela, Odisha 769008 India e-mail@nitrkl.ac.in 08-February-2017

[4] Abigail See, Peter J. Liu, Christopher D. Manning - Get To The Point: Summarization with Pointer-Generator NetworksarXiv:1704.04368v2 [cs.CL] 25 Apr 2017

[5] Mehdi Allahyari, Seyedamin Pouriyeh, Mehdi Assefi, Saeid Safaei, Elizabeth D. Trippe, Juan B. Gutierrez, Krys Kochut - Text Summarization Techniques: A Brief Survey, arXiv:1707.02268v3 [cs.CL] 28 Jul 2017

# THANK YOU