# Automatic Condensation of Electronic Document

**Madhusudan M. Oza[1] Pratibha P. Dhundale[2] Vishal S. Baviskar[3] Falgun R. Gahane[4] Amardip R Sabale[5]**

[1,2,3,4,5]Department of Computer Engineering

[1,2,3,4,5]SSBT's COET, Bambhori, Jalgaon. North Maharashtra University, Jalgaon (M.S.), India

*Abstract—* In today's world there are huge amount of data is available on Internet. Everyone is used to with retrieval of this data/info from different sites, applications in direct or indirect way. Many times it is necessary to summarize the articles or documents in short but also important format. Summarization via manual method is an older technique. Each human may have different summary however, some important aspect of text is always present in everyone's summary. In practical, humans have limitations of speed and time. It reduces the performance as per time of work and time of input articles are increasing. It is pretty good if same is done by Computer which has intelligence! Text Summarization is a way to generate a text, which contains the important portion of information of the original text or texts. Several techniques are generated depending upon many parameters to find the summary as the type, position and format of the sentences in an input text, formats of different words occurrence of a particular word in a text etc. The solution to be discussed is about to use Automatic Method for summarization.

*Key words:* TRM, Integer Array for Processing

## I. INTRODUCTION

Summary is a text that is formed from one or more text, that contains a important portion of the information from the original text, and that is not more than half of the original text.[3]

This trouble-free definition derives three important properties that characterize automatic summarization [3]:

1) Summaries may be formed from a solo document or several documents
2) Summaries should protect vital information
3) Summaries should be short

'Text' here includes multimedia documents, on-line documents, hypertexts etc. of the several types of summary that have been famous, analytic summaries (that give an idea of what the text is about, with lack of actual content) and informative summery (that give some reduced version of the content) are frequently referenced. Extracts are summaries formed by reusing portions (sentences, words, etc.) of the key text, while abstracts are shaped by re-producing the extracted content.[1]

Need to note that, this solution is based on abstraction technique.

## II. PROPOSED WORK

### A. Objective:

Proposed system should work on Paragraphs for the Automatic Condensation of Electronics Document's Text. It should take all words as input. Then it will create weight for each word. This weight is based on occurrence of a word in whole text document. By using similarity function system should create value in between 0 to 1 for all possible pairs of words. Now consider each paragraph as node. Based on numerical similarity joined them to create task relationship map. Select a typical threshold value (min value) and select the paragraphs above the threshold value.

Now perform the text traversal using convenient techniques from paragraph to paragraph to select the important sentences. These sentences will logically combine for an output file.

### B. Workflow:

Salton et al. employs the technique to produce intra document links between paragraphs of a document, and get a text relationship map (TRM) according to those links. Each node on the TRM indicates a paragraph and it is represented by a vector of weighted terms. If two nodes have high similarity then, a link is created between them, which is typically computed as the inner product between the vectors of the corresponding paragraphs [2].

Simply, two nodes, they are said to be semantically related if there is a link in them. to measure paragraphs importance, bushiness of a paragraph is defined. The bushiness of a paragraph is the number of links connecting it to other paragraphs. For example, the bushiness of P2 in Fig. 1 (finding bushy node) is 5 because it linked with P1, P3, P4, p5 and P6.
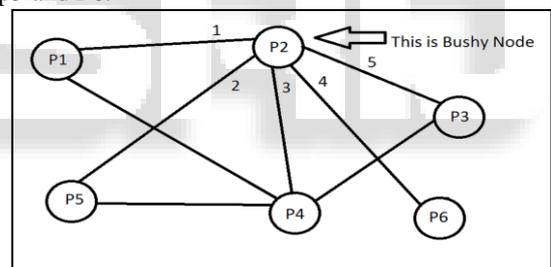


Fig. 1: Text Relationship Map (TRM)

A highly bushy node linked with many other nodes; simply, it has a lot of overlapping vocabulary with other nodes; thus, a highly bushy node is likely to discuss main topics that are covered in many other paragraphs.

As for text summarization, we proposed method to generate the summary according to the bushiness of paragraphs: global bushy path. It identifies paragraphs with high bushiness but traverse them in different text order. System focuses on Bushy Method.
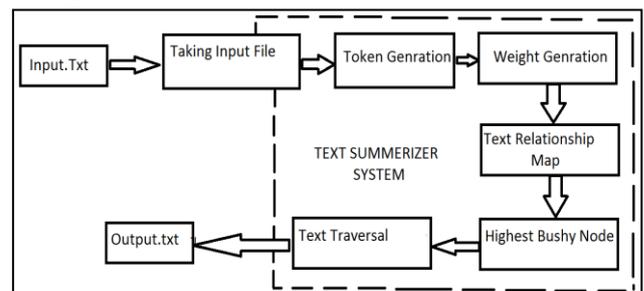
### C. Algorithm:



Fig. 2: Detailed Workflow

Above fig. 2, shows detailed workflow of text summarizer system. Following algorithm is a backbone of workflow of the system:

1) Take Input in txt format
2) Generate tokens for each paragraph i.e. separate each valid word from remaining document. This separate word is called as token.
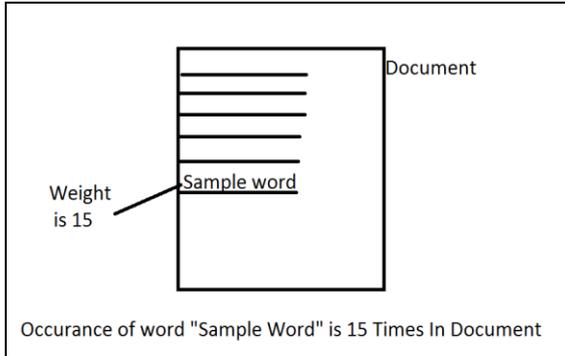3) Find occurrence of each word .i.e. how many time the 'word' occurred in whole text. (fig3)

Fig. 3: Occurrence of Word

4) Generate Links between Nodes based on similarity of words in each node. This graphical presentation is called Text Relationship Map(TRM ) (fig 1)
5) Find the highest bushy node by checking which node is having maximum links in TRM (fig 1)
6) Create a similarity function that checks total occurrence of a 'word' between nodes and compare this occurrence with threshold value. If, it is greater than threshold value then word is valid for consideration. Store this word in separate file with its occurrence. ex. If word "Develop" is have total occurrences in whole document is 57% and threshold value is 40% then "Develop" is a valid word to take in consideration.

In this proposed solution 40% threshold value is used.

7) After storing all valid words as per above discussion, extract sentences from actual document which contains this valid words. Here first priority is to a highest bushy node then, remaining less bushy nodes i.e. priority to extract sentences from nodes is from highest bushy node to lowest bushy node. This step is called as Text Traversal. Out of the final summary 40% summary is taken from highest bushy node. Cause it contains most overlapping vocabulary. Remaining 60% summary is taken from all less bushy nodes.
8) After extracting sentences by using above step, rearrange sentences in proper order. This re-arrangement is required because in step 7 we are focusing on priority of bushy nodes to extract sentences but, the generated output must be in original sequence as provided in input text.
9) Put generated summary in output file.

## III. RESULT AND DISCUSSION

This section contains sample snapshots of systems important processing and its explanation. These snapshots are from system developed by our team.
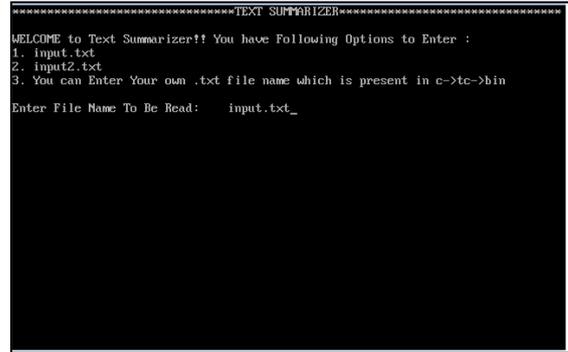
### A. Welcome Screen:

Fig. 4: Welcome Screen

Above Fig. 4, is snapshot of Welcome screen. This screen is providing direction to user that how to use the software. Provided two options are input files named input.txt and input2.txt respectively. These two files are present in bin folder of TurboC++ directory. There is third option that user can add its own input file by providing name of that file. The file must be present in bin folder of TurboC++ directory.
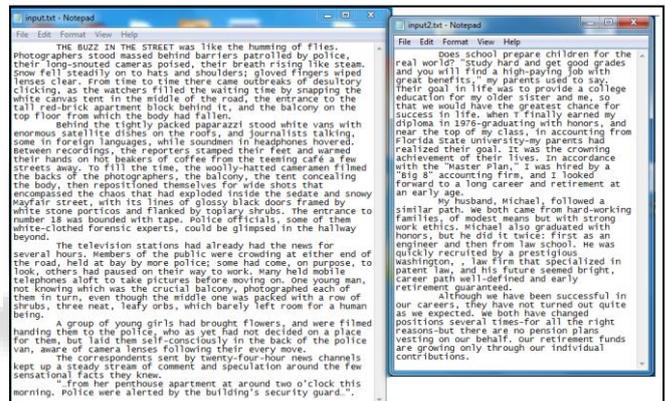
### B. Input files:

Fig. 5: Two Sample Input Files

Above Fig. 5 is showing two input files respectively named input.txt and input2.txt.
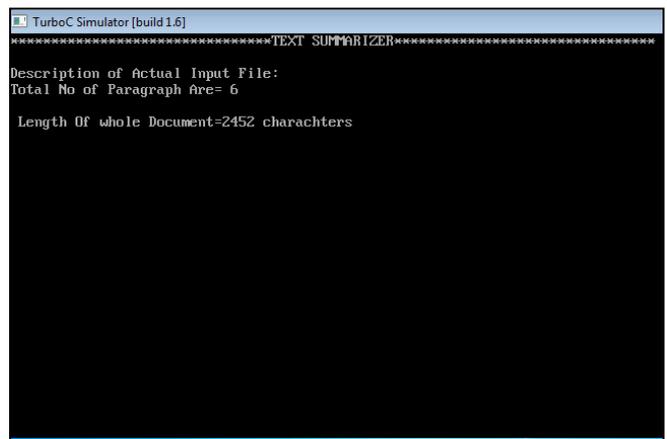
### C. Input Files Information:

Fig. 6: Description of Input.Txt File

Above fig. 6, is snapshot after providing input.txt file name. This screen indicates description of input.txt file. It shows total number of paragraphs and total number of characters in input.txt file including spaces and special symbols.

*D. Integer Array for Processing:*



Fig. 7: Sorted Array of Integer Numbers

Above fig 7 is snapshot which is showing sorted array of integer numbers. This array includes integer occurrence of all words in input.txt file. Next line in snapshot shows 20% and 30% length of array. This length is used in text traversing i.e. out of all array 20% top occurrence words are taken in consideration from highest bushy node and 30% top occurrence words are taken in consideration from non-bushy node.

*E. Generated Summary:*



Fig. 8: Resulted Summary

Above fig 8 is a snapshot of Resulted summary. After summary, it includes total number of sentences of actual input file and total number of sentences in summary.

## IV. CONCLUSION AND FUTURE SCOPE

The proposed system focus on paragraph extraction mechanism. Early mentioned bushy method is with little bit complex but powerful algorithm. System requirement is economical for both development and deployment side. All the risks from different point of view are to taken in consideration. Hence, this system is feasible to use. Result of generated summary contains near about all important sentences from actual input file. Result is also short as compare with actual input file.

Hence, from above discussion we can say that this system is useful and reliable.

After implementation of this proposed system there are huge future scopes for summarizer.
Some of them listed below:
1) Text Summarizer can be developed for all kind of document such as- .docx, .doc, .PDF...etc.
2) In future Text Summarizer can be a very important module of Robotics field. like, A robot have to learn new culture of society by reading different books in less time period.
3) Text summarizer needs to develop with consideration of similarity of words. Like word "use" and "utilize" have same meaning hence, occurrence of word "use" need to increase with 2.

## V. ACKNOWLEDGMENT

REFERENCES

[1] Salton, G., Singhal, A., Mitra, M., & Buckley, C. (1997). Automatic text structuring and summarization. Information Processing & Management, 33(2), 193207.
[2] Salton, G., Singhal, A., Mitra, M., & Buckley, C. (1997). Automatic text summarization by Paragraph Extraction
[3] Alok Ranjan Pal, Projjwal Kumar Maiti and Diganta Saha, An Approach To Automatic Text Summarization Using Simplified Lesk Algorithm And Wordnet, International Journal of Control Theory and Computer Modeling (IJCTCM) Vol.3, No.4/5, September 2013.
[4] Young, S. R., & Hayes, P. J. (1985). Automatic classification and summarization of banking telexes. In proceedings of the 2nd Conference on Artificial Intelligence Applications (pp. 402408).
[5] Maheedhar Kolla(2002). Automatic Text Summarization Using Lexical Chains: Algorithms And Experiments, Department of Mathematics and Computer Science University of Lethbridge, Canada.
[6] Jen-Yuan Yeh, Hao-Ren Ke, Wei-Pang Yang, I-Heng Meng , Text summarization using a trainable summarizer and latent semantic analysis, Information Processing and management 41 (2005) 7595.
[7] Banko, M., Mittal, V., Kantrowitz, M., and Goldstein, J.1999. Generating extraction based summaries from handwritten summaries by aligning text spans. In Proceedings of PACLING99.