

A Novel Gravity Optimization Algorithm for Extractive Arabic Text Summarization

Mustafa J. Hadi *, Ayad R. Abbas , Osamah Y. Fadhil 

Department of Computer Sciences, University of Technology, Baghdad, Iraq

*Corresponding Author.

Received 11/09/2022, Revised 11/03/2023, Accepted 13/03/2023, Published Online First 20/07/2023



© 2022 The Author(s). Published by College of Science for Women, University of Baghdad.

This is an Open Access article distributed under the terms of the [Creative Commons Attribution 4.0 International License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Abstract

An automatic text summarization system mimics how humans summarize by picking the most significant sentences in a source text. However, the complexities of the Arabic language have become challenging to obtain information quickly and effectively. The main disadvantage of the traditional approaches is that they are strictly constrained (especially for the Arabic language) by the accuracy of sentence feature functions, weighting schemes, and similarity calculations. On the other hand, the meta-heuristic search approaches have a feature that tolerates imprecision, gets prohibited results, and is not strictly bound by the above restrictions. This paper used the Gravitational Optimization Algorithm (GOA), a powerful metaheuristic approach based on the law of gravity, to address the challenge of extractive summarizing Arabic texts. The objective function of the GOA algorithm is derived based on sentence significance, such as its length, similarity degree, position, statistical term frequency, and named entity ownership. Essex Arabic Summaries Corpus (EASC) was used to evaluate the proposed method and measured by the Recall-Oriented Understudy for Gisting Evaluation (ROUGE). The proposed approach achieved 68.04% Recall, 58.49% Precision, and 60.05% F1-measure using ROUGE-1, higher than standard summarizers and metaheuristic approaches.

Keywords: Abstractive Summarization, Extractive Summarization, Arabic Text Summarization, Similarity Graph, Gravitational Optimization Algorithm.

Introduction

In recent years, Internet users from the Arab world have increased rapidly. However, digital content in the Arabic language is still lacking perfect .development plans Text summarization generates the most informative sentences or a summary from minous texts to reduce volureading time and accelerate information search. Text summarization can be classified into two techniques: Extractive summarization and abstractive summarization. The first technique identifies the meaningful sentences input text and reproduces them from the .as a summary

In contrast, the second technique interprets the input text and generates new summary text using advanced Natural Language Process (NLP) techniques¹. These two techniques can be applied to single-document or document summarization-multi². In addition, according to the language, summarization systems can be classified into two types. The first is monolingual summarization systems, which work only inone language. The second systems, is multilingual summarization covering more than one language³.

In general, Arabic text summarization approaches and methodologies are still immature due to many challenges in the Arabic language: For example, Arabic text meaning depends on the variation, presence or dialectal v-context, cross absence of diacritics, and the evaluation process of Arabic summarization systems^{4,5}. One of the most effective methods for solving text summarization is metaheuristic search algorithms like cuckoo search⁶, ant colony⁷, artificial bee colony⁸, Particle Swarm Optimization (PSO)⁹, and genetic algorithm (GA)¹⁰. These optimization algorithms can be helpful in optimization problems to select appropriate sentences from the text and build a representative summary.

The difficult problem facing researchers in dealing with the Arabic language is that it is a highly inflectional and derivational language and the preprocessing tools of the Arabic language are still lacking improvement.

There is an essential disadvantage of the traditional approaches used in summarization is that they are strictly constrained (especially for the Arabic language) to the accuracy of sentence features, weighting schemes, and similarity computations. On the contrary, the metaheuristic search approaches are tolerated imprecision, get prohibited promising results, and are not strictly bound by the aforementioned biased restrictions. Most metaheuristic search approaches deal with a continuous (real point vectors) model. The challenge for many studies is how to apply these approaches to an environment with discontinuous elements (summarization as an example). In order to accomplish this task, many researchers modify the original metaheuristic search approaches by a

significant change in the algorithm structure or in its equations. In fact, the unprofessional changes in the structure or the equation of an algorithm may take the algorithm's goal away and get unbalanced solutions. Therefore, modifying an algorithm and investing it without negatively changing its natural path is a great achievement in itself.

This research studies the GOA algorithm and applies it in the summarization environment. However, a big challenge is reducing the difference between using real item space and discrete item space. Therefore, challenge this has been successfully tackled by proposing a new method that combines NLP with GOA (as a metaheuristic approach) augmented by a constructed neighborhood area based on a text similarity graph.

The main two contributions of this work can be summarized as follows:

1- Investigate the ability of a novel GOA algorithm to address the performance problems in the summarization environment in terms of time and solution quality.

2- Addressing the challenges of the poor performance of the available Arabic text summarization systems due to the fact that the Arabic language is a highly inflectional and derivational language and due to the preprocessing tools of the Arabic language are still imperfect tools.

The remainder of this paper is organized as follows: Section 2 reviews prior work in the areas of the Arabic text summarization method. Section 3 introduces the methodology for GOA. Section 4 presents the proposed Arabic text summarization. Section 5 presents the experimental results, and the last section presents the conclusions.

Materials and Methods

Related Works:

This section only covers Arabic text summarization studies that use abstractive and extractive approaches. Since 2015, only four studies have focused on using the abstractive technique.

Typically, this type of summarization is more complex to implement, although sometimes it is effective. For example, some authors¹¹ proposed a textual graph-based model to remove multi-document redundancy and generate a coherent summary using concatenating related sentences. Unfortunately, the experimental results were only on

the reduction ratio but neglected the enhancement of the accuracy of the r that, summarization model. After Other authors¹² introduced an abstractive Arabic text summarizer based on the Rough set theory. It starts by segmenting the input text and applying a -rule based sentence reduction technique

Nevertheless, this requires human intervention to evaluate the proposed method. In 2020, two studies took advantage of deep learning. The first¹³ used a deep neural network learning methodology that deals with long texts more efficiently by identifying focus points in the text. However, the accuracy s did not exceed 60%. The second result¹⁴ proposed

an abstractive Arabic text summarization model based on sequence-to-sequence RNN encoder-decoder architecture.

Furthermore, five studies have focused on using the extractive technique in the last four years. For instance, some studies^{15, 16} proposed metaheuristic searches like PSO and Firefly (FF) to extract summaries for single Arabic documents on the EASC corpus. A study¹⁷ introduced two summarizing techniques, including score-based and supervised machine learning using only a single document. Each sentence is evaluated using a novel formulation that considers sentence diversity, relevance, and coverage based on a combination of semantic and statistical features.

ASDKGA is presented as a single-document text summarizing approach that combines statistical features, domain expertise, and genetic algorithms to extract key ideas from Arabic political documents on the EASC corpus¹⁸.

Moreover, other studies have used hybrid integration techniques including both abstractive and extractive methods to provide an informative and coherent summary of a long document^{19,20}.

In fact, and based on current studies, the preprocessing tools of the Arabic language are still problematic. As a result, this paper used the GOA, a powerful metaheuristic approach based on the law of gravity, to address the challenge of summarizing Arabic texts. The proposed method exceeds the previous methods in terms of performance because of its ability to access areas considered forbidden within the research space.

Method:

Gravitational Optimization Algorithm (GOA)

GOA is one of the newest heuristic algorithms²¹. The algorithm is based on gravity and mass interactions at a low level. The solutions in the GOA population are referred to as agents; these agents interact with one another through gravity. Therefore, this represents the global movements of the agent, while the agent with the highest mass represents the algorithm's exploitation step. The solution with the higher mass is the best. The gravitational constant G is calculated using Eq.1 at iteration t .

$$G_0 e^{-\frac{\alpha t}{T}} \quad 1$$

Where G_0 and α are initialized at the beginning of the search, and their values are decreased as the search progresses. The total number of iterations is denoted by T .

The masses of the objects obey Newton's law of gravity using Eq.2:

$$F = G \frac{M_1 M_2}{R^2} \quad 2$$

F is the gravitational force magnitude, and G is the gravitational constant. M_1 is the first object's mass; M_2 is the second object's mass; and R is the distance between the two objects M_1, M_2 .

When a force F is applied to an object, the object moves with acceleration a . Whereas a depends on the applied force and the mass M , as shown in Eq.3 below:

$$a = \frac{F}{M} \quad 3$$

The Eq.4 and Eq.5 are used to calculate the velocity and position of the agents in the next iteration ($t+1$), respectively:

$$v_i(t+1) = rand_i \times v_i(t) + a_i(t) \quad 4$$

$$x_i(t+1) = rand_i \times x_i(t) + a_i(t) \quad 5$$

Where $rand_i$ is the random number in the range $[0, 1]$.

Proposed Text-Summarization Model

In the context of the summarization problem, it can be thought that the population is a complete text whose elements are sentences. However, if each real point represents a sentence of a text best recognized after a series of iterations, how is this sentence represented in a real point vector? This indeed needs an innovative method to drop sentences in real point vectors. This task has been addressed by using GOA metaheuristic search tool which has features of ease of implementation, convergence stability, and low computational cost. The proposed text-summarization model has multi-steps explained in Fig.1.

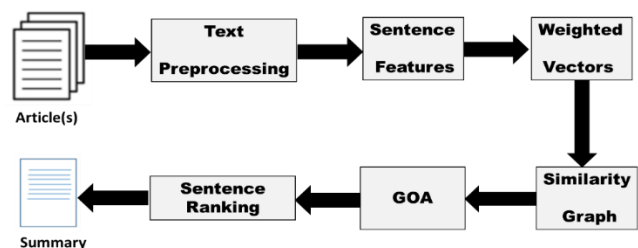


Figure 1. The proposed text-summarization model

The text-summarization algorithm of the proposed model is shown as follows:

- 1: Input: Text to be summarized
- 2: Initialization procedure:
 - 2.1: Preprocessing: Tokenization, stemming, and stop word removal.
 - 2.2: Similarity graph creation: Each sentence has several sentences intersected with several identical words.
- 3: Initialize population procedure:
 - 3.1: Solutions (sentences) are randomly selected from the similarity graph concerning population size.
 - 3.2: Set the solution similarity values as points in the search space.
- 4: Compute the fitness for each solution using the fitness function.
- 5: Update the gravitational G and K best constants.
- 6: Compute the masses, total forces, acceleration, and velocities.
- 7: Update positions in the population depending on the axes directions of the velocity, vector and after that, it can:
 - 7.1: If it is a premise, then accept the new solution
Else accept another solution based on probability:
If $\text{random}() > G$: Select a random solution from the K best solutions
Else: Select a solution from the whole state space
- 8: Repeat Step 4 to Step 7 until producing the same fitness values or reaching a predefined .limit

Text Preprocessing Steps

Text preprocessing is a procedure which can be divided mainly into four text operations:

- 1-Tokenize the raw text to extract the terms .
- 2-Lexical analysis of the terms with the objective of treating digits, hyphens, punctuation marks, and the case folding.
- 3-Elimination of stop words with the objective of filtering out words with very low discrimination values.
- 4-Stemming of the remaining terms for allowing the retrieval of documents containing syntactic variations of query terms.

There are many available Arabic stemmers. In this work, ISRI stemmer is used to stem the Arabic words. ISRI (stands for Information Science Research Institute) stemmer is a new

root-extraction stemmer without a root dictionary. This feature makes ISRI stemmer more capable of stemming rare and new words.

Calculating Fitness Function

The structure of sentence features measures each sentence's score in the text to rank each sentence. For example, the following statistical sentence features from $f1$ to $f5$ are used to allocate a score or fitness to each sentence:

1. Sentence length (f_1): The longest sentence contains essential information; it can be calculated by the number of words in the current sentences divided by the max sentence length.
2. Similarity degree (f_2): Sentence i is nearest to sentences with t cosine similarities. The more similarity the sentence gets, the better it is.
3. Sentence position (f_3): Usually, the informative sentences in a text covered by writers at the beginning and end of any article show the importance of sentences. In contrast, the middle sentence is relative using Eq.6.

$$\text{Max} \left[\frac{1}{l}, \frac{1}{n-i+1} \right] \quad 6$$

where n = number of sentences in the document n and i = position of the sentence

4. Statistical term frequency (f_4): Average TF-IDF for all the words in the sentence.
5. Named entity ownership (f_5): The more a sentence has a named entity, the better it is. For example, the following is Eq.7 for calculating the fitness function.

$$\text{FitnessFunction} = \sum_{k=1}^5 f_k \quad 7$$

K is the number of statistical sentence features equal to five features from $f1$ to $f5$.

Building the Similarity Graph

An efficient search space structure based on a text similarity graph augmented the gravitational optimization algorithm. The similarities between data points can be organized in graphs for solving a range of practical problems. Let $G = (V, E)$ is a graph, V represents a set of vertices v_i and E represents a set of edges e_{ij} Let $x_1 \dots x_n$ is a set of data points, the similarity between all pairs of data points x_i and x_j is noted by $w_{ij} \geq 0$. In the G graph, each data point x_i is represented by vertex v_i , and two vertices are connected if the similarity w_{ij} between them is

positive. The edge e_{ij} , then, is weighted by w_{ij} . The weighted adjacency matrix of the graph is the matrix $W_{n \times m} = w_{ij} = 1, \dots, n, i, j = 1, \dots, m$. If $w_{ij} = 0$, the vertices v_i and v_j are not connected²².

Practical Example

The following example illustrates the main procedure of the proposed summarization based on the GOA method. Let us have a text with seven sentences numbered 0, 1, ..., 6, as shown in Table 1.

Table 1. An example of seven Arabic sentences.

Id	Arabic Sentence	English Sentence
0	تعود كلمة الحاسوب في أصلها إلى كلمة حساب، وقد عرّف جهاز الحاسوب بأنه عبارة عن آلة حسابية عالية الدقة.	The word computer comes from arithmetic, and the computer was defined as a high-precision arithmetic machine.
1	تجمع الحواسيب بين ما يُعرف بالبرمجيات والمعدات مكونة معاً أجهزة الحاسوب الإلكترونية.	Computers combine what is known as software and hardware to form electronic computers.
2	تُعرف البرمجيات بأنها أحد المكونات الرئيسية لجهاز الحاسوب والتي تتكون من مجموعة من الأوامر البرمجية.	Software is defined as one of the main components of a computer, which consists of a set of program commands.
3	أما المكونات المادية فهي تلك الأجهزة والمعدات المادية التي يتكون منها جهاز الحاسوب.	As for the hardware components, they are those hardware and physical equipment that makes up a computer.
4	للحاسوب القدرة على حل العمليات الحسابية بسرعة كبيرة جداً والقدرة على التعامل مع عمليات حسابية مُعقدة وبدقة متناهية.	The computer can solve arithmetic operations quickly and deal with complex arithmetic operations with extreme accuracy.
5	تمتلك الحواسيب الإلكترونية القدرة على تخزين البيانات ومعالجتها أو استرجاعها وتقاس سرعتها بالميغاهيرتز.	Electronic computers can store, process, or retrieve data, and their speed is measured in megahertz.
6	العديد من الشركات العراقية قد أغنت الأسواق العراقية بمختلف أنواع أجهزة الحاسوب ذات الدقة العالية والسرعة الفائقة.	Many Iraqi companies have enriched the Iraqi market with various types of computers with high accuracy and speed.

As it can be seen in Table 2, sentence 0 is nearest (cosine similarity) to sentence 4 with 0.131 than others, and the second closest to sentence 6 with

0.069 than others, and so on. It should be noted that the summation of the cosine similarities gives specific importance to each sentence.

Table 2. Cosine similarity matrix of seven sentences.

ij	0	1	2	3	4	5	6
0	1.000	0.002	0.057	0.022	0.131	0.000	0.069
1	0.002	1.000	0.026	0.128	0.001	0.120	0.024
2	0.057	0.026	1.000	0.047	0.001	0.000	0.001
3	0.022	0.128	0.047	1.000	0.001	0.000	0.021
4	0.131	0.001	0.001	0.001	1.000	0.095	0.015
5	0.000	0.120	0.000	0.000	0.095	1.000	0.019
6	0.069	0.024	0.001	0.021	0.015	0.019	1.000

The key parameters of seven sentences are: Data size= 7; max_iterations= 4; Population size= 3, and Initial population= [(0.131, 0.069), [4, 0, 6, 2, 3, 1, 5], [1.0, 0.163, 1.0, 0, 0.281], [0.047, 0.128], [0, 0.684, 0.175, 0.25], [5, 4, 6, 0, 2, 3, 0.219], ([0.12, 0.095], [1, 5, 4, 6, 3, 2, 0], [0.5, 0.248, 0.684, 0, 0.234])].

The following computational steps summarize the GOA:

- Select random solutions (sentences) from the similarity graph, let 0, 3, and 5.
- Sentence 0 is nearest to 4, 6, 2, 3, 1, 5. It is nearest to 4 by 0.131 and nearest to 6 by 0.069. The others 2, 3, 1, 5 can be alternatives (local neighbors). Here, the vector [0.131, 0.069] is assumed as a point in the state space. And the list [1.0, 0.163, 1.0, 0, 0.281] contains the features of

this point, i.e. [Sentence length, Similarity degree, Sentence position, Statistical term frequency, Named entity affection]

6. The rest of the solutions, 3 and 5, have the same procedure.

According to the results in Table 3, In iteration 0, the K =3 can be obtained using the following Equation:

$$K = f \left(\text{pupolation} - (\text{pupolation} - 1) \times \left(\frac{\text{iteration}}{\text{float}(\text{maxiteration})} \right) \right) \quad 8$$

Table 3. The results of seven sentences.

New	Prob.	Delta	Move	Velocity	Old	G	k	Iter
]0.128 ,0.120[□ (3, 1, 5)	L	N	L]0.174 ,0.096[]0.131 ,0.069[□ (4, 0, 6)	1.000	3	0
]0.057 ,0.047[□ (0, 2, 3)	-	P	R]0.180 ,0.503]0.128 ,0.047[□ (1, 3, 2)			
]0.131 ,0.095[□ (0, 4, 5)	-	P	R]0.317 ,0.720[]0.120 ,0.095[□ (1, 5, 4)			
]0.057 ,0.047[□ (0, 2, 3)	L	N	R]0.174 ,0.413[]0.128 ,0.120[□ (3, 1, 5)	0.607	2	1
]0.131 ,0.095[□ (0, 4, 5)	L	N	R]0.266 ,0.705[]0.057 ,0.047[□ (0, 2, 3)			
]0.069 ,0.024[□ (0, 6, 1)	L	N	L]0.591,0.229[]0.131 ,0.095[□ (0, 4, 5)			
]0.131 ,0.069[□ (4, 0, 6)	-	P	L]0.326 ,0.251[]0.057 ,0.047[□ (0, 2, 3)	0.368	2	2
]0.069 ,0.024[□ (0, 6, 1)	L	N	R]0.142 ,0.733[]0.131 ,0.095[□ (0, 4, 5)			
]0.128 ,0.120[□ [3, 1, 5(H	N	R]0.310 ,0.443[]0.069 ,0.024[□ (0, 6, 1)			
]0.069 ,0.024[□ (0, 6, 1)	H	N	R]0.176 ,0.197[]0.131 ,0.069[□ (4, 0, 6)	0.223	1	3
]0.131 ,0.069[□ (4, 0, 6)	-	P	L]0.427 ,0.120[]0.069 ,0.024[□ (0, 6, 1)			
]0.120 ,0.095[□ (1, 5, 4)	H	N	R]0.272 ,0.310[]0.128 ,0.120[□ [3, 1, 5(

In the beginning, all agents apply the force, then K is decreased linearly, and at the end, there will be just one agent using force on the others. The constant gravitational G=1 is computed using Eq 1. G value is decreased with time to control the search accuracy. The population [0] = ([0.131, 0.069], [4, 0, 6, 2, 3, 1, 5], [1.0, 0.163, 1.0, 0, 0.281]) is found from the initial population. The velocity= [0.174, 0.096] moving left (L). The candidate solution is ([0.131, 0.095], [0, 4, 5, 6, 1, 3, 2], [0.333, 0.168, 0.947, 0, 0.244]).

Both old and new fitnesses are computed using Eq.7, where $fitness_{old} = \sum([1.0,0.163,1.0,0,0.281]) \times 0.2 = 0.489$, while $fitness_{new} = \sum([0.333,0.168,0.947,0,0.244]) \times 0.2 = 0.338$. After that, the delta is obtained by finding the difference between old and new fitnesses; then, the Delta value is -0.150, negative (N). Because the

probability (0,1) is less than the value of G, the random so Prob.=L. That means alution is accepted from whole space (except itself) = ,3] ,[0.12 ,0.128]) The same computational steps are ([.4 ,0 ,6 ,2 ,5 ,1 [1] applied on populationand population [2], the new random solutions can be obtained and accepted as ([0.057, 0.047], [0, 2, ,0.131]) and [5 ,4 ,6 ,1 ,3 respectively. After four ,([2 ,3 ,1 ,6 ,5 ,4 ,0] ,[0.095 iterations, thethree best solutions are obtained (0, 5, and 6), as shown in Table 4.

Table 4. Summary of seven sentences.

Id	Arabic Sentence	English Sentence
0	تعود كلمة الحاسوب في أصلها إلى كلمة حساب، وقد عرّف جهاز الحاسوب بأنه عبارة عن آلة حسابية عالية الدقة.	The word computer comes from arithmetic, and the computer was defined as a high-precision arithmetic machine
5	تمتلك الحواسيب الالكترونية القدرة على تخزين البيانات ومعالجتها او استرجاعها وتقاس سرعتها بالميجاهيرتز.	Electronic computers can store, process, or retrieve data, and their speed is measured in megahertz
6	العديد من الشركات العراقية قد أغنت الاسواق العراقية بمختلف انواع اجهزة الحاسوب ذات الدقة العالية والسرعة الفائقة.	Many Iraqi companies have enriched the Iraqi market with various types of computers with high accuracy and speed

Results and Discussion

The EASC corpus (Essex Arabic Summaries Corpus) was used to test the performance of the proposed method. It is an Arabic natural language resource. It contains 153 Arabic articles and 765 human-generated extractive summaries of articles. The number of sentences in EASC equal 2360 and the number of words equals 41493. EASC is publicly available for advancing research on Arabic text summarization. The summaries were generated using Mechanical Turk. In this paper

-was used to compute the effectiveness of auto-generated summaries. ROUGE scores are reported summarizing using three commonly used metrics (Precision, Recall, and F1-measure) compared with several standard summarizers like Text Rank, SumBasic, KLSum, and LSA methods. The experimental results in Fig 2 show the Recall, Precision, and F1-measure using the evaluation of ROUGE-1 with 68.04%, 58.49% and 60.05% respectively, and they are higher than the TextRank, SumBasic, KLSum, and LSA.

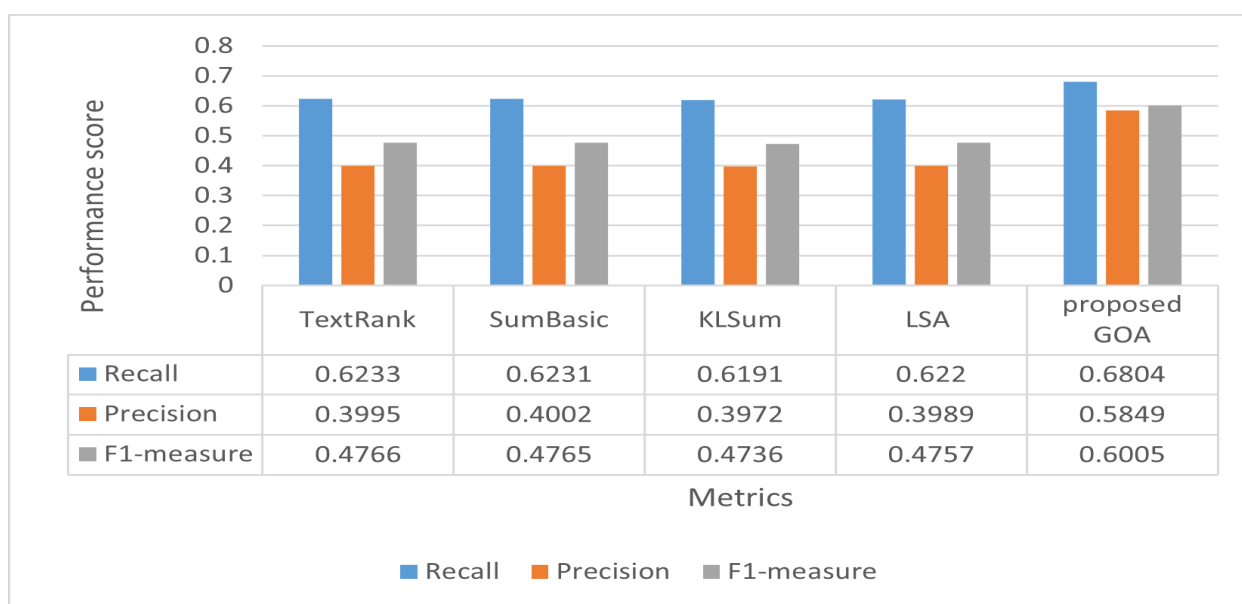


Figure 2. Performance evaluation of Rouge-1.

Likewise, all experimental results in Fig 3 show that the Recall, Precision, and F1-measure using the evaluation of ROUGE-2 was 60.95%, 52.07%, and

53.48% higher than the TextRank, SumBasic, KLSum, and LSA.

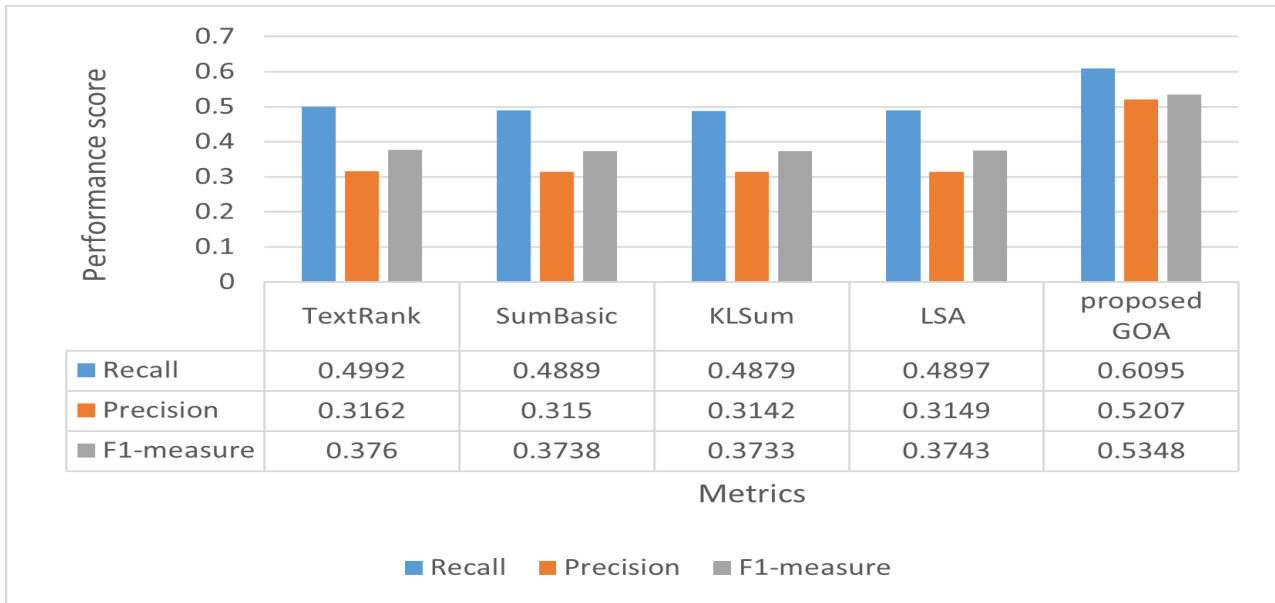


Figure 3. Performance evaluation of Rouge-2.

Finally, all experimental results in Fig 4 show that the Recall, Precision, and F1-measure using the evaluation of ROUGE-SU4 was 61.33%, 53.39%,

and 54.60% higher than the TextRank, SumBasic, KLSum, and LSA.

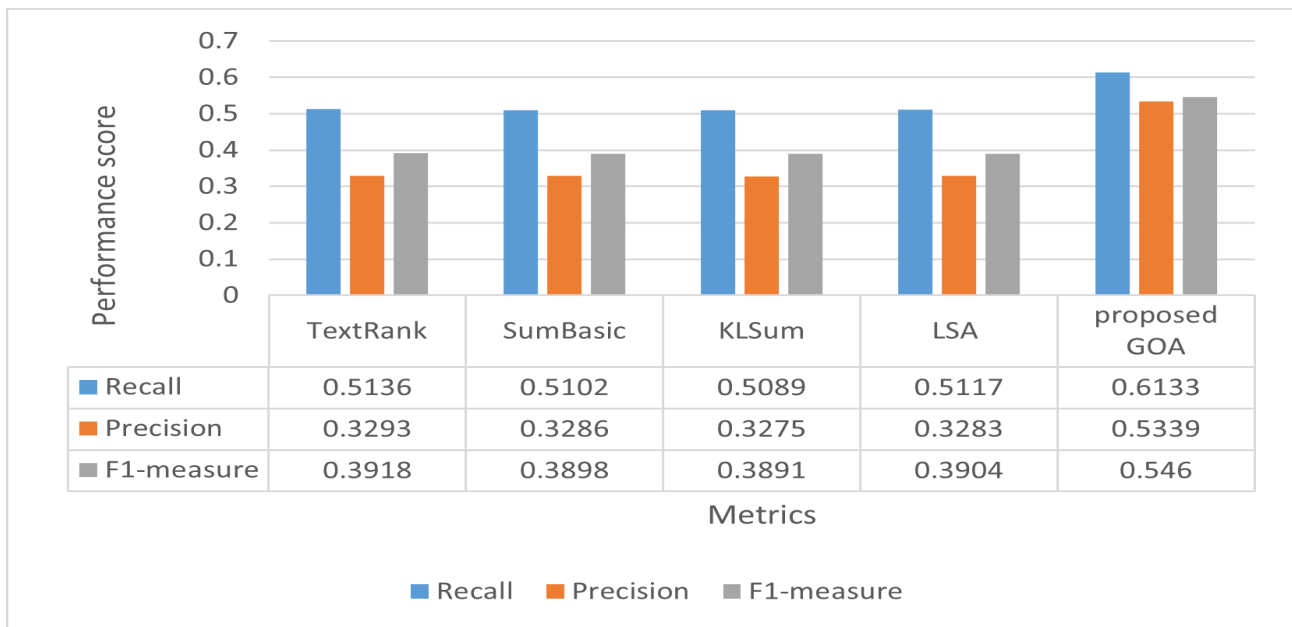


Figure 4. Performance evaluation of Rouge-SU4.

The novelty in this paper is that the GOA is augmented by an efficient search space structure based on a text similarity graph. This graph structure has a significant role in feeding the proposed algorithm to find the optimal solutions and

improve the convergence speed. The algorithm will look for promising solutions in advanced stages during its search process and within a reasonable time.

The proposed algorithm has an objective function computed based on the significant features of the sentences (such as the length, position, term frequency, similarity degree, and named entity recognition).

This GOA approach is compared with metaheuristic approaches like GA¹⁰, PSO^{15, 23}, and FF^{16, 24}, as shown in Table 5. All these approaches were evaluated on ESAC corpus, using ROUGE-1 and

ROUGE-2 metrics, except ROUGE-SU4 has been used in this paper. Although all previous approaches did not specify the average values of these metrics based on the number of documents, our approach has scored higher Recall, Precision 1-using ROUGE measure values than the values obtained by -and F1 the GA, PSO, and FF. At the sametime, ROUGE-2 scored a higher Recall value than the GA, PSO, and FF values.

Table 5. Comparisons against other summarization approaches.

F1-measures	Precision	Recall	ROGUE	Approach
0.5476	0.5658	0.5713	ROUGE-1	GA
0.4465	0.4597	0.4710	ROUGE-2	GA
0.5532	0.5882	0.5444	ROUGE-1	PSO
0.4538	0.4814	0.4483	ROUGE-2	PSO
0.5732	0.5732	0.6014	ROUGE-1	FF
0.6005	0.5849	0.6804	ROUGE-1	GOA
0.5348	0.5207	0.6095	ROUGE-2	GOA
0.5461	0.5339	0.6133	ROUGE-SU4	GOA

In this work, the experimental tests of GOA approach have an explicit superiority over the other approaches. GOA approach has a few parameters. All calculated by their own equations, Eqs (1-8). The best maximum iteration is 100 and the population size P_{Size} is computed as parentage of the original data size using the following Equation:

$$P_{Size} = \text{int}(\text{round}(\frac{Data_{Size}}{100}) * \text{Percentage})$$

9

Conclusion

This paper proposes a new method combining NLP and a metaheuristic approach to summarize Arabic text with single documents. Three phases are applied in text summarization: text preprocessing, building a similarity graph, and GOA. The experimental results are compared with several standard summarizers and metaheuristic approaches. The proposed approach has higher metrics values than standard summarizers (TextRank, SumBasic, KLSum, and LSA) or metaheuristics (GA, PSO, and FF). In addition, a summarization environment has been successfully used with a discrete item space dropped on continuous item space by using GOA after reinforcing it with a constructed neighborhood area based on a text similarity graph. This graph structure has a significant role in feeding the proposed algorithm to find the optimal solutions and

The percentage used in this work is 30%.

The main limitation of the proposed method, although it is superior to other methods, is that its results are still affected by the ambiguity present in the Arabic words. The so-called Arabic Word Sense Disambiguation (WSD) is not yet complete and available, and the inclusion of an available Arabic WSD has an uncertain resultant and more time-consuming outcome. The process of optimizing the Arabic WSD remains a major challenge in the literature.

improve the convergence speed. The algorithm looked for promising solutions in advanced stages during its search process and within a reasonable period. The proposed algorithm, the graph structure and the GOA algorithm style make the advantage to reaching fruitful areas (promise sentences) were statistically forbidden because of inaccurate similarity calculations for unperfect Arabic features. The proposed system achieves a Although, the real challenge to address the bias resulting from the unperfect features obtained from unperfect Arabic preprocessing tools. However, there is still no perfect Arabic stemmer, no perfect word sense Arabic semantic disambiguation, and no perfect analysis. The ROUGE evaluation metrics reveal that the proposed GOA-based method is superior to the other standard methods in accuracy and less

computational effort. Because of a few abstractive models available in the Arabic language. Future work should try to optimize the abstractive

Arabic text summarization model using deep learning that will automatically generate a summary from a long text.

Author's Declaration

- Conflicts of Interest: None.
- We hereby confirm that all the Figures and Tables in the manuscript are ours. Furthermore, any Figures and images, that are not ours, have been

- included with the necessary permission for re-publication, which is attached to the manuscript.
- Ethical Clearance: The project was approved by the local ethical committee in University of Technology.

Author's Contribution Statement

M. J. H.: design, acquisition of data, analysis, interpretation; A. R. A.: editing, revision, and proofreading; O. Y. F.: revision

References

1. Gupta V, Lehal GS. A Survey of Text Summarization Extractive Techniques. *J Emerg Technol Web Intell.* 2010; 2: 258-268. <http://www.jetwi.us/uploadfile/2014/1226/20141226030617764.pdf>
2. Mani K, Verma I, Meisheri H, Dey L. Multi-document summarization using distributed bag-of-words model. *IEEE/WIC/ACM Int Conf Web Intell (WI).* 2018; 672-675. <https://doi.org/10.1109/WI.2018.00-14>
3. To HQ, Nguyen KV, Nguyen NL-T, Nguyen AG-T. Monolingual versus Multilingual BeRTology for Vietnamese Extractive Multi-Document Summarization. *arXiv preprint. arXiv.* 2021. 2108.13741. <https://doi.org/10.48550/arXiv.2108.13741>
4. Al-Saleh AB, Menai MEB. Automatic Arabic text summarization: a survey. *Artif Intell Rev.* 2016; 45: 203-234. <https://doi.org/10.1007/s10462-015-9442-x>
5. Al Qassem LM, Wang D, Al Mahmoud Z, Barada H, Al-Rubaie A, Almoosa NI. Automatic Arabic summarization: a survey of methodologies and systems. *Procedia Comput Sci.* 2017; 117: 10-18. <https://doi.org/10.1016/j.procs.2017.10.088>
6. Mirshojaei SH, Masoomi B. Text summarization using cuckoo search optimization algorithm. *J Comput Robot.* 2015; 8: 19-24. https://jcr.qazvin.iau.ir/article_683_e08cfc39b8a850adf76246e0096d3d22.pdf
7. Hassan OF. Text summarization using ant colony optimization algorithm. *Sudan University of Science and Technology,* 2015. <https://repository.sustech.edu/handle/123456789/11173>
8. Sanchez-Gomez JM, Vega-Rodríguez MA, Pérez CJ. Extractive multi-document text summarization using a multi-objective artificial bee colony optimization approach. *Knowl Based Syst.* 2018; 159: 1-8. <https://doi.org/10.1016/j.knosys.2017.11.029>
9. Gamal M, El-Sawy A, AbuEl-Atta AH. Hybrid Algorithm Based on Chicken Swarm Optimization and Genetic Algorithm for Text Summarization. *Int J Intell Eng Syst.* 2021; 14: 319-331. <http://www.inass.org/2021/2021063027.pdf>
10. Jaradat YA, Al-Taani AT. Hybrid-based Arabic single-document text summarization approach using genetic algorithm. *2016 7th Int Conf Inf Commun Syst.* 2016; 85-91. <https://doi.org/10.1109/IACS.2016.7476091>
11. Alwan MA, Onsi HM. A Proposed Textual Graph Based Model for Arabic Multi-document Summarization. *Int J Adv Comput Sci Appl.* 2016; 7: 435-439. <https://dx.doi.org/10.14569/IJACSA.2016.070656>
12. Azmi AM, Altmami NI. An abstractive Arabic text summarizer with user controlled granularity. *Inf Process Manag.* 2018; 54: 903-921. <https://doi.org/10.1016/j.ipm.2018.06.002>
13. Al-Maleh M, Desouki S. Arabic text summarization using deep learning approach. *J Big Data.* 2020; 7: 1-17. <https://doi.org/10.1186/s40537-020-00386-7>
14. Suleiman D, Awajan R. Deep learning based abstractive Arabic text summarization using two layers encoder and one layer decoder. *J Theor Appl Inf Technol.* 2020; 98: 3233-3244. <file:///home/uu/Downloads/5Vol98No16.pdf>
15. Al-Abdallah RZ, Al-Taani AT. Arabic single-document text summarization using particle swarm optimization algorithm. *Procedia Comput Sci.* 2017; 117: 30-37. <https://doi.org/10.1016/j.procs.2017.10.091>
16. Al-Abdallah RZ, Al-Taani AT. Arabic text summarization using firefly algorithm. *Amity Int Conf Artif Intell* 2019; 61-65. <https://doi.org/10.1109/AICAI.2019.8701245>
17. Qaroush A, Farha IA, Ghanem W, Washaha M, Maali E. An efficient single document Arabic text

- summarization using a combination of statistical and semantic features. J King Saud Univ- Comput Inf Sci. 2021; 33: 677-692. <https://doi.org/10.1016/j.jksuci.2019.03.010>
18. Al-Radaideh QA, Bataineh DQ. A Hybrid Approach for Arabic Text Summarization Using Domain Knowledge and Genetic Algorithms. Cognit Comput. 2018; 10: 651-669. <https://doi.org/10.1007/s12559-018-9547-z>
19. Ali ZH, Hussein AK, Abass HK, Fadel E. Extractive multi document summarization using harmony search algorithm. Telkomnika, Telecomm Comput, Electro Cont. 2021; 19: 89-95. <http://doi.org/10.12928/telkomnika.v19i1.15766>
20. Elmadani KN, Elgezouli M, Showk A. BERT Fine-tuning for Arabic Text Summarization. arXiv.2020; 2004. 14135. <https://doi.org/10.48550/arXiv.2004.14135>
21. Rashedi E, Nezamabadi-Pour H, Saryazdi S. GSA: a gravitational search algorithm. Inf Sci. 2009; 179: 2232-2248. <https://doi.org/10.1016/j.ins.2009.03.004>
22. Hassan AKA, Hadi MJ. Distributed Information Retrieval Based on Metaheuristic Search and Query Expansion. J Kufa Math Comput; 2017; 4.3: 4-11. <https://doi.org/10.31642/JoKMC/2018/040302>
23. Iqbal Z, Ilyas R, Chan HY, Ahmed N. Effective Solution of University Course Timetabling Using Particle Swarm Optimizer based Hyper Heuristic Approach. Baghdad Sci J. 2021; 18: 1465- 1475. [https://doi.org/10.21123/bsj.2021.18.4\(Suppl.\).1465](https://doi.org/10.21123/bsj.2021.18.4(Suppl.).1465)
24. Al-Behadili HNK. Improved Firefly Algorithm with Variable Neighborhood Search for Data Clustering. Baghdad Sci J. 2022; 19: 409- 421. <https://doi.org/10.21123/bsj.2022.19.2.0409>

طريقة جديدة في تحسين الجاذبية لتلخيص النص العربي بالاستخلاص

مصطفى جاسم هادي، ايداد روضان عباس ، اسامة يونس فاضل

قسم علوم الحاسوب، الجامعة التكنولوجية، بغداد، العراق.

الخلاصة

يحاكي نظام تلخيص النص التلقائي كيفية تلخيص البشر من خلال اختيار الجمل الأكثر أهمية في النص المصدر. ومع ذلك ، أصبحت تعقيدات اللغة العربية صعبة للحصول على المعلومات بسرعة وفعالية. يتمثل العيب الرئيسي في الأساليب التقليدية في أنها مقيدة بشكل صارم (خاصة بالنسبة للغة العربية) من خلال دقة وظائف ميزات الجملة ومخططات الترتيب وحسابات التشابه. من ناحية أخرى ، تتميز مناهج البحث المسماة metaheuristic بميزة تتسامح مع عدم الدقة ، وتحصل على نتائج محظورة ، ولا تلتزم بشكل صارم بالقيود المذكورة أعلاه. استخدمت هذه الورقة خوارزمية تحسين الجاذبية (GOA) ، وهي منهج ماورائي قوي قائم على قانون الجاذبية ، لمواجهة التحدي المتمثل في تلخيص النصوص العربية. يتم اشتقاق الوظيفة الموضوعية لخوارزمية GOA بناءً على أهمية الجملة ، مثل طولها ودرجة التشابه والموضع وتكرار المصطلح الإحصائي وملكية الكيان المحدد. تم استخدام مجموعة الملخصات العربية من Essex (EASC) لتقييم الطريقة المقترحة وتم قياسها من خلال الاستدعاء الموجه نحو الاسترداد لتقييم التلاعب (ROUGE). حقق النهج المقترح 68.04٪ استرجاع ، 58.49٪ دقة ، 60.05٪ قياس F1 باستخدام ROUGE-1 ، أعلى من الملخصات القياسية والنهج المسماة metaheuristic.

الكلمات المفتاحية: تلخيص تجريدي، تلخيص استخلاصي، تلخيص النص العربي، مخطط التشابه، خوارزمية تحسين الجاذبية.