# Extractive multi-document text summarization based on graph independent sets

Taner Uçkan [a,*], Ali Karcı [b]

[a] Van Yüzüncü Yıl University, Computer Programming Department, 65000 Van, Turkey
[b] İnönü University, Department of Computer Engineering, 44000 Malatya, Turkey

## ARTICLE INFO

## ABSTRACT

We propose a novel methodology for extractive, generic summarization of text documents. The Maximum Independent Set, which has not been used previously in any summarization study, has been utilized within the context of this study. In addition, a text processing tool, which we named KUSH, is suggested in order to preserve the semantic cohesion between sentences in the representation stage of introductory texts. Our anticipation was that the set of sentences corresponding to the nodes in the independent set should be excluded from the summary. Based on this anticipation, the nodes forming the Independent Set on the graphs are identified and removed from the graph. Thus, prior to quantification of the effect of the nodes on the global graph, a limitation is applied on the documents to be summarized. This limitation prevents repetition of word groups to be included in the summary. Performance of the proposed approach on the Document Understanding Conference (DUC-2002 and DUC-2004) datasets was calculated using ROUGE evaluation metrics. The developed model achieved a 0.38072 ROUGE performance value for 100-word summaries, 0.51954 for 200-word summaries, and 0.59208 for 400-word summaries. The values reported throughout the experimental processes of the study reveal the contribution of this innovative method.

© 2019 Production and hosting by Elsevier B.V. on behalf of Faculty of Computers and Artificial Intelligence, Cairo University. This is an open access article under the CC BY-NC-ND license (http://creativecommons.org/licenses/by-nc-nd/4.0/).

## 1. Introduction

With the rapid development of Internet technology, the volume of electronic documents on the Internet have increased with extraordinary speed. In today's era of rapid data growth, people can acquire and share information instantaneously from various sources. The Internet already provides access to billions of documents. As this increases every second, an exponential growth of data has been seen over a short period of time. Search engines are able to list the most related documents or web pages by utilizing several user inputs. However, even the most developed search engines lack the capacity to synthesize the information they retrieve. Therefore, this rampant growth of information raises the problem of information management, necessitating tools that can process data as a means of coping with such problems [1]. Despite the latest progress seen in document summarization, the problem has not yet been fully resolved. Automated document summarization aims to find methods to detect the most important information in text, and to subsequently condense it for ease of use by readers [2]. Moreover, it is necessary to obtain certain features by processing data for the purpose of shortening the time spent accessing information [3]. These issues have increased interest in the field of automated text summarization systems [1].

This abundance of data availability undoubtedly enhances human life, but at the same time makes the quick access of accurate information increasingly difficult [4]. Automated text summarization technologies are therefore increasingly studied by researchers so as to achieve greater efficiencies through enhanced or new methods [5]. However, despite all of the studies conducted in the field of document summarization, the need for improvement and innovation has not diminished [2]. Automated document summarization is a significant subtopic of natural language processing (NLP), which has the objective to present long text documents in a compressed and comprehensible form [6]. Document summarization

* Corresponding author.
*E-mail addresses:* taneruckan@yyu.edu.tr (T. Uçkan), ali.karci@inonu.edu.tr (A. Karcı).

techniques can generally be divided into two categories: extractive and abstractive. Extractive document summarization consists of three stages: representation of texts, sentence scoring, and sentence selection. Abstractive document summaries interpret the main content of documents by using natural language generation approaches, and then re-expresses it to form a summary [7,8]. They can also be categorized as single-document and multi-document, depending on the number of documents to be summarized. The first studies were conducted as single-document summarization. Multi-document summarization studies have started to be performed, and methods have been developed for application to more than one textual document [6,9,10]. Summaries can also be categorized as either generic or query-focused [11]. The vast majority of studies conducted by researchers have been on the basis of generic summarization. In this type of summarization, a few limited assumptions about the purpose of forming the summary are made, and the general content is maintained while trying to cover as much information as possible. On the other hand, query-focused summarization is aimed at summarization based on queries determined by the user and information related to the subject in the text is returned [12–14]. In the proposed framework, graph independent sets [15] for Automated Document Summarization are utilized within a unique approach using intensive Graph Theory techniques.

As a stage of the proposed summarization system, a new text processing tool called KUSH (named after its developers: Karcı, Uçkan, Seyyarer, and Hark) is used to ensure that the connections and relationships between the sentences are transferred to the representative graphs in the most accurate way. This innovative software tool has had a very positive effect on performance values in terms of transferring semantically distinguishable and more accurately measurable relationships between the sentences and the corresponding graph.

In the current study, the performance of the proposed document summarization method was tested using two publicly available datasets; DUC-2002 and DUC-2004. The performance results are presented and compared with various existing methods. The innovative contributions of the proposed method can be summarized as follows:

- Using the KUSH text processing tool, graphs are obtained by eliminating uncertainties regarding words and their meanings, thus achieving maximum graph representation.
- Two important concepts were combined to model an extractive text summarizer: Textual Graph and Maximum Independent Sets.
- Maximum Independent Set was determined on the graph of the texts obtained after the text developed processing tool.
- In the next stage, the nodes forming the Independent Set were removed from the representative graphs. In this way, it was possible to determine the main ideas and concepts that should be included in the summary of the documents to be summarized.
- No previous document summarization study has utilized Maximum Independent Sets. However, the current study found that it provides a quite robust and simple framework as a result of summarizing the texts based on a mathematical approach.

The remainder of the study is organized as follows: In Section 3, information is provided on the general stages of the proposed summarization method, Textual Graph, Proposed Text Preprocessing Tool, and graph independent sets. In Section 4, information is given on the dataset and evaluation metrics used, the experimental results of the proposed method for summarizing the texts are presented, and the proposed model is compared with the state-of-the-art methods. Finally, in Section 5, the experimental results are discussed and interpreted.

## 2. Related work

Scoring of phrases or sentences and obtaining summaries is the most common method used in automated extractive summarization. Sentence scoring is adopted in the majority of methods applied today. Scoring methods are classified as word scoring, sentence scoring, and graph scoring [16]. In the word scoring methods, a scoring is made considering the importance of the sentences containing the frequency of a word in the text [17–19], with words such as proper nouns, places and objects that are considered as a determinant being scored higher [20,7]. In the text scoring methods, the formal properties of the words (emboldened, italicized, underlined) are taken into account [21]. In addition, sentences starting with phrases such as "Briefly," "Finally" and "As a result" in the text are defined as sign phrases, and the sentences following these statements are noted as being important sentences [18]. Similarly, evaluation is based on the title of the text to be summarized. Sentences that contain words found in the title are considered to be added to the summary, and their importance levels are increased accordingly [22]. Sentence scoring methods also take into account the dimensions of sentences, attaching greater importance to sentences of a larger size [21,23]. Points are assigned to sentences by determining the position of the sentence and whether or not it involves numerical values [18,20,24].

The authors of Ref. [25] described the design and evaluation of extractive summarization approach as a way to assist learners with reading difficulties. Graph-based representations are frequently used in text analysis methods as they provide very effective solutions. In Ref. [9], the authors proposed TextRank, which contains graph-based representation for summarization using the intersections of the text contents. Similarly, LexRank was introduced in Ref. [13] utilizing an eigenvector centrality-based algorithm, one of the node centrality methods. Both the TexRank and LexRank algorithms were inspired by the PageRank [26] algorithm, a document summarization framework presented to obtain central sentences in a document using mutual information between term and sentence sets [27].

In Ref. [28], a multi-layered representation of documents, sentences and words was employed. The authors of Ref. [29] described documents with graphs in their study by utilizing link generation for automated document summarization. They defined the structure by revealing the text relationships in the documents, and evaluating the summaries by comparing them with those created by human hand. In Ref. [30], a graph-based approach was presented in order to provide semantic continuity, with nodes corresponding to the terms of the documents and the edges reflecting semantic relationships between these nodes. Basically, a graph diameter calculation is performed for all the nodes in the graph, and the shortest and longest paths are described as the weakest and strongest bonds. While graph structures and documents were defined in Ref. [31], nodes and edges were created based on local similarities.

Random Walk has been used to obtain summaries of the main documents. In Ref. [32], a summarization system for the biomedical field was proposed. Using a system called Unified Medical Language, a graph was obtained based on concepts and relationships with a semi-dictionary-based application, and then the PageRank algorithm was applied. In Ref. [12], the authors proposed a new graph based on reinforced random walk.

Although most researchers have focused on extractive text summarization, some have worked effectively on abstractive summarizing. In Ref. [33], abstractive multi-sentence summarization was performed using the RNN structure. Similarly, Refs. [34,35] were also works in the field of abstractive summarization.

In the current study, we introduce an innovative and extremely simple text summarization system by taking graph-based

automated document summarization studies one step further. The aforementioned examples from the literature are presented in Table 1 as a historical view of this research field.

As can be seen from the literature, there have been a number of graph-based studies. The differences in terms of the proposed study begin primarily at the preprocessing phase. Using the preprocessing tool developed for this study, paragraphs are preprocessed in order that relations between the sentences are defined to the maximum level. The other and most significant difference of the proposed method is the process of finding independent clusters by analyzing the nodes in the graphs. When determining the maximum independent set clusters, the degree of independence of the nodes in the graphs are taken into account. In extractive summarization studies, the aspiration is that the sentences to be selected should have maximum coverage value while determining the sentences to be included in the summary. In other words, it is expected that the rate of carrying the main idea of a sentence to be selected will be high. This requires that the sentence has a higher number of words compared to the other sentences, i.e., it is not an independent sentence. By using maximum independent sets, independent nodes are found and sentences far from the main idea of a paragraph are selected, thereby increasing the likelihood of sentences being selected for the summary.

Today, the problem of finding maximum independent sets still remains as a NP-hard problem. The process of finding maximum independent sets has not yet reached an optimal result, and can be seen in the current study as a weakness. New methods to arrive at a solution for this problem will be pursued meticulously.

## 3. Proposed summarization method

A block diagram of the stages of the Document Summarization framework proposed for text summarization is shown in Fig. 1.

In this study, a graph-based generic, extractive and multi-document summarization method is presented to extract appropriate summaries from the given texts. The proposed document summarization method consists of three main stages. In the first stage, non-discriminatory stop words (e.g., pronouns, prepositions, conjunctions) were removed from the dataset. A number of preprocesses are performed by the developed KUSH text preprocessing tool. When forming graphs of close-meaning words that differ from each other in terms of their spelling but are semantically

derived from a common word root, they are prevented from being processed as if different words. In the second stage, word commonalities between the phrases are represented mathematically and graphically. In addition, this stage includes determination of the nodes forming the Maximum Independent Set and the removal of the sentences corresponding to these nodes from the main graph. The final stage is the weighting of the sentences that make up the documents using the eigenvector node centrality method and the selection of important sentences. At this stage of the proposed method, the Top *N* phrases are combined and 200- and 400-word summaries created separately. The success of the framework was then tested in detail using a number of different ROUGE performance metrics.

### 3.1. Text preprocessing and KUSH

Most dense datasets are not configured, with few of the relevant datasets being of a structured format. Structured datasets are those that can be expressed in rows and columns of a table or utilize a number of tags. In the current study, it was possible to provide a certain structure to the data being studied by disregarding complex and time-consuming processes. Some level of preprocessing is necessary for datasets that do not have a specific integrity but are required to be structured. It is predicted that converting texts that will be used as input to the system into workable formats and isolating them from unnecessary data that are not discriminatory will increase the performance of the system because dense clusters are formed in natural language [3].

In text summarization processes, the data to be studied requires some level of preprocessing. Stop words (e.g., pronouns, prepositions, conjunctions) are expressions that have no distinguishing features and should be removed from the dataset prior to attempting summarization. For this reason, non-representative words in the sentences are removed from the original datasets in order to perform text summarization by transferring the data to the required format so that it can be studied. Thus, the processing load is reduced and continues with words having categorical meaning. In the DUC-2002 and DUC-2004 datasets used in the scope of the current study, the text to be expressed with graphs are stored in files with the .txt extension. In this study, the normalization steps, which refers to discarding stop words, spaces, and unwanted characters, was carried out using Python language library and tags.

**Table 1**
Text summarizing studies in the literature.

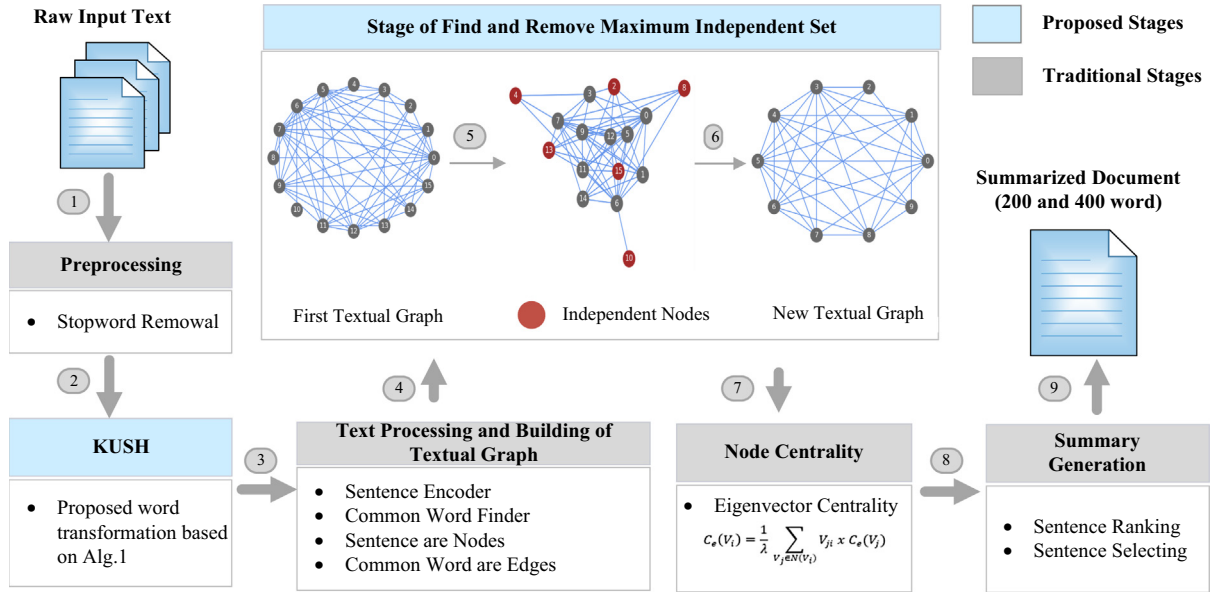| Authors, Publication Year | Summarization Type | Single/Multi | Model |
|---|---|---|---|
| [17] Luhn, 1958 | Extractive | Single | Statistical Method |
| [29] Salton et al., 1997 | Extractive | Multi | Hypertext Link Generation Algorithms |
| [13] Erkan et al., 2004 | Extractive | Multi | Graph-based Method & Eigenvector Centrality |
| [36] Mihalcea et al., 2005 | Extractive & Generic | Single | Graph-based Method & PageRank |
| [30] Medelyan, 2007 | Keyphrases Extraction | Single | Graph-based Method & Lexical Chains |
| [37] Fattah et al., 2009 | Extractive | Multi | Ga, Mr, Ffnn, Pnn, Gmm |
| [38] Shardan et al., 2010 | Extractive | Single | Semantic Nets, Fuzzy Logic & Evolutionary Programming |
| [31] Chen et al., 2011 | Key Term Extraction | Single | Graph-based Method & Random Walk |
| [32] Plaza et al., 2012 | Extractive | Multi | Graph-based Method & Umls |
| [23] Abuobieda et al., 2012 | Extractive | Single | Feature Selection & Genetic Concepts |
| [34] Moawad et al., 2012 | Abstractive | Single | Semantic Graph Reducing |
| [39] Gupta, 2013 | Extractive | Multi | Machine Learning |
| [28] Canhasi et al., 2014 | Query-Focused | Multi | Graph-based Method & Matrix Factorization |
| [35] Linqing Liu et al., 2014 | Abstractive | Multi | Reinforcement Learning |
| [20] Student et al., 2015 | Extractive | Single | Sentence Extraction, GA & NN |
| [27] Parveen et al., 2015 | Extractive | Single | Graph-based Method & ILP |
| [33] Nallapati, et al., 2016 | Abstractive | Multi | Recurrent Neural Networks |
| [12] Xiong & Ji, 2016 | Query-Focused | Multi | Hyper Graph-based Ranking |
| [19] Nasr Azadani et al., 2018 | Extractive | Single | Graph Clustering & Frequent Itemset |

**Fig. 1.** Schematic outline of proposed document summarization model.

These normalization steps constitute the first part of the data pre-processing and preparation stage.

In addition, when forming graphs of close-meaning words that differ from each other in terms of their spelling but semantically derived from a common word root, they are processed as if they are different words, which makes it difficult to identify connections and relationships between the sentences. It is predicted that this problem will significantly affect summarization performance, hence a text processing tool was developed within the framework of the proposed model. This software tool, which we named KUSH, was developed using C# on the .NET platform. The steps for the proposed KUSH preprocessing tool are given in Algorithm 1, and the pseudocode of the Best Alternative Search function is given in Algorithm 2.

**Algorithm 1.** Proposed KUSH algorithm

```
Input: {text} - texts in daily language
Output: {kush_text} - texts processed by KUSH algorithm
1    "text read";
2    sentences = [], words = [], alternative = [];
3    best_alternative = 0, kush_text = {}
4    sentences [] = text.Split(.);
5    for c to len(sentences) do
6        sentences [c] = clear(c);
7    end for
8    words[] = sentences.Split(.);
9    for k to boyut(words) do
10       alternative [k] = alternative_search(words [k], words);
11       for u to alternative do
12           best_alternative = best_alternative_search
         (k, alternative);
13       end for
14       words [k] = best_alternative;
15   end for
16   for k to words do
17       kush_text+ = k+" ";
18   end for
19   return kush_text;
```

**Algorithm 2.** Pseudocode of best_alternative_search algorithm

```
1. input: Word Vectors taken from Raw Text
2. Output: Upgraded word vector with best alternatives
3. n:Controls how many first characters of the two words to be
   compared must be equal
4. for i:0 to length(WordVector) step 1 do:
5.    for j:0 to length(WordVector) step 1 do:
6.       if i.th word contains j.vector
#It is checked whether the word in index j is in any subsection
   of the word in index i
7.          if i.th sub n charter = j.th n character (default n = 2)
//At least n of the beginning characters of the two words to be
   compared must be the same.
//Thus, it is che;cked whether or not the words come from the
   same origin.
//Otherwise, a word can be found in the middle or end of
   another word and an incorrect change applied.
8.          Replace i.th value with j.th value
9.       end if
10.     end if
11.  end for j
```

### 3.2. Textual graph

Many different methods have been proposed to represent information [40,41]. In this study, weighted and non-directional graphs are used for information representation. Prior to creating the graphs, the data passes through a preprocessing and preparation phase. After data preprocessing and preparation, the proposed Document Summarization model creates graphs corresponding to the text in order to summarize the data in textual format. The operations performed in this section corresponds to the "Text Processing and Building of Textual Graph" stage as shown in Fig. 1. A simple text example is shown in Fig. 2 along with the corresponding graph created for that text. The Sentence-Word Graph is obtained when performing the conversion in which the nodes are represented by sentences and the edges by common word numbers. The relationship levels of the sentence pairs are revealed by
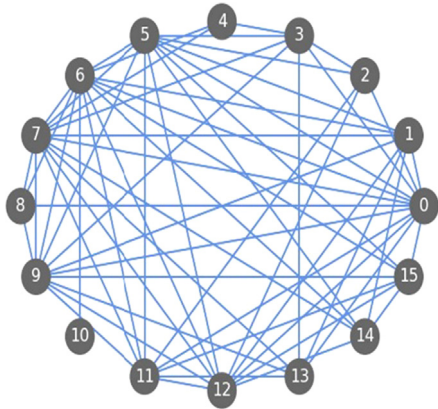
**Fig. 2.** Textual diagram of document d061 in DUC-2002 dataset after conversion with KUSH tool.

calculating the number of word intersections of each sentence with all other sentences. In this way, we can express the sentences and the relationships between them graphically with high level of representation.

A node is added to the representative graph for each sentence in the texts. For the edges between the nodes, the edge weight was added by considering the number of intersecting words of the phrases in the text. The KUSH algorithm regulates the semantic accuracy of the relationship between sentences, interconnecting the nodes based on meaningful relationships. Thus, all the sentences were associated with each other in terms of their common content that covers all combinations, ensuring that the relationships between the sentences can be accurately transferred to the graph. The steps for the Textual graph are given in Algorithm 3.

**Algorithm 3**. Textual Graph algorithm

---

1. Input: Text taken from preprocessing stage
2. Result: A is the Adjacency matrix
3. S_V = []#The input paragraph is separated to obtain sentence vectors(S_V).
4. A = [length(S_V), Length(S_V)] # Adjacency matrix is a square matrix in sentence vector dimension.
5. for i:0 to length(V) step 1 do:
6.   wordvector$_i$ []; #The word vector of $i$.th sentence is obtained by using the space character between the words
7.     for j:1 to length(V) step 1 do:
8.         wordvector$_j$ [];#The word vector of $j$.th sentence is obtained by using the space character between the words
9.           if count(wordvectori_i ∩ wordvectori_j > 0) #If the number of common words between sentences i and j.th is greater than zero, this value is assigned to the adjacency matrix as edge weight.
10.               A[i,j] = count(wordvectori_i ∩ wordvectori_j)
11.             Else : A[i,j] = 0 # otherwise, the adjacency matrix is assigned a value of zero.
12.               end if
13.       end for j
14. Graph is created using the Adjacency matrix A

---

### 3.3. Graph independent set

Many problems in graph theory are based on the existence of sub-diagrams that correspond to certain constraints. One of the most emphasized problems is the maximum independent set of a graph. In a G Graph, the required independent vertex sets are indicated by S, which contains the vertex values of V(G). There is no direct connection between the vertex values in the S set on the G graph. Similarly, a maximum independent set is not a suitable subset of another independent set on G graph. In other words, each vertex in an independent set S has at least one endpoint not in S, and each vertex not in S is adjacent to at least one vertex in S [42].

Erdös and Moser first raised the problem of determining the maximum number of maximum independent set values, n vertices and the graphs of these vertices [43]. After this problem was raised, various approaches have since been presented. In their study, Chan and Har-Peled presented Approximation Algorithms to find the maximum weighted maximum independent set [44].

Description: A Graph $G = (V; E)$ is a set of vertices (V) together with a set of edges (E). It is thought that the complementary $\overline{E} = \{(i, j) \notin E; i, j \in V, i \neq j\}$ graph of this G graph $\overline{G} = (V, \overline{E})$ meets the conditions. An independent set is a set of nodes in the graph. There is no direct neighborhood between the nodes in this set. If a $S \subset V$ subset is an independent set, $\forall i, j \in S$ and the edges of these nodes should meet the conditions of $(i, j) \notin E$. A maximum independent set is an independent maximum cardinality set. Finding a maximum independent set is equivalent to finding a maximum clique or minimum vertex cover, and these three problems are classified as NP-hard problems [45].

Although Maximum Independent Set (MIS) has been used in many fields in graph theory, NP-Hard remains a current problem. An independent set is a set of nodes in a graph, such that no two vertices in the set are connected by an edge. That is, each edge in the undirected graph has at most one endpoint. There may be more than one maximal set in a graph. Adding another node to the maximal clusters disrupts the independent set's property as it would require that the set contains an edge. All maximal independent sets in a graph are determined and the maximal sets with the most nodes among the specified maximal independent sets are determined as the maximum independent set. For instance, in Fig. 3, {n2, n3} and {n2, n3, n5} are a maximal independent set, while {n2, n3, n5, n8} is the maximum independent set, and four is the independence number of the graph.

The current study aims to remove sentences which are not to be summarized, and to reduce the number of sentences to be worked on prior to the sentence selection. Separation of unnecessary and meaningless sentences is performed by using graph independent set. The document, for which a summary is to be created, is first converted into a graph and then an independent set of this graph is determined. The nodes in the independent set are removed from the main document and the remaining sentences are summarized using the node-centered values. Fig. 4 explains each step of the process.

As can be seen in Fig. 4, the conversion of a small document to graphs is followed by finding the independent sets and then removing them from the document. The number of sentences in the document and the number of nodes in the first graph created are assumed to be n. The number of independent set nodes obtained from the main graph is m; while the probability that the kth node in the main sentence is present in the summary without any scoring is $P_k = \left(\frac{1}{n}\right)$, it is $P_k = \left(\frac{1}{n-m}\right)$ after the independent sets are removed from the main sentence. Thus, before calculating node centralities, the probability of the presence of sentences, which are considered to be semantically strong in the paragraph, in the summary is increased to a certain extent. In this way, when the Eigen Centrality Values are obtained from the initial graph, the process is performed by considering more sentences and this increases the probability of the presence of sentences to be excluded from the summary. By removing the sentences in the independent set from the main document, the sentences, which have high eigenvector centrality values but should be excluded from the summary, are determined and thus more accurate summaries are produced.
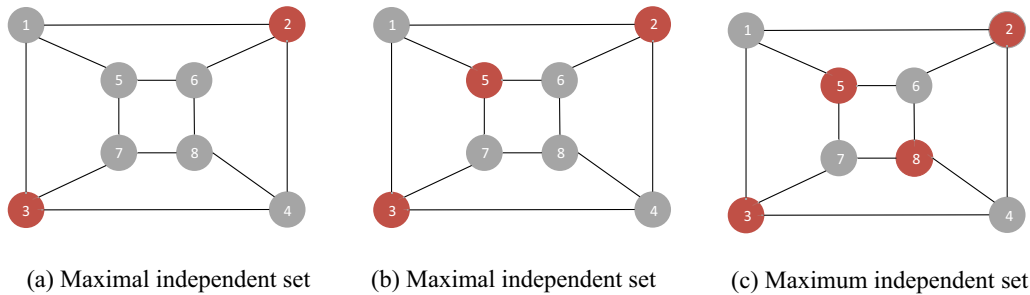
(a) Maximal independent set    (b) Maximal independent set    (c) Maximum independent set

**Fig. 3.** Illustration of difference between Maximal Independent Set and Maximum Independent Set.



**Find And Remove Unvaluable Sentences From Main Document**

First Textual Graph          Graph Independent Set          Last Textual Graph After Removing Independent Set From Main Document
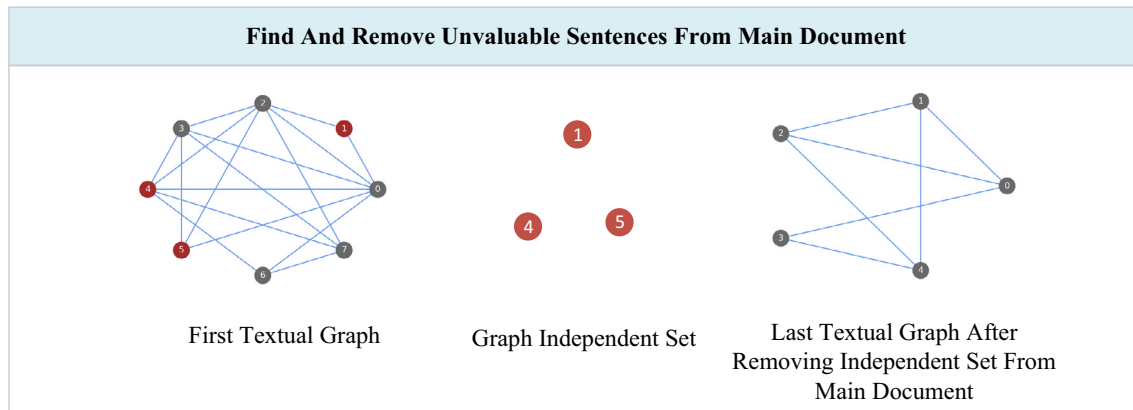
**Fig. 4.** Finding and removing the independent set from a simple graph.

The proposed method can be briefly described as follows. G graph is obtained from the given text (T) file. A graph is $G = (V, E)$; the maximum independent set of this graph is included in the method proposed in this study. $S \subset V$ is this set; since the maximum independent set is S, a graph is defined as $G2 = (V_2, E_2)$.

$T \rightarrow G$
$E_2 \subset E$ and $S \subset V$.
$E_2 = \{(u,v)|u \in V\text{-}S$ and $v \in V\text{-}S\}$ and $V_2 = V\text{-}S$.
$G_2 \rightarrow L$ Matrix (L is Laplacian Matrix of $G_2$).
$L \rightarrow$ eigenvector and eigenvalue
eigenvector and eigenvalue $\rightarrow$ Summarization.

Finding the maximum independent set is a problem that is in need of a reliable approach, but a definitive and effective method has yet to be developed. Algorithm 4 was used in the current study to find the maximum independent set.

**Algorithm 4**. Maximum Independent Set Algorithm

---

Input: Textual Graph (G) (Section 3.2)
Output: Maximum Independent Set (S)
1. V: Set of Vertices
2. E: Set of Edges
3. G ← (V,E)
4. S: Independent set of G and Initial value S = $\varnothing$
5. if $\forall$v $\in$ and degree(v) = 0 then result S = V.
6. Other cases;
7. u $\in$ V select node.
8. n1 = Maximum_Independent_Set((G-{u})(See Alg. 5.)
9. n2 = Maximum_Independent_Set((G-{u}-{neighbors (u)})
10. S = maximum(n1,n2)
11. Return S

---

The pseudocode of the Maximum Independent Set algorithm can be summarized as:

**Algorithm 5**. Maximum_Independent_Set Functions Algorithm

---

Maximum_Independent_Set(G)

---

1. Start with the set of vertices of the graph, V and an empty set for the MIS, S
2. While V$\neq\varnothing$:
3.1. Find a vertex of minimum degree v $\in$ V
3.2. Add it to S
3.3. Remove it and its neighbors from V
4. End While
5. Return S

---

*3.4. T.GISEC method: proposed document summarization algorithm*

The proposed method obtains the Graph from the Text. The independent set is obtained from the graph and, after the independent set is removed, the eigenvalue is obtained and summarized with eigenvalue centrality. In the proposed summarization approach, based on the assumption that the sentences that should be excluded from the summary will be represented as Independent Set, the mentioned nodes are determined and removed from the representative graph. Thus, in the experimental process leading to the summary, a method which has not previously been used in any document summarization study is presented. In order to find effective nodes, ineffective nodes were first identified. As shown in Algorithm 6, the remaining nodes after the ineffective nodes have been removed, eigenvector node centrality is weighted with [5] metric. For the obtained Top *N* node scores, 200- and 400-word summaries were obtained.

**Algorithm 6**. Proposed Document Summarization algorithm

---

**Input:** Processed Text From Kush Preprocessing Step
(see Section 3.1)
**Output:** Summary

---

1    V: Set of Vertices
2    E: Set of Edges
3    $G \leftarrow (V,E)$
4    **I ← Find All Independent Set of G Graphs**
5    **for each** $I_i$ **in** $I$
6       Remove $i_i^{th}$ sentences from source document (see Algorithm 2)
7    **End For**
8    Calculate New Vertex and New Edges
9    Create New Textual Graph $G_{new} \leftarrow (V_{new}, E_{new})$ (see Section 3.2)
10   Calculate Node Centrality
11   Sort all Node Centrality Values of $G_{new}$ by Descending
12   Create summary of (100, 200 and 400 Words), beginning from the highest Eigen Centrality Value
13   **Return** Summary

---

## 4. Experimental results

### 4.1. Description of dataset

In this study, two Document Understanding Conference datasets (DUC-2002 and DUC-2004) were used to test the accuracy of the proposed method. The DUC-2002 dataset contains documents for both abstractive and extractive summarization; however, the extractive summarization files were used based on the summarization method proposed in the current study. NIST produced 60 reference sets, with each containing documents, single-document summaries, and multi-document abstracts/extracts, and defined criteria such as event sets and biographical sets.

For the DUC-2002 dataset, three tasks were defined. Task 1, fully automated summarization of multiple newswire/newspaper documents (articles) on a single subject, with 60 sets of approximately 10 documents each provided as system input. Task 2, automated summarization of multiple newswire/newspaper documents (articles) on a single subject. Task 3, one or two pilot projects with extrinsic evaluation.

For the DUC-2004 dataset, five tasks were defined. Task 1, to create very short summaries with a maximum length of 75 characters for 500 newswire/newspaper documents (articles); these summaries can be construed as headlines, although participants were allowed to use any format (including keyword lists). Task 2, to produce summaries with a maximum length of 665 characters for 50 clusters of 10 documents each; these summaries were general and not focused on any particular aspect of the documents. Task 3, similar to Task 1, except that the document set consisted of 24 clusters of 10 documents each. Task 4, same document set as Task 3, except with general multi-document summaries for each of the 24 clusters. Task 5, similar to Task 2, with summaries for 50 clusters of 10 documents. In this study we selected Task 2 for 665 bytes summaries [46]. Characteristics of the two datasets are presented in Table 2.

### 4.2. Evaluation metrics

Recall-Oriented Understudy for Gisting Evaluation (ROUGE) performance metrics are the most popular assessment metrics used in summarization systems. ROUGE is a performance metric in which summaries are automatically evaluated. These metrics

**Table 2**
Characteristics of the datasets.

| Description | DUC-2002 | DUC-2004 |
|---|---|---|
| Number of clusters | 59 | 50 |
| Number of documents in each cluster | ~10 | 10 |
| Number of documents | 567 | 500 |
| Summary length | 200 and 400 words | 665 bytes |

are evaluated on the basis of n-gram, word sequences. It is a metric based on the intersections between summaries created by summarization systems and ideal summaries created by human hand. In order to evaluate the performance of the proposed framework, the ROUGE assessment toolkit was utilized, which is based on n-gram statistics and highly correlated to human evaluation. The scores generated by the ROUGE metrics range from 0 to 1. The higher the score, the more shared content is available with the model summary; and for the proposed method, the automated summary was considered to be better and more informative. In Refs. [47,48], Lin revealed a high correlation between ROUGE scores and the scores given by individuals. In the current study, we used ROUGE-N (N-1, N-2), ROUGE L, ROUGE-W-1.2 and ROUGE-SU metrics to evaluate the performance of the proposed Document Summarization method.

ROUGE-N evaluates the number of n-grams common to the proposed summary and the model summary.

$$ROUGE - N = \frac{\sum_{C \in \{ReferenceSummaries\}gram_n \in S} \sum Count_{match}(gram_n)}{\sum_{C \in \{ReferenceSummaries\}gram_n \in S} \sum Count(gram_n)} \quad (1)$$

where $n$ equals the length of $gram_n$ n-gram, and $Count_{match}(gram_n)$ is the maximum number of $n$-grams intersecting in the potential summary and reference summary. As can clearly be seen in Eq. (1), Rouge-N is a measurement associated with recall. This is because the denominator of the equation is the sum of the number of n-grams in the reference summary. For example, (N-1) measures the number of uni-grams shared between two summaries. Similarly (N-2) focuses on the number of bi-grams shared between the proposed summary and the model summary. Similarly, it calculates the longest common word sequence in the two sub-word sequences given in the ROUGE-L value. X and Y are two given word sequences. Calculations of the ROUGE-L value of the sequences are made in Eqs. (2)–(4).

$$R_{lcs} = \frac{LCS(X, Y)}{m} \quad (2)$$

$$P_{lcs} = \frac{LCS(X, Y)}{n} \quad (3)$$

$$F_{lcs} = \frac{(1 + \beta^2)R_{lcs}P_{lcs}}{R_{lcs} + \beta^2 P_{lcs}} \quad (4)$$

ROUGE-W-1.2 calculates the longest consecutive matches between the proposed summary and the model summary. ROUGE-SU4 measures the number of skip-bigrams common to the two summaries.

### 4.3. Experimental works

All experimental processes were performed using a computer with an Intel Core i7-7700 CPU 3.60 GHz and 16 Gb memory using .net and Python.

In the DUC-2002 dataset, the texts are stored in files with the .txt extension. The normalization steps, meaning the removal of non-discriminatory words, expressions, spaces and undesired characters called stop words, were performed using the Python

language library and tags. Thus, the processing load was reduced and the Document Summarization process was conducted with words having categorical meaning.

After the preprocessing stage, the developed software tool called KUSH was used to provide the most accurate transfer of relationships between word phrases and textual graphs. Owing to its simple and efficient algorithm, the software assigned substitutes, chosen from the text to be summarized, for the changed words. As a result, graphs with high levels of representation were obtained.

Maximum Independent Set, which has not previously been used in any summarization study, was utilized for the current study. Based on the assumption that sentences corresponding to the nodes in the independent set should be excluded from the summary, the nodes forming the Independent Sets on the graphs were determined and removed from the graph. Thus, prior to the effect of the nodes on the global graph being quantified, a limitation was applied to the summary. This limitation prevented the repetition of word groups being included in the summary, thereby resulting in more comprehensive summaries being generated. In addition, the experimental processes affected the approach for removing the Independent Sets, which has been used for the first time in a summarization study, as an encouraging step to include sentences with minimum words that intersect with each other in the main text. The values reported throughout the experimental processes of the study clearly demonstrate the contributions of this innovative method.

In the final stage, the centrality values of the obtained nodes were calculated with Node centrality calculations. Eigenvector Centrality was used to weight the nodes forming the graphs in the study. In order to thoroughly evaluate the performance of the proposed summarization system, the experimental process was conducted separately for summaries consisting of 200 and 400 words. The performance of the proposed Document Summarization framework was evaluated using the ROUGE performance metrics, as shown in Table 3.

The metrics used were Rouge-1, Rouge-2, Rouge-3, Rouge-4, Rouge-L, Rouge-W-1-2, Rouge-S*, and Rouge-SU*. In addition, each metric type has been reported separately with a focus on Recall, Precision, and F-Score. For summaries of 200 and 400 words, the values have been presented separately. The first column of Table 3 shows the Summarization performance metric types, while the other columns show performance evaluation scores with an average reported by the recommended framework for 200- and 400-word summaries.

### 4.4. Comparison with state-of-the-art methods

The Summarization performance results obtained with the proposed Document Summarization framework were compared with previously published studies. As shown in Tables 4–6, 100-, 200- and 400-word summaries were compared with seven different competitive methods in order to show consistency with previous summarization approaches.

Luhn [17] used statistical information obtained from word frequencies and distributions to calculate the significance of sentences with machine learning techniques, forming summaries according to the highest sentence scores obtained. Landauer et al. [49,50] presented a new approach using only general mathematical methods, without utilizing any knowledge pre-acceptance. By transferring the connected structure of a text to graphs, Mihalcea realized an iterative, extractive and unsupervised application called TextRank [51,52] that scores units based on the significance of text units. An important feature of the TextRank system, which is an iterative study based on Google's PageRank [26], is that there is no predefined and manually configured structure. In Ref. [13], Erkan et al. proposed a stochastic and graph-based approach called LexRank to calculate the significance of textual units within NLP. In their approach, they calculated the significance of sentences on the representative graph based on the eigenvector centrality (node centrality-based) measure. In their experimental work, the authors showed that LexRank performed better in most cases than both degree-based methods and centroid-based methods. SumBasic [53], a fundamental method frequently used in the literature, selects the top *N* sentence to be included in the summary using an effective probability function. In Ref. [54], the authors added sentences to the summary with a specified probability model, in a greedy way for multi-document summarization.

The reason for choosing the aforementioned methods is that they use mathematical, statistical or graph-based theories, and are thereby similar to the summarization method presented in the current study. In addition, all of the selected methods are unsupervised document Summarization methods, which is also the same as in the proposed framework.

Table 4 presents the ROUGE performance metric values of the 200-word summaries reported by Random, Luhn, LSA, TextRank, LexRank, SumBasic, KL-Sum and the document summary method proposed in the current study on the DUC-2002 dataset. The highest values for each row in the table are highlighted in bold font. The proposed method is shown to have outperformed all of the competitive methods, without exception. Table 4 clearly demonstrates that the proposed approach produced excellent results when compared to the competitive methods in terms of 200-word summaries.

Fig. 5 displays the ROUGE-1, ROUGE-2, ROUGE-L, ROUGE-W-1.2, and ROUGE-SU4 average scores for all eight approaches, including the proposed method. All of the comparisons in the graph are summarized separately for the 200-word summaries based on the Recall, Precision and F-Score values. The proposed method, based on a Rouge-1 and F-Score assessment, reported a 25% higher value than Random, 8% more than Luhn, 33% more than LSA, 8% more than TextRank, 5% more than LexRank, 9% more than
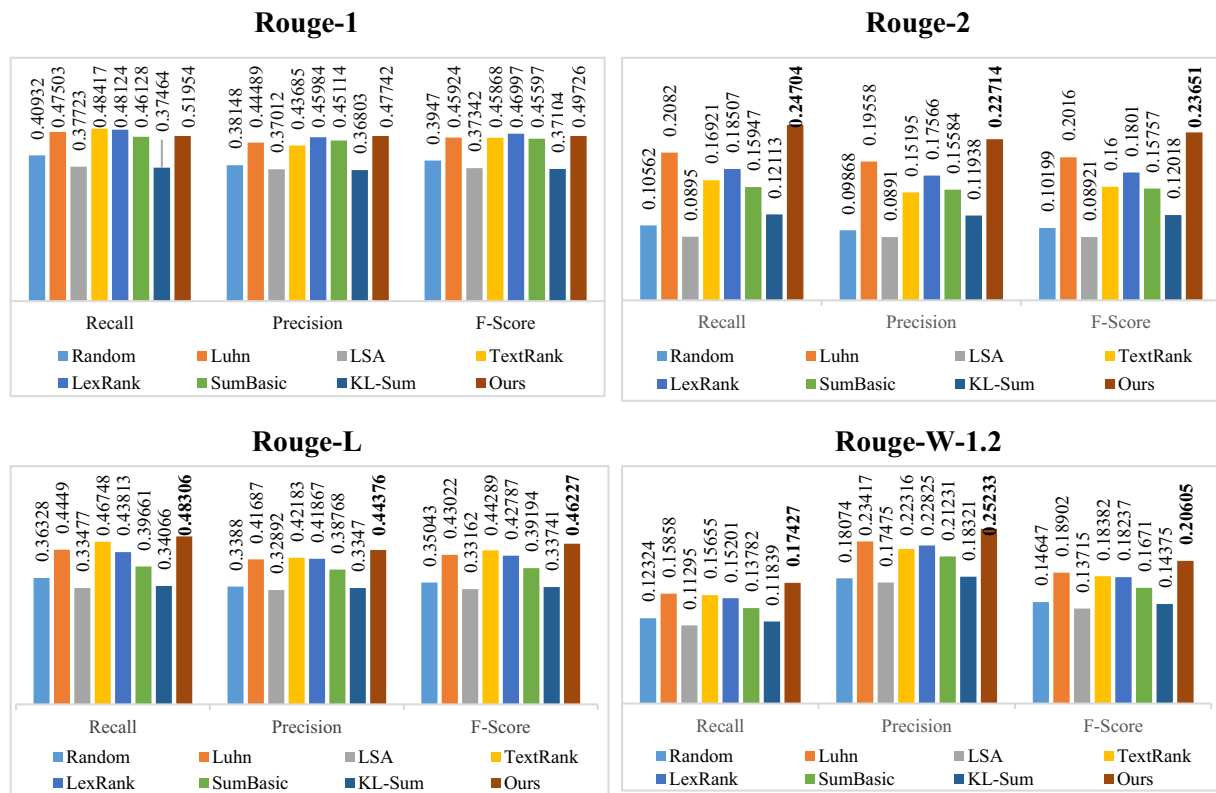
**Table 3**
ROUGE metric scores for the proposed document summarization method.

| ROUGE evaluation methods | Average | | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | Summary for 200-words | | | Summary for 400-words | | |
| | Recall | Precision | F-Score | Recall | Precision | F-Score |
| Rouge-1 | 0.51954 | 0.47742 | 0.49726 | 0.59208 | 0.56983 | 0.58060 |
| Rouge-2 | 0.24704 | 0.22714 | 0.23651 | 0.32471 | 0.31207 | 0.31819 |
| Rouge-3 | 0.18677 | 0.17210 | 0.17900 | 0.26191 | 0.25137 | 0.25647 |
| Rouge-4 | 0.17145 | 0.15794 | 0.16430 | 0.24083 | 0.23096 | 0.23574 |
| Rouge-L | 0.48306 | 0.44376 | 0.46227 | 0.56361 | 0.54225 | 0.55258 |
| Rouge-W-1.2 | 0.17427 | 0.25233 | 0.20605 | 0.17856 | 0.27067 | 0.21510 |
| Rouge-S* | 0.23114 | 0.19577 | 0.21143 | 0.32032 | 0.29671 | 0.30775 |
| Rouge-SU* | 0.23388 | 0.19823 | 0.21403 | 0.32163 | 0.29798 | 0.30905 |

**Table 4**
Comparison with state-of-the-art methods of proposed document Summarization (200-word summaries).

| ROUGE evaluation methods | | Summary for 200-words (average) (DUC-2002) | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Methods | | | | | | | |
| | | Random | Luhn [17] | LSA [49,50] | TextRank [52,51] | LexRank [13] | SumBasic [53] | KLSum [54] | Proposed method |
| Rouge-1 | Rec. | 0.40932 | 0.47503 | 0.37723 | 0.48417 | 0.48124 | 0.46128 | 0.37464 | 0.51954 |
| | Prec. | 0.38148 | 0.44489 | 0.37012 | 0.43685 | 0.45984 | 0.45114 | 0.36803 | 0.47742 |
| | F-Sco. | 0.39470 | 0.45924 | 0.37342 | 0.45868 | 0.46997 | 0.45597 | 0.37104 | 0.49726 |
| Rouge-2 | Rec. | 0.10562 | 0.20820 | 0.08950 | 0.16921 | 0.18507 | 0.15947 | 0.12113 | 0.24704 |
| | Prec. | 0.09868 | 0.19558 | 0.08910 | 0.15195 | 0.17566 | 0.15584 | 0.11938 | 0.22714 |
| | F-Sco. | 0.10199 | 0.20160 | 0.08921 | 0.16000 | 0.18010 | 0.15757 | 0.12018 | 0.23651 |
| Rouge-L | Rec. | 0.36328 | 0.44490 | 0.33477 | 0.46748 | 0.43813 | 0.39661 | 0.34066 | 0.48306 |
| | Prec. | 0.33880 | 0.41687 | 0.32892 | 0.42183 | 0.41867 | 0.38768 | 0.33470 | 0.44376 |
| | F-Sco. | 0.35043 | 0.43022 | 0.33162 | 0.44289 | 0.42787 | 0.39194 | 0.33741 | 0.46227 |
| Rouge-W-1.2 | Rec. | 0.12324 | 0.15858 | 0.11295 | 0.15655 | 0.15201 | 0.13782 | 0.11839 | 0.17427 |
| | Prec. | 0.18074 | 0.23417 | 0.17475 | 0.22316 | 0.22825 | 0.21231 | 0.18321 | 0.25233 |
| | F-Sco. | 0.14647 | 0.18902 | 0.13715 | 0.18382 | 0.18237 | 0.16710 | 0.14375 | 0.20605 |
| Rouge-SU* | Rec. | 0.15112 | 0.19901 | 0.12783 | 0.21541 | 0.20321 | 0.17981 | 0.13361 | 0.23388 |
| | Prec. | 0.13124 | 0.17444 | 0.12342 | 0.17468 | 0.18494 | 0.17190 | 0.12913 | 0.19823 |
| | F-Sco. | 0.14020 | 0.18556 | 0.12527 | 0.19200 | 0.19309 | 0.17548 | 0.13098 | 0.21403 |



**Fig. 5.** Graphical comparison of proposed method with state-of-the-art methods (200-word summaries).

SumBasic, and 34% more than KL-Summ summarization methods. Similar evaluation rates were obtained for all other ROUGE performance measurement methods, as can be seen in the relevant tables and figures that follow. This confirms that the proposed document summarization approach can be considered as a rather effective and efficient summarization system.

To emphasize the applicability of the proposed method, the experimental processes were repeated for 400-word summaries; again, by examining in detail all the competitive methods compared to the performance of the document summarization approach proposed in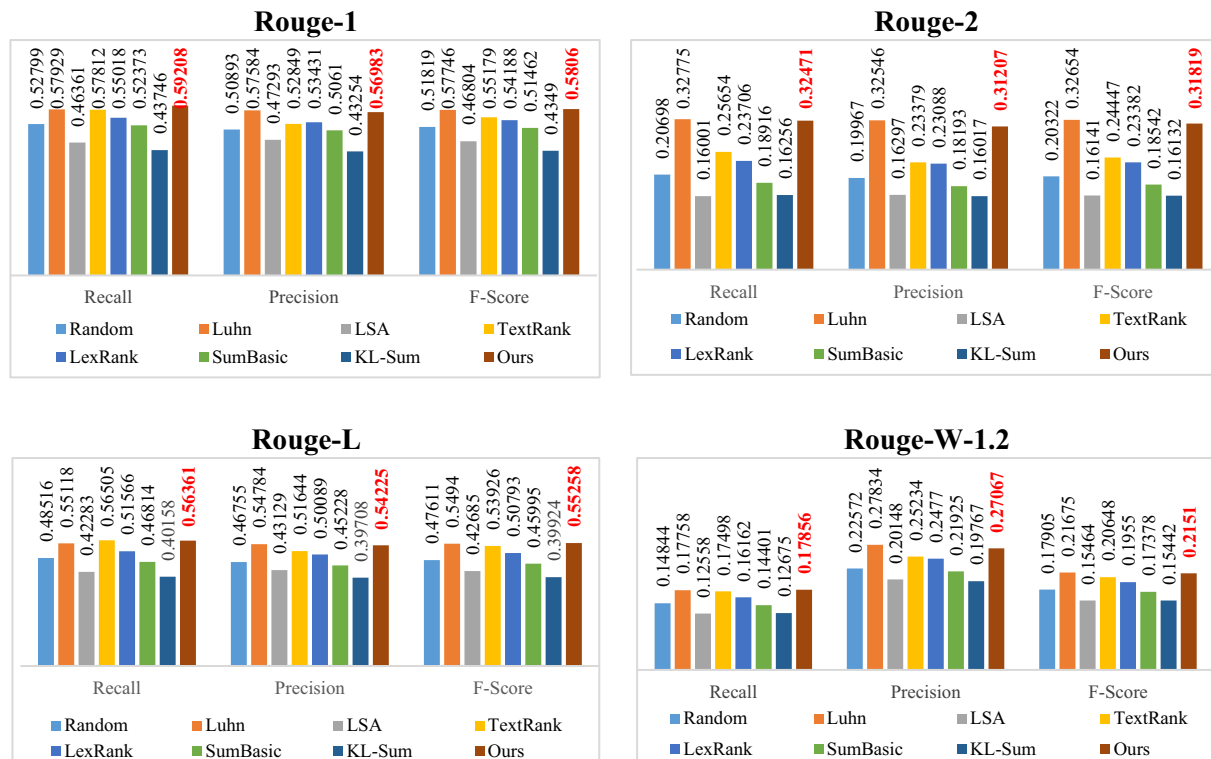 the current study. The results of this comparison are presented in Table 5, with the highest values in each row highlighted in bold font. The proposed method, based on a Rouge-1 and F-Score assessment, reported a 12% higher value than Random, 0.5% more than Luhn, 24% more than LSA, 5% more than TextRank, 7% more than LexRank, 12% more than SumBasic, and 33% more than KL-Summ summarization methods. It is clear, therefore, that many of the competitive methods have been left behind by the proposed method on the basis of 400-word summaries.

Fig. 6 illustrates the ROUGE-1, ROUGE-2, ROUGE-L, ROUGE-W-1.2, and ROUGE-SU4 average scores for all eight approaches, including the proposed method. All of the comparisons in the

**Table 5**
Comparison with state-of-the-art methods of proposed document Summarization (400 word summaries).

| ROUGE evaluation methods | | Summary for 400-words (DUC-2002) | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Methods | | | | | | | |
| | | Random | Luhn [17] | LSA [49,50] | TextRank [52,51] | LexRank [13] | SumBasic [53] | KLSum [54] | Proposed method |
| Rouge-1 | Rec. | 0.52799 | 0.57929 | 0.46361 | 0.57812 | 0.55018 | 0.52373 | 0.43746 | 0.59208 |
| | Prec. | 0.50893 | 0.57584 | 0.47293 | 0.52849 | 0.53431 | 0.50610 | 0.43254 | 0.56983 |
| | F-Sco. | 0.51819 | 0.57746 | 0.46804 | 0.55179 | 0.54188 | 0.51462 | 0.43490 | 0.58060 |
| Rouge-2 | Rec. | 0.20698 | 0.32775 | 0.16001 | 0.25654 | 0.23706 | 0.18916 | 0.16256 | 0.32471 |
| | Prec. | 0.19967 | 0.32546 | 0.16297 | 0.23379 | 0.23088 | 0.18193 | 0.16017 | 0.31207 |
| | F-Sco. | 0.20322 | 0.32654 | 0.16141 | 0.24447 | 0.23382 | 0.18542 | 0.16132 | 0.31819 |
| Rouge-L | Rec. | 0.48516 | 0.55118 | 0.42283 | 0.56505 | 0.51566 | 0.46814 | 0.40158 | 0.56361 |
| | Prec. | 0.46755 | 0.54784 | 0.43129 | 0.51644 | 0.50089 | 0.45228 | 0.39708 | 0.54225 |
| | F-Sco. | 0.47611 | 0.54940 | 0.42685 | 0.53926 | 0.50793 | 0.45995 | 0.39924 | 0.55258 |
| Rouge-W-1.2 | Rec. | 0.14844 | 0.17758 | 0.12558 | 0.17498 | 0.16162 | 0.14401 | 0.12675 | 0.17856 |
| | Prec. | 0.22572 | 0.27834 | 0.20148 | 0.25234 | 0.24770 | 0.21925 | 0.19767 | 0.27067 |
| | F-Sco. | 0.17905 | 0.21675 | 0.15464 | 0.20648 | 0.19550 | 0.17378 | 0.15442 | 0.21510 |
| Rouge-SU* | Rec. | 0.25020 | 0.29304 | 0.19823 | 0.31213 | 0.27351 | 0.24097 | 0.17922 | 0.32163 |
| | Prec. | 0.23189 | 0.28977 | 0.20537 | 0.26032 | 0.25871 | 0.22412 | 0.17485 | 0.29798 |
| | F-Sco. | 0.24052 | 0.29118 | 0.20140 | 0.28307 | 0.26541 | 0.23197 | 0.17688 | 0.30905 |



**Fig. 6.** Graphical comparison of proposed method the state-of-the-art methods (400-word summaries).

graph were created separately for 400-word summaries based on Recall, Precision, and F-Score values.

In our experiments conducted with 100-word summaries from the DUC-2004 dataset, it can be seen in Table 6 that the proposed algorithm surpassed other competitive methods for all performance criterions based on ROUGE-2 values. Best values were also reported for the ROUGE-1, ROUGE-L, ROUGE-SU, and the ROUGE-W-1.2 precision metric types. The highest values for each row in the tables are highlighted in bold font. Also, Fig. 7 displays the ROUGE-1, ROUGE-2, ROUGE-L, ROUGE-W-1.2, and ROUGE-SU4 average scores for all eight systems, including the proposed method. All of the comparisons in the graph were created
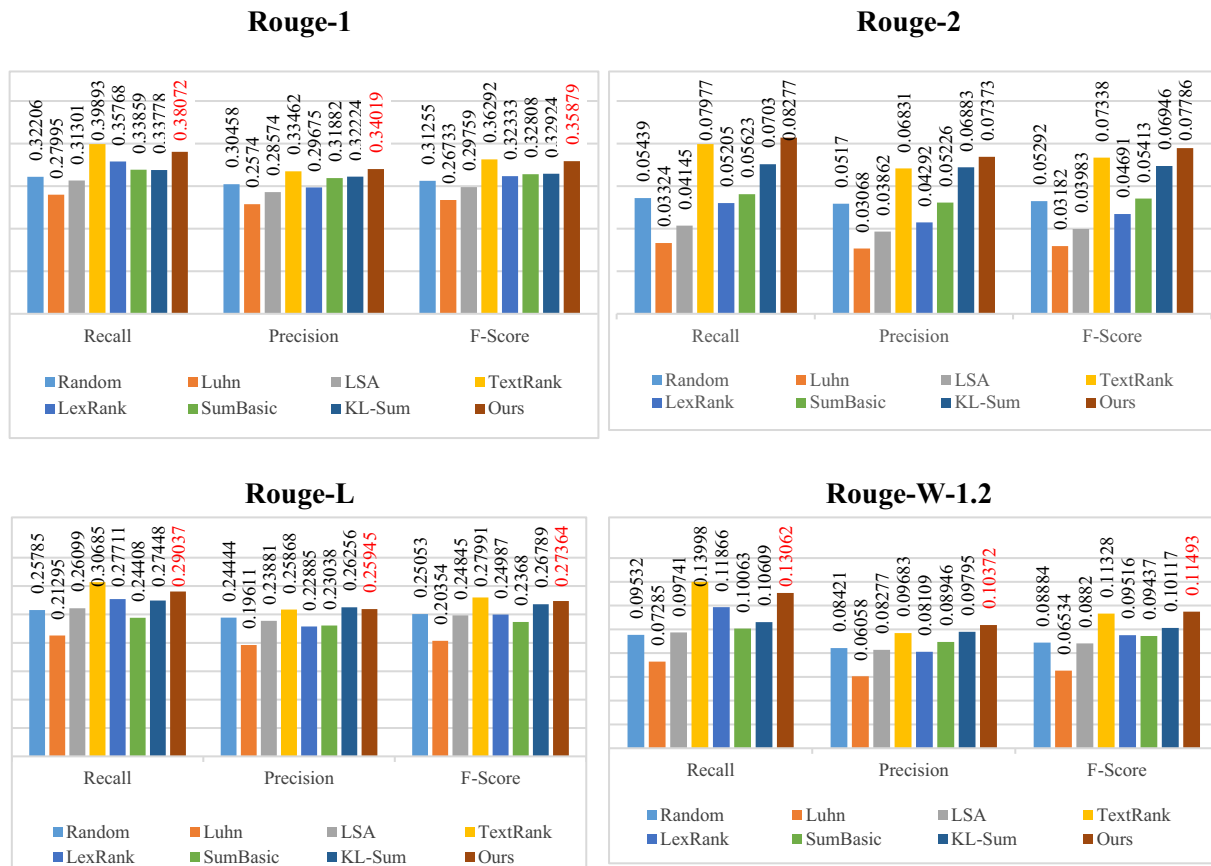
separately for the 100-word summaries based on the Recall, Precision, and F-Score values.

The experimental study was examined on the basis of 100-, 200- and 400-word summaries on the DUC-2002 and DUC-2004 datasets for all approaches. The selected state-of-the-art methods and the proposed Document Summarization framework are similar as in they are all mathematical-statistical and graph-based approaches. When Tables 4 and 5 are examined, the proposed approach outperformed all of the competitive methods and reported excellent results on the basis of 200-word summaries. The proposed method also offers the best ROUGE values alongside Luhn on the basis of 400-word summaries.

**Table 6**
Comparison with state-of-the-art methods of proposed document Summarization (100 word summaries).

| ROUGE evaluation method | | Summary for 100-words (DUC-2004) | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Method | | | | | | | |
| | | Random | Luhn [17] | LSA [49,50] | TextRank [52,51] | LexRank [13] | SumBasic [53] | KLSum [54] | Proposed method |
| Rouge-1 | Rec. | 0.35768 | 0.27995 | 0.31301 | 0.39893 | 0.32206 | 0.33859 | 0.33778 | 0.38072 |
| | Prec. | 0.29675 | 0.25740 | 0.28574 | 0.33462 | 0.30458 | 0.31882 | 0.32224 | 0.34019 |
| | F-Sco. | 0.32333 | 0.26733 | 0.29759 | 0.36292 | 0.31255 | 0.32808 | 0.32924 | 0.35879 |
| Rouge-2 | Rec. | 0.05205 | 0.03324 | 0.04145 | 0.07977 | 0.05439 | 0.05623 | 0.07030 | 0.08277 |
| | Prec. | 0.04292 | 0.03068 | 0.03862 | 0.06831 | 0.05170 | 0.05226 | 0.06883 | 0.07373 |
| | F-Sco. | 0.04691 | 0.03182 | 0.03983 | 0.07338 | 0.05292 | 0.05413 | 0.06946 | 0.07786 |
| Rouge-L | Rec. | 0.27711 | 0.21295 | 0.26099 | 0.30685 | 0.25785 | 0.24408 | 0.27448 | 0.29037 |
| | Prec. | 0.22885 | 0.19611 | 0.23881 | 0.25868 | 0.24444 | 0.23038 | 0.26256 | 0.25945 |
| | F-Sco. | 0.24987 | 0.20354 | 0.24845 | 0.27991 | 0.25053 | 0.23680 | 0.26789 | 0.27364 |
| Rouge-W-1.2 | Rec. | 0.09350 | 0.07427 | 0.08836 | 0.10436 | 0.08957 | 0.08427 | 0.09315 | 0.09908 |
| | Prec. | 0.14127 | 0.12498 | 0.14797 | 0.16113 | 0.15517 | 0.14533 | 0.16340 | 0.16197 |
| | F-Sco. | 0.11212 | 0.09286 | 0.11011 | 0.12622 | 0.11334 | 0.10656 | 0.11841 | 0.12278 |
| Rouge-SU | Rec. | 0.11866 | 0.07285 | 0.09741 | 0.13998 | 0.09532 | 0.10063 | 0.10609 | 0.13062 |
| | Prec. | 0.08109 | 0.06058 | 0.08277 | 0.09683 | 0.08421 | 0.08946 | 0.09795 | 0.10372 |
| | F-Sco. | 0.09516 | 0.06534 | 0.08820 | 0.11328 | 0.08884 | 0.09437 | 0.10117 | 0.11493 |



**Fig. 7.** Graphical comparison of proposed method with state-of-the-art methods (100-word summaries).

The performance of the proposed system was tested in detail using different datasets and performance metrics. As has been reported in the presented tables and figures, the proposed method produced generally successful results. This outcome is partly related to the creation of rich textual charts using the proposed text processing software (named as KUSH). However, the primary reason for the success is the subtraction of nodes that make up the independent sets prior to the node-weighting phase. Thus, the selection of appropriate sentences to be selected for summarization is more accurate as a result. This provides for graph-based preprocessing, allowing for work on simpler graphs, which in turn improves coverage, which is an important feature of the rich summary.

The comparative results show that the unique and innovative Document Summarization framework based on the graph independent set produces superior results to the previous studies when tested against the DUC datasets.

## 5. Conclusion

In order to lessen the time necessary to access today's ever-increasing vast capacity of information, data needs to be analyzed in order to obtain a number of properties. This challenging task has naturally increased the interest in automated summarization systems. The current study therefore proposes a new, unsupervised, extractive, and generic text summarization system.

In the proposed approach, text irregularities brought about by spoken language have been eliminated using a software tool called KUSH that was developed in this study. The software tool prepares texts for summarizing using a unique algorithm, and has been shown to positively affect results in terms of revealing semantically distinguishable and measurable relationships between sentences.

Sentences represent the first finite set in representative graphs, while the number of common words represents the other finite set. In this way, texts are represented by nonoriented and weighted graphs.

As researchers, our anticipation was that the set of sentences corresponding to the nodes in the independent set should be excluded from the summary. Based on this, the nodes forming the Independent Set on the graphs were identified and removed from the graph. Thus, prior to the quantification of the effect of the nodes on the global graph, a limitation was applied to the summary. Within the scope of the proposed Document Summarization process, the nodes forming Independent Sets were determined on the graphs representing the texts to be summarized. A preliminary assessment was performed between the nodes prior to the quantitative information being obtained from the nodes forming the textual graph. With the proposed model, a summarization approach was established with an infrastructure created based on a mathematical model, which has not previously been seen in any existing summarization study. The effectiveness of the model has been demonstrated in terms of Recall, Precision, and F-Scores using the Rouge-1, Rouge-2, Rouge-3, Rouge-4, Rouge-L, Rouge-W-1-2, Rouge-S*, and Rouge-SU* performance metrics. Experimental processes were also repeated for 100-, 200- and 400-word summaries. The proposed system's success was verified against two datasets, with the following results obtained during the experimental processes.

With a 95% confidence interval, superior results were reported for the proposed summarization method over all existing state-of-the-art methods, without exception, in terms of 200-word summaries; and higher results obtained over most state-of-the-art methods in terms of 100- and 400-word summaries.

Prior to the weight of the nodes on the global graph being mathematically sorted, a unique approach is performed, which has not been used in any other document summarization study. The approach isolates the nodes forming independent nodes from the main graph and then creates the summary. This limitation prevents the repetition of word groups from being included in the summary, and thereby generating much more comprehensive summaries. In addition, the experimental processes have been used as a preliminary assessment step which encourages the removal of the Independent Sets, and first used in a summarization study. So, in the summaries, the sentences with maximum relations with each other are included.

The values reported throughout the experimental processes of the study demonstrate the contribution of this innovative method. The aforementioned stage of the presented method has been shown to be very effective and consistent, thanks to its mathematical grounding. In addition, the proposed method depends on the correct formation of the graphs, which is still at an early stage of development. The method presented includes a robust fiction. However, the text requires a preprocessing step and a correct graph transformation, which affects the system's performance considerably. While maximum independent sets of charts are found, the size of the chart used is decisive on the system's performance, and as such, where there are very large graphs, there may be issues related to insufficient memory. However, in the proposed method, no memory problems were experienced due to the dataset dimensions used during the experimental processes (average of 200 sentences). It is one of our goals to develop the proposed method by testing its success on different graph-building methods.

In terms of the reported performance results, the proposed model is expected to be quite significant and notable for researchers.

## Conflict of interest statement

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## References

[1] Ge Yao J, Wan X, Xiao J. Recent advances in document summarization. Knowl Inf Syst 2017;53(2):297–336.
[2] Ermakova L, Cossu JV, Mothe J. A survey on evaluation of summarization methods. Inf Process Manage 2019;56(5):1794–814.
[3] Hark C, Seyyarer A, Uçkan T, Karci A. Doğal dil işleme yaklaşimlari ile yapisal olmayan dökümanlarin benzerliği. In: IDAP 2017 – international artificial intelligence and data processing symposium. p. 1–6.
[4] Mao X, Yang H, Huang S, Liu Y, Li R. Extractive summarization using supervised and unsupervised learning. Expert Syst Appl 2019;133:173–81.
[5] Hark C, Uçkan T, Seyyarer A, Karci Abubekir. Metin Özetleme İçin Çizge Tabanlı Bir Öneri. IDAP 2018 – international artificial intelligence and data processing symposium, 2018.
[6] Joshi A, Fidalgo E, Alegre E, Fernández-Robles L. SummCoder: an unsupervised framework for extractive text summarization based on deep auto-encoders. Expert Syst Appl 2019;129:200–15.
[7] Gambhir M, Gupta V. Recent automatic text summarization techniques: a survey. Artif Intell Rev 2017;47(1):1–66.
[8] Tan J, Wan X, Xiao J. Abstractive document summarization with a graph-based attentional neural model. In: Proceedings of the 55th annual meeting of the association for computational linguistics. p. 1171–81.
[9] Mihalcea R, Tarau P. A language independent algorithm for single and multiple document summarization. In: Proc. IJCNLP 2005, 2nd int. join conf. nat. lang. process.. p. 19–24.
[10] Sarkar K, Saraf K, Ghosh A. Improving graph based multidocument text summarization using an enhanced sentence similarity measure. In: 2015 IEEE 2nd int. conf. recent trends inf. syst. ReTIS 2015 – proc.. p. 359–65.
[11] Van Lierde H, Chow TWS. Query-oriented text summarization based on hypergraph transversals. Inf Process Manage 2019;56(4):1317–38.
[12] Xiong S, Ji D. Query-focused multi-document summarization using hypergraph-based ranking. Inf Process Manage 2016;52(4):670–81.
[13] Erkan G, Radev DR. Lexrank: graph-based lexical centrality as salience in text summarization. J Artif Intell Res 2004;22:457–79.
[14] Kaynar O, Görmez Y, Işık YE, Demirkoparan F. Comparison of graph based document summarization method. In: 2017 International conference on computer science and engineering (UBMK). p. 598–603.
[15] Boppana R et al. Approximating maximum independent sets by excluding subgraphs. BIT 1992;32(February):180–96.
[16] Ferreira R et al. Assessing sentence scoring techniques for extractive text summarization, 2013.
[17] Luhn HP. The automatic creation of literature abstracts. IBM J Res Dev 1958;2 (2):159–65.
[18] Shardan R, Kulkarni U. Implementation and evaluation of evolutionary connectionist approaches to automated text summarization, 2010.
[19] Nasr Azadani M, Ghadiri N, Davoodijam E. Graph-based biomedical text summarization: an itemset mining and sentence clustering approach. J Biomed Inform 2018;84:42–58.
[20] Student PG, Coe DM. A comparative study of Hindi text summarization techniques: genetic algorithm and neural network, 2015.
[21] Gupta V. Hybrid algorithm for multilingual summarization of Hindi and Punjabi documents. In: Mining intelligence and knowledge exploration. Springer; 2013. p. 717–27.
[22] Gupta V, Lehal GS. A survey of text summarization extractive techniques. J Emerg Technol Web Intell 2010;2(3):258–68.
[23] Abuobieda A, Salim N, Albaham AT, Osman AH, Kumar YJ. Text summarization features selection method using pseudo genetic-based model. In: Proc. – 2012 int. conf. inf. retr. knowl. manag. CAMP'12. p. 193–7.
[24] Fattah MA, Ren F. GA, MR, FFNN, PNN and GMM based models for automatic text summarization. Comput Speech Lang 2009;23(1):126–44.

[25] Nandhini K, Balasundaram SR. Improving readability through extractive summarization for learners with reading difficulties. Egypt Inform J 2013;14 (3):195–204.

[26] Brin S, Page L. The anatomy of a large-scale hypertextual web search engine. Comput Networks ISDN Syst 1998;30(1–7):107–17.

[27] Parveen D, Ramsl H-M, Strube M. Topical coherence for graph-based extractive summarization. In: Proceedings of the 2015 conference on empirical methods in natural language processing. p. 1949–54.

[28] Canhasi E, Kononenko I. Weighted archetypal analysis of the multi-element graph for query-focused multi-document summarization. Expert Syst Appl 2014;41(2):535–43.

[29] Salton G, Singhal A, Mitra M, Buckley C. Automatic text structuring and summarization. Inf Process Manage 1997;33(2):193–207.

[30] Medelyan O. Computing Lexical Chains with Graph Clustering. In: Proceedings of the ACL 2007 student research workshop. p. 85–90.

[31] Chen Y-N, Huang Y, Yeh C-F, Lee L-S. Spoken lecture summarization by random walk over a graph constructed with automatically extracted key terms. In twelfth annual conference of the international speech communication association, 2011.

[32] Plaza L, Stevenson M, Díaz A. Resolving ambiguity in biomedical text to improve summarization. Inf Process Manage 2012;48(4):755–66.

[33] Nallapati R, Zhou B, dos Santos C, Gulçehre Ç, Xiang B. Abstractive text summarization using sequence-to-sequence RNNs and beyond. In: *CoNLL 2016 – 20th signll conf. comput. nat. lang. learn. proc.* p. 280–90.

[34] Moawad IF, Aref M. Semantic graph reduction approach for abstractive text summarization. In: Proc. – ICCES 2012 2012 int. conf. comput. eng. syst.. p. 132–8.

[35] Linqing Liu HL, Yao Lu, Yang Min, Qu Qiang, Zhu Jia. Convolutional neural networks for sentence classification. In: EMNLP 2014 – 2014 conf. empir. methods nat. lang. process. proc. conf.. p. 1746–51.

[36] Sinha R, Mihalcea R. Unsupervised graph-based word sense disambiguation using measures of word semantic similarity. In: ICSC 2007 international conference on semantic computing. p. 363–9.

[37] Fattah MA, Ren F. GA, MR, FFNN, PNN and GMM based models for automatic text summarization, 2008.

[38] Author C, Shardanand Prasad R, Kulkarni U. Implementation and evaluation of evolutionary connectionist approaches to automated text summarization. J Comput Sci 2010;6(11):1366–76.

[39] Gupta V. LNAI 8284 – Hybrid algorithm for multilingual summarization of Hindi and Punjabi documents, 2013.

[40] Quezada-Sarmiento PA, Macas-Romero J del C, Roman C, Martin JC. A body of knowledge representation model of ecotourism products in southeastern Ecuador. Heliyon 2018;4(12).

[41] Quezada-Sarmiento PA, Enciso-Quispe LE, Jumbo-Flores LA, Hernandez W. Knowledge representation model for bodies of knowledge based on design patterns and hierarchical graphs. Comput Sci Eng 2018;PP(c):1.

[42] Oh S. Maximal independent sets on a grid graph, 2017.

[43] Jou MJ, Chang GJ. The number of maximum independent sets in graphs. Taiwan J Math 2000;4(4):685–95.

[44] Chan TM, Har-Peled S. Approximation algorithms for maximum independent set of pseudo-disks, 2011.

[45] Feo TA, Resende MGC, Smith SH. A greedy randomized adaptive search procedure for maximum independent set. Oper Res 2008;42(5): 860–78.

[46] NIST. Document understanding conferences. NIST [Online]. Available: https://www-nlpir.nist.gov/projects/duc/ [Accessed: 08-May-2019].

[47] Lin CY. Rouge: a package for automatic evaluation of summaries. In: Proc. work. text Summ. branches out (WAS 2004). p. 25–6.

[48] Lin C-Y, Hovy E. Automatic evaluation of summaries using N-gram co-occurrence statistics. In: Automatic evaluation of summaries using n-gram co-occurrence statistics. p. 150–7.

[49] Landauer TK, Foltz PW, Laham D. An introduction to latent semantic analysis. Discourse Process 1998;25(2–3):259–84.

[50] Landauer TK, Dutnais ST. A solution to Plato's problem: the latent semantic analysis theory of acquisition, induction, and representation of knowledge. Psychol Rev 1997;104(2):211–40.

[51] Mihalcea R. Language independent extractive summarization. In: Proc. ACL 2005 interact. poster demonstr. sess. – ACL '05, no. June. p. 49–52.

[52] Mihalcea R, Tarau P. TextRank: bringing order into texts. Proceedings of the ACL 2004 on interactive poster and demonstration sessions, 2004. pp. 20-es.

[53] Vanderwende L, Suzuki H, Brockett C, Nenkova A. Beyond SumBasic: task-focused summarization with sentence simplification and lexical expansion. Inf Process Manage 2007;43(6):1606–18.

[54] Haghighi A, Vanderwende L. Exploring content models for multi-document summarization 2009;(June),362.