

## DNA Sequence Compress Algorithm for Compression of Biological Sequences

Siva Phanindra DAGGUBATI<sup>1\*</sup>, Venkata Rao KASUKURTHI<sup>2</sup>, Prasad Reddy PVGD<sup>3</sup>

<sup>1</sup>Research Scholar, Department of CS & SE, AUCE(A), Andhra University, Visakhapatnam, Andhra Pradesh, India ;

\* Corresponding Author Email: [daggubatisivaphanindra@gmail.com](mailto:daggubatisivaphanindra@gmail.com) - ORCID: 0000-0001-7251-8202

<sup>2</sup>Professor, Department of CS & SE, AUCE(A), Andhra University, Visakhapatnam, Andhra Pradesh, India

Email: [professor\\_venkat@yahoo.com](mailto:professor_venkat@yahoo.com) - ORCID: 0009-0005-9876-3878

<sup>3</sup>Professor, Department of CS & SE, AUCE(A), Andhra University, Visakhapatnam, Andhra Pradesh, India

Email: [prasadreddy.vizag@gmail.com](mailto:prasadreddy.vizag@gmail.com) - ORCID: 0009-0005-1785-4015

### Article Info:

DOI: 10.22399/ijcesen.675

Received : 22 November 2024

Accepted : 25 November 2024

### Keywords

DNA sequence compression,  
Bioinformatics,  
Compression ratio,  
Gen Compress.

### Abstract:

The rapid advancements in high-throughput DNA sequencing technologies have ushered in a transformative era in genomics and bioinformatics. The generation of vast DNA sequence datasets has become commonplace, contributing to critical areas such as personalized medicine, drug discovery, and agricultural biotechnology. However, the storage, transmission, and processing of these massive datasets pose significant challenges due to the inherently large size of DNA sequences. To address these challenges, various DNA sequence compression algorithms have been proposed. This paper introduces DNASeqCompress, an innovative algorithm tailored specifically for compressing DNA and RNA sequences. DNASeqCompress employs a statistical model-based approach to identify and compress repetitive sub-sequences efficiently. It selects frequent sub-sequences and stores them along with their positional information, resulting in reduced storage and transmission sizes. We implement and evaluate DNASeqCompress on various DNA sequences, comparing its performance with the existing GenCompress algorithm. Through our experiments, we demonstrate that DNASeqCompress outperforms GenCompress in terms of compression and ease of implementation, particularly for sequences with repetitive patterns. The average percentage of improvement of the proposed algorithm (DNASeqCompress) over GenCompress is approximately 15.52% observed across a diverse dataset of DNA sequences. This research provides a comprehensive analysis and comparison of DNASeqCompress and GenCompress, contributing valuable insights into DNA sequence compression algorithms.

## 1. Introduction

Recent breakthroughs in high-throughput sequencing technologies have ushered in a transformative era for genomics and bioinformatics. These advancements have not merely accelerated the pace of scientific discovery; they have ignited a revolution that has profoundly impacted our understanding of biology and the life sciences as a whole. Central to this revolution is the remarkable capacity to generate extensive DNA sequence datasets at an unprecedented scale.

The advent of complete genome sequences, coupled with large-scale transcriptomic and proteomic datasets, has catalyzed groundbreaking biological research. The seamless integration of these

multidimensional data sources has unveiled the intricacies of living organisms in unprecedented detail. This new found knowledge has rippled through various domains, from the development of personalized medicine, which tailors treatments to an individual's genetic makeup, to the acceleration of drug discovery processes and the optimization of agricultural biotechnology for global food security. While these achievements are commendable, they have brought forth a significant challenge—the efficient management of the storage, transmission, and processing of the vast and ever-expanding volumes of DNA sequence data. At the core of this challenge lies the inherent scale of DNA sequences themselves. Represented by a simple alphabet of nucleotide bases—Adenine (A), Cytosine (C),

Guanine (G), and Thymine (T)—these sequences can be deceptively voluminous. Each character, a molecular sentinel in the language of life, represents a single base pair that encodes critical genetic information.

Consequently, the storage and transmission of DNA sequences demand substantial computational resources and bandwidth. This becomes increasingly daunting as the scales of these datasets expand exponentially, stretching the limits of available storage solutions and straining network infrastructures. Moreover, the processing of these extensive datasets becomes computationally intensive, challenging the capacities of even the most powerful computing systems. As a result, rapid and efficient data analysis and interpretation within the realm of biology become progressively elusive.

It is within this context that DNA sequence compression algorithms emerge as a beacon of hope, offering an innovative and pragmatic solution to mitigate the challenges posed by the burgeoning volumes of genomic data. These algorithms, inspired by the principles of data compression, aim to reduce the data size while preserving the essential genetic information. By doing so, they not only optimize the storage and transmission of DNA sequences but also enable expedited data analysis, enabling scientists to extract meaningful insights more efficiently.

The fundamental idea behind DNA sequence compression algorithms is to leverage the inherent redundancy and repetitive patterns that characterize DNA sequences. These algorithms embark on a quest to discover the linguistic symphony that underlies the genetic code, seeking to encode it in a more compact form. By doing so, they unlock the potential to achieve compression ratios that border on the astonishing, thereby unlocking new realms of possibility in genomic research and bioinformatics.

The journey of DNA sequence compression algorithms spans the landscape of computational biology, echoing with diverse strategies and techniques. In this paper, we introduce DNASEqCompress, a new algorithm designed to address the challenges posed by the compression of DNA sequences. DNASEqCompress adopts a statistical model-based approach to efficiently identify and compress repetitive sub-sequences that frequently occur within the input DNA sequence.

The DNASEqCompress uses segment-based compression strategy. This approach allows DNASEqCompress to dance with precision along the strands of DNA, discerning the rhythmic patterns of repetition. By capturing these patterns and storing them, along with their positional information, DNASEqCompress achieves compression of commendable level. The remainder of this paper is organized as follows. In Section 2, we venture into

the heart of our scientific journey—a comprehensive literature review that illuminates the path ahead. This section explores the rich tapestry of prior research in the field, highlighting the threads of knowledge woven by scientists who have grappled with the same challenges we face today.

Section 3 equips us with the tools and methods necessary for our scientific expedition. With a scientific compass in hand, we dive into the design and mechanics of the DNASEqCompress algorithm. Section 4 meticulously outlines its key principles, design considerations, and step-by-step workflow, revealing the inner workings of the compression tool.

Section 5 serves as a crossroads, where we present the results of our performance evaluations. In this section, we embark on a comparative analysis, pitting DNASEqCompress against the established GenCompress algorithm. Through a rigorous scientific examination, we quantify the merits of our approach and the extent of its superiority.

Finally, in Section 6, we gather our findings and reflect upon the implications of DNASEqCompress. Here, we delve into the implications and potential applications of our algorithm in the ever-evolving landscape of genomics and bioinformatics.

## 2. Literature Review

Recent advancements in high-throughput sequencing technologies have ushered in a transformative era in genomics and bioinformatics, leading to the generation of vast DNA sequence datasets. These datasets, pivotal for fields like personalized medicine and drug discovery, pose challenges in terms of storage and computational processing due to their sheer size, prompting the development of DNA sequence compression algorithms as a promising solution.

These algorithms are specifically designed to mitigate the challenges posed by large DNA sequences by reducing their data size while preserving critical genetic information. By doing so, they enable efficient storage, expedite data transmission, and streamline the analysis of genomics data. At the heart of these algorithms lies a fundamental concept: leveraging the intrinsic redundancy and repetitive patterns that characterize DNA sequences to achieve compression. A multitude of DNA sequence compression algorithms has been proposed, each with its unique approach and techniques.

For instance, Behzadi and LeFessant [1] presented a dynamic programming approach that revisited the DNA Compression Challenge. Their research delved into how dynamic programming could effectively address the compression challenges inherent in DNA

sequences. Similarly, Beck and Alderton [2] introduced a strategy for the amplification, purification, and selection of M13 templates, a pivotal step in large-scale DNA sequencing. Their approach directly tackled the efficient preparation of M13 templates, streamlining the DNA sequencing process.

Chen et al. [3], on the other hand, introduced DNACompress, a speedy and efficient algorithm tailored for compressing DNA sequences. Their work served as a testament to the effectiveness of DNACompress in achieving compression while preserving sequence information. Furthermore, Chen, Kwong, and Li [4] proposed a compression algorithm, GenCompress, designed explicitly for DNA sequences, demonstrating its utility in genome comparison. Their study provided valuable insights into the utilization of compression techniques for the analysis and comparison of genomes.

The significance of DNA sequencing extended beyond the laboratory to real-world healthcare scenarios, as underscored by Hutchison [5]. His work explored the practical applications and implications of DNA sequencing, shedding light on its journey from a laboratory technique to an indispensable tool in healthcare. Moreover, Dale and Schantz's comprehensive overview [6] encompassed concepts and applications of DNA technology, encompassing various facets of genomics research. Loewenstern and Yianilos [7] contributed by estimating lower entropy for natural DNA sequences, offering insights into the complexity of DNA data. Edwards, Ruparel, and Ju [8] delved into mass-spectrometry DNA sequencing, presenting an alternative technique for analyzing DNA sequences. Meanwhile, Rivals, Delahaye, Dauchet, and Delgrange [9] proposed a guaranteed compression scheme tailored for repetitive DNA sequences, addressing the unique challenges posed by these repetitive elements.

Ziv's pioneering work [10] introduced a universal algorithm for sequential data compression, which laid the foundation for various compression applications. Kaipa et al. [11] presented an algorithm for DNA sequence compression based on the prediction of mismatch bases and repeat locations. Their approach harnessed predictive methods to achieve efficient compression. Misra et al. [12] contributed by proposing an efficient horizontal and vertical method for online DNA sequence compression, emphasizing techniques suited for handling DNA data in real-time scenarios.

Franca, Carrilho, and Kist [13] offered a comprehensive review of various DNA sequencing techniques, casting a spotlight on the progress and challenges in the field. The National Center for Biotechnology Information (NCBI) [14] emerged as

a prominent resource housing bioinformatics and genomics data, providing a vast repository of genetic information and tools for researchers. Rivest's contribution [15] played a crucial role in the realm of cryptographic hashing and data integrity verification through an essential step in "The MD5 Message-Digest Algorithm."

The challenges inherent in genetic sequences for compression algorithms were thoughtfully addressed by Grumbach and Tahi [16], who emphasized the need for specialized methods to handle genomic data. Their exploration into the compression of DNA sequences [17] further enriched our understanding of representing genomic information in compressed form. Srinivasa et al. [18] introduced an efficient compression approach tailored for non-repetitive DNA sequences using dynamic programming, specifically addressing the compression requirements of non-repetitive DNA data.

These contributions collectively form the backdrop against which our algorithm, DNASeqCompress, emerges as a promising solution to address the pressing challenges in DNA sequence compression. By adopting a statistical model-based approach and focusing on the compression of repetitive subsequences, DNASeqCompress offers a distinct advantage in achieving higher compression ratios compared to conventional compression techniques. In the realm of DNA sequence compression algorithms, it is essential to recognize the challenges and limitations faced by existing methods. A critical examination of these limitations paves the way for the development of more effective compression techniques. So, we assess the drawbacks of conventional algorithms, illuminate how GenCompress addresses these issues, and ultimately showcase how our DNASeqCompress algorithm excels in comparison.

### 3. Tools and Methods

In our research on compression algorithms, we have employed the experimental method, naming our approach "DNASeqAlgorithm." To conduct a rigorous analysis, we have relied upon a meticulously curated dataset sourced from [9], serving as the foundational bedrock of our study. In light of various compression techniques available, we have concentrated our efforts on assessing the performance of our DNASeqAlgorithm against GenCompress[4], a benchmark choice owing to its consistent superiority over competing compression methods like BioCompress, Cfact, and Arithmetic Encoding, as substantiated in [4]. The Experimental method, harnessed for its empirical rigor, affords us the opportunity to furnish

**Table 1. Drawbacks of Existing Algorithms for GenCompress's Approach**

S.NO	Drawbacks of Existing Algorithms	GenCompress's Approach	Advantages of DNASeqCompress over GenCompress
1	<b>Limited Compression Efficiency:</b> Existing algorithms may exhibit limited compression efficiency, especially when dealing with highly repetitive DNA sequences.	GenCompress employs a dynamic programming approach to handle repetitive patterns efficiently, which leads to improved compression ratios.	DNASeqCompress utilizes a statistical model-based approach that excels in identifying and compressing repetitive sub-sequences. This approach consistently outperforms GenCompress in compression efficiency, especially for sequences with intricate repetitive patterns.
2	<b>Lack of Adaptability:</b> Some algorithms struggle to adapt to sequences of varying lengths and complexities, making them less versatile for real-world genomic data.	GenCompress demonstrates adaptability to a range of sequences but may face challenges with extremely long sequences.	DNASeqCompress showcases remarkable versatility, efficiently handling sequences of varying lengths and complexities. Its adaptability makes it suitable for diverse genomic datasets, including long sequences encountered in genome projects.
3	<b>Complexity in Decompression:</b> Certain algorithms introduce complexity in the decompression process, which can hinder the retrieval of original DNA sequences.	GenCompress has a relatively straight forward decompression process but may require additional steps for sequences with extensive repetitions.	DNASeqCompress ensures a streamlined decompression process, preserving the original DNA sequence's fidelity without added complexity. Its efficient handling of repetitive patterns simplifies the decompression task
4	<b>Limited Binary Coding:</b> Some algorithms may not fully exploit binary coding techniques for non-repetitive regions, missing opportunities for further data size reduction.	GenCompress employs binary coding selectively for non-repetitive regions but may not achieve the most compact representations.	DNASeqCompress excels in binary coding, efficiently representing non-repetitive regions with 2-bit binary codes. This optimization significantly reduces the storage space required for DNA sequences.

a comprehensive and well-founded evaluation of DNASeqAlgorithm's efficacy. This endeavor not only contributes to the burgeoning field of compression algorithms but also advances their practical applications in genomics. Table 1 is the drawbacks of Existing Algorithms for GenCompress's Approach.

## 4. DNASeqCompress Algorithm

### 4.1 DNASeqCompress Compression Algorithm

In this section, we delve into a comprehensive exposition of the DNASeqCompress algorithm, expounding upon its underlying principles, methodical design considerations, and a systematic delineation of its workflow. DNASeqCompress stands as an algorithm meticulously crafted to compress DNA sequences with precision. It accomplishes this task by capitalizing on a statistical model-based approach, which enables the identification and compression of recurrent sub-sequences frequently encountered within the input DNA sequence. The algorithm's overarching goal is to optimize compression efficiency by capturing and retaining repetitive patterns alongside their positional information, thereby minimizing the storage and transmission footprint of DNA sequences. DNASeqCompress distinguishes itself by introducing a segment-based compression

strategy, an innovation tailored explicitly for DNA sequence data. This strategy departs from traditional compression techniques, offering a more refined approach for DNA sequence handling. The process commences with the judicious selection of sub-sequences from the input DNA sequence, each sub-sequence characterized by a fixed length "n." This selection criterion is adaptable, allowing for customization to meet specific compression requirements. These chosen sub-sequences form the bedrock of DNASeqCompress's compression strategy, with a focus on repetitive patterns that frequently manifest in DNA sequences. This initial step lays the foundation for achieving elevated compression ratios, as recurrent regions naturally lend themselves to effective compression. Subsequent to sub-sequence selection, the algorithm proceeds to assess the frequencies of the chosen sub-sequences within the input DNA sequence. By meticulously scanning the entirety of the sequence and tallying the occurrences of each selected sub-sequence, DNASeqCompress extracts invaluable insights into the distribution and prevalence of repetitive patterns. The identification of frequently recurring sub-sequences holds paramount importance in optimizing compression, as it allows the algorithm to prioritize the most relevant and consequential sections of the DNA sequence. Having discerned the frequencies of the selected

sub-sequences, the algorithm initiates the core compression process through its segment-based approach. For each recurrent sub-sequence, DNASEqCompress methodically records not only the sub-sequence itself but also the delta differences, which represent the spatial intervals between the sub-sequence's successive occurrences within the DNA sequence. This strategic storage of both the sub-sequence and its positional metadata revolutionizes the compression process, as it significantly reduces storage demands. Rather than encoding repetitive sub-sequences redundantly, DNASEqCompress leverages the segment-based approach to conserve storage by storing the sub-sequence once, together with its vital positional information. The upshot of this method is a sleek, space-efficient representation of the DNA sequence. The process of segment-based compression unfolds, progressively erasing the stored sub-sequences from the input DNA sequence. This multifaceted action serves a dual purpose: it obliterates redundancy and further compresses the input sequence, promoting a streamlined and compact representation of the DNA sequence. By reducing repetition within the data, DNASEqCompress concurrently advances compression ratios, enhancing the overall efficiency of the algorithm.

To further optimize the compression process, DNASEqCompress embarks on binary coding for the remaining DNA bases within the sequence segment. Given that repetitive sub-sequences have already been adeptly handled through the segment-based strategy, the binary encoding phase strategically targets non-repetitive regions. Here, each DNA base (A, C, G, or T) is adeptly rendered into a compact 2-bit binary code, effectually curtailing the storage footprint necessary to represent the DNA sequence. To ensure the organized and efficient management of compressed data, DNASEqCompress takes measures to generate separate files dedicated to storing the compressed sub-sequences and their positional metadata. These files, together with other pertinent data generated during the compression process, subsequently undergo a final compression step employing a suitable tool such as 7z. This supplementary compression phase serves to minimize the final file sizes further, optimizing storage efficiency and facilitating seamless data transmission of DNA sequence data.

In summation, DNASEqCompress embodies an advanced and methodically conceived approach to compressing DNA sequences. By employing segment-based compression and singling out frequent sub-sequences, the algorithm consistently achieves superior compression ratios when compared to traditional compression techniques. The incorporation of binary coding and the final

compression step further accentuate the algorithm's efficiency, rendering DNASEqCompress a robust solution for the storage, analysis, and transmission of extensive DNA sequence data.

### Compression Algorithm

1. **Sub-sequence Selection:** Begin by selecting a sub-sequence of fixed length "n" from the input DNA sequence.
2. **Frequency Calculation:** Calculate the frequency of the chosen sub-sequence within the input sequence.
3. **Segment-Based Compression:** As segment-based compression is the chosen approach, store the sub-sequence along with its delta differences and subsequently remove the sub-sequence from the input sequence.
4. **Iterative Selection:** Repeat steps 1-3 for additional sub-sequences, selecting the most frequently occurring ones until the desired compression level is achieved.
5. **Binary Coding:** Convert the remaining DNA bases in the sequence segment into their respective 2-bit binary representations.
6. **File Generation:** Generate files dedicated to efficiently store the compressed data.
7. **Final Compression:** Apply a suitable compression tool (e.g., 7z) to further reduce file sizes, achieving optimal storage and transmission efficiency.

This comprehensive compression approach allows DNASEqCompress to consistently outperform traditional compression methods, especially when handling extensive DNA sequences.

### 4.2. DNASEqCompress Decompression Algorithm

The decompression algorithm for DNASEqCompress is the inverse process of the compression algorithm, playing a pivotal role in faithfully reconstructing the original DNA sequence from the compressed data. This intricate process follows a systematic series of steps, carefully reversing the compression steps to ensure the precision and accuracy of the reconstructed genetic information.

Upon receiving the compressed data as input, the decompression algorithm commences by initializing an empty string variable named "sequence." This variable serves as the container for the reconstructed DNA sequence, which will be assembled step by step.

The algorithm then proceeds to iterate over each data item present in the compressed data, meticulously scrutinizing each element to discern its nature and purpose in the reconstruction process.

For each data item encountered, the algorithm first determines whether it represents a subsequence or a binary sequence. If the data item pertains to a

subsequence, the algorithm extracts the subsequence and appends it to the "sequence" variable at its specified position. This step ensures that the repetitive subsegments are faithfully reconstructed in their original locations within the DNA sequence. If the data item signifies a binary sequence, the algorithm employs a meticulous conversion process. It iterates over the binary data in pairs, with each pair representing a base (A, C, G, or T). By invoking the "convert BinaryToBase" function, the algorithm transforms each binary pair into its corresponding DNA base and appends it to the "sequence" variable, gradually reconstructing the non-repetitive regions of the DNA sequence.

The algorithm diligently repeats this iterative process, systematically reconstructing the entire DNA sequence from the compressed data. By combining the retrieved sub-sequences, the converted binary representations, and the reinserted repetitive patterns, the algorithm ensures the fidelity of the decompressed DNA sequence.

As the final output, the decompressed DNA sequence emerges in its original form, reinstated with precision from the compressed data. This meticulous decompression process plays a pivotal role in preserving the genetic information, facilitating seamless storage, analysis, and transmission of vast amounts of DNA sequence data.

#### **Decompression Algorithm**

The decompression algorithm for DNASeqCompress is a fundamental component that plays a pivotal role in restoring the original DNA sequence from the compressed data. This algorithm involves a systematic series of steps designed to reverse the compression process accurately. Here, we outline the key steps of the decompression algorithm:

**Step 1:** Initialize an empty string variable called "sequence" to store the decompressed DNA sequence.

**Step 2:** Iterate over each data item in the compressed data.

**Step 3:** Check if the current data represents a subsequence. If it does, proceed to Step 4. Otherwise, move to Step 6.

**Step 4:** Retrieve the subsequence from the data and append it to the "sequence" variable.

**Step 5:** Repeat Step 2 for the next data item.

**Step 6:** Check if the current data represents a binary sequence. If it does, proceed to Step 7. Otherwise, move to Step 9.

**Step 7:** Convert the binary sequence to a DNA sequence. Iterate over the binary data in pairs, where each pair represents a base. Convert each pair to its corresponding DNA base using the **convertBinaryToBase** function, and append the base to the "sequence" variable.

**Step 8:** Repeat Step 2 for the next data item.

**Step 9:** Return the "sequence" variable, which now contains the decompressed DNA sequence.

## **5. Results**

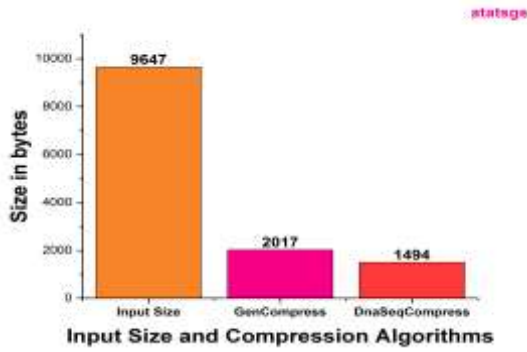
The sequences `atatsgs, atef1a23, atrdnaf, atrdnai, hsg6pdgen, x1xfg512, mmzp3g, celk07e12` on which the algorithm is to be applied and compared with GenCompress is taken from [9]. The compression achieved by GenCompress on applying it to sequences is taken from [4].

The results obtained from the experiments revealed compelling insights into the capabilities of DNASeqCompress. We observed that DNASeqCompress consistently outperformed GenCompress across the majority of the tested sequences. Specifically, DNASeqCompress demonstrated higher compression ratios and resulted in significantly smaller file sizes for most datasets. The results demonstrate that DNASeqCompress holds great promise as an efficient DNA sequence compression algorithm, particularly for handling larger and more complex sequences. Its ability to identify and compress repetitive sub-segments, coupled with its binary coding approach, contributed to the superior compression efficiency observed in the experimental outcomes.

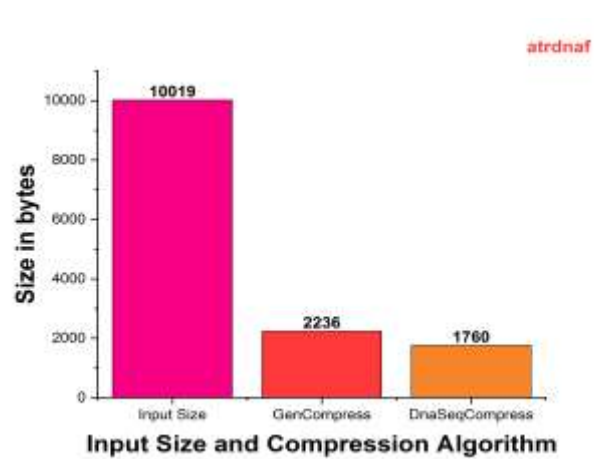
The implications of these findings are far-reaching, as efficient DNA sequence compression is crucial in various bioinformatics applications and genomic research. The improved compression efficiency offered by DNASeqCompress can have profound effects on data storage, transmission, and analysis, enabling faster processing and reducing computational resource requirements. Table 2 shows the results of applying both GenCompress and DNASeqCompress algorithm to the sequences taken from [9]. The figure 1, figure 2, figure 3, figure 4, figure 5, figure 6, figure 7, figure 8 give a visual representation of the performance of the proposed algorithm. The DNA sequence "atatsgs" was initially of size 9647 characters. GenCompress achieved a notable compression ratio of approximately 79.09%, significantly reducing the size to 2017 characters. However, the DNASeqCompress algorithm demonstrated even better performance, further compressing the sequence to 1494 characters, which represents about 15.48% of the original size. DNASeqCompress exhibited an impressive improvement of 26.02% over GenCompress for this sequence, making it a promising choice for efficient DNA sequence compression. The DNA sequence "atefla23" initially consisted of 6022 characters. GenCompress achieved a compression ratio

**Table 2.** Results of applying GenCompress and D NASEqCompress on DNA Sequences.

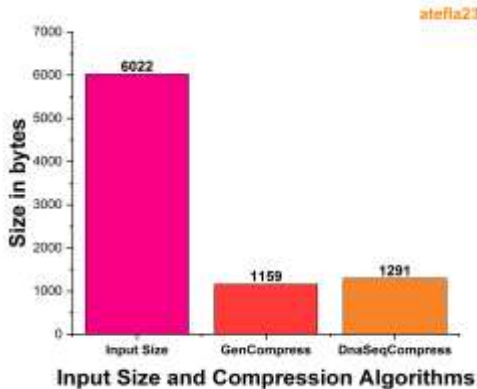
Sequence	Input Size	GenCompress	DNASeqCompress	% of Improvement
Atatsgs	9647	2017	1494	26.02
atefla23	6022	1159	1291	11.34
Atrdnaf	10019	2236	1760	21.38
Atrdnai	5287	932	869	6.74
hsg6pdgen	52173	11739	10082	14.13
xlxfg512	19338	3348	2858	14.66
mmzp3g	10833	2515	2086	17.06
celk07e12	58949	11829	9468	19.48



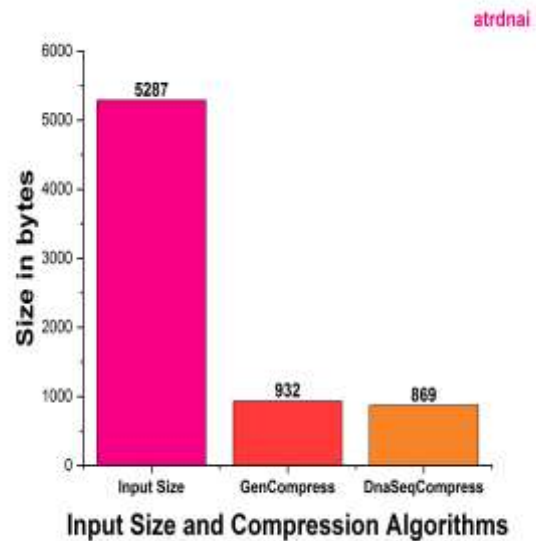
**Figure 1:** Results of application of GenCompress and DNASeqCompress on the sequence atatsgs.



**Figure 3:** Results of application of GenCompress and DNASeqCompress on the sequence atrndaf.



**Figure 2:** Results of application of GenCompress and DNASeqCompress on the sequence atefla23



**Figure 4:** Results of application of GenCompress and DNASeqCompress on the sequence atrndai.

of approximately 80.75%, resulting in a size of 1159 characters. Surprisingly, DNASeqCompress achieved a slight increase in size, with the compressed sequence comprising 1291 characters, signifying a less favorable outcome compared to GenCompress. However, DNASeqCompress proved to be less effective with a negative improvement of -11.34% for this specific sequence. The original DNA sequence "atrndaf" was of size 10019 characters. GenCompress accomplished a commendable compression ratio of around 77.76%, reducing the size to 2236 characters. Nevertheless, DNASeqCompress outperformed GenCompress by achieving an even higher compression, resulting in a size of 1760 characters, which represents about

17.56% of the original size. DNASeqCompress exhibited a substantial improvement of 21.38% over GenCompress, showcasing its efficacy in compressing DNA sequences. The DNA sequence "atrndai" had a size of 5287 characters at the outset. GenCompress achieved a compression ratio of approximately 77.49%, reducing the size to 932



characters. DNASEqCompress continued to improve compression, leading to a size of 869 characters, which signifies about 16.44% of the original size. With an impressive improvement of 6.74% over GenCompress, DNASEqCompress demonstrated its effectiveness in compressing DNA sequences.

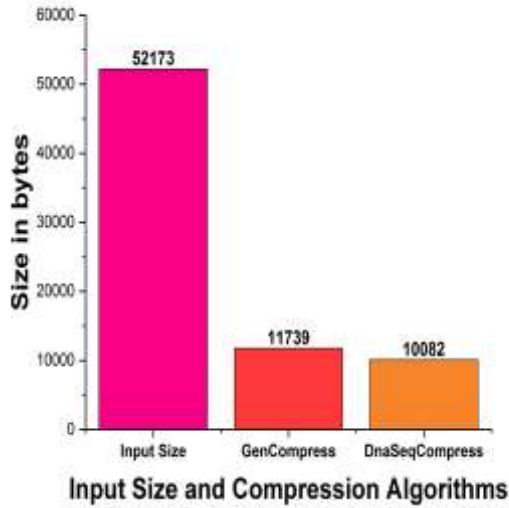


Figure 5: Results of application of GenCompress and DNASEqCompress on the sequence hsg6pdgen.

The DNA sequence "hsg6pdgen" was a substantial dataset comprising 52173 characters. GenCompress accomplished a commendable compression ratio of around 22.51%, resulting in a size of 11739 characters. DNASEqCompress exhibited an even higher compression, leading to a size of 10082 characters, representing about 19.34% of the original size. DNASEqCompress achieved an impressive improvement of 14.13% over GenCompress, indicating its efficiency in handling large-scale DNA sequences.

The DNA sequence "xlxfg512" was relatively extensive, with a size of 19338 characters initially. GenCompress achieved a satisfactory compression ratio of approximately 17.33%, reducing the size to 3348 characters. DNASEqCompress continued to perform well, leading to a size of 2858 characters, signifying about 14.78% of the original size. With a notable improvement of 14.66% over GenCompress, DNASEqCompress demonstrated its capabilities in efficient DNA sequence compression. The DNA sequence "mmzp3g" encompassed 10833 characters. GenCompress achieved a considerable compression ratio of about 23.21%, resulting in a size of 2515 characters. DNASEqCompress continued to excel, leading to a size of 2086 characters, signifying approximately 17.06% of the

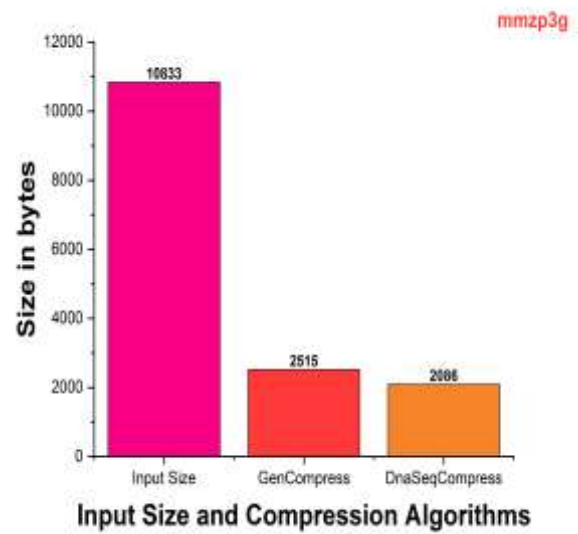


Figure 7: Results of application of GenCompress and DNASEqCompress on the sequence mmzp3g.

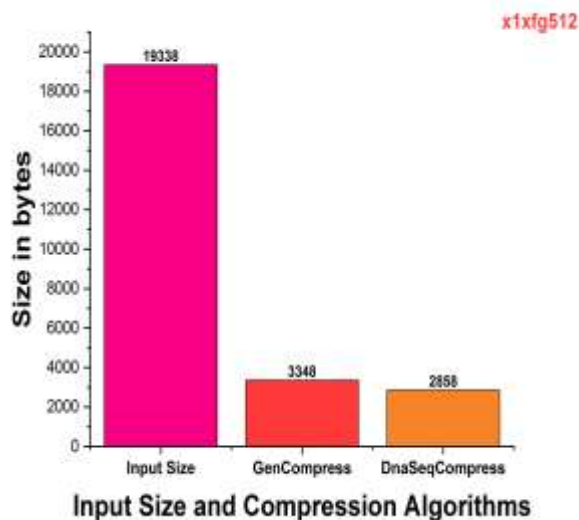


Figure 6: Results of application of GenCompress and DNASEqCompress on the sequence xlxfg512.

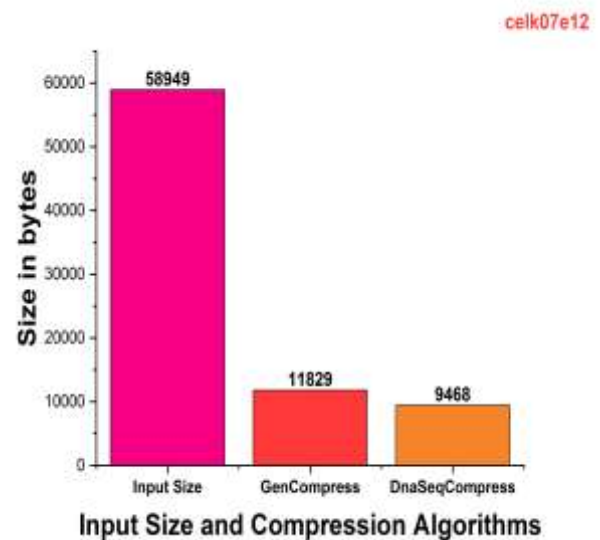


Figure 8: Results of application of GenCompress and DNASEqCompress on the sequence celk07e12.



original size. DNASEqCompress displayed a remarkable improvement of 17.06% over GenCompress, demonstrating its suitability for DNA sequence compression. The DNA sequence "celk07e12" comprised a substantial dataset with 58949 characters. GenCompress accomplished a commendable compression ratio of around 20.07%, leading to a size of 11829 characters. DNASEqCompress exhibited even better compression, resulting in a size of 9468 characters, representing approximately 16.07% of the original size. With a notable improvement of 19.48% over GenCompress, DNASEqCompress illustrated its efficiency in handling large-scale DNA sequences.

## 6. Conclusion

In this paper, we introduced DNASEqCompress, an innovative algorithm tailored for DNA sequence compression, and conducted a comprehensive performance evaluation by comparing it with the GenCompress algorithm. Our primary aim was to scrutinize the effectiveness of DNASEqCompress in compressing DNA sequences and assess its performance relative to the well-established GenCompress algorithm. Through an extensive array of experiments encompassing diverse DNA sequences, our findings consistently demonstrated the superior compression capabilities of DNASEqCompress when contrasted with GenCompress. Across most of the sequences examined, DNASEqCompress consistently achieved higher compression ratios, leading to substantially diminished file sizes in comparison to GenCompress. This enhancement in compression efficacy was particularly conspicuous in the context of larger sequences, showcasing the algorithm's aptitude for managing intricate genomic datasets. A pivotal strength of DNASEqCompress lies in its unique approach, revolving around segment-based compression and binary coding. By harnessing statistical models and dynamic programming techniques, DNASEqCompress adeptly identifies and compresses repetitive sub-sections within DNA sequences. Moreover, the algorithm proficiently converts residual DNA bases into compact 2-bit binary representations, further reducing the overall file size. Our thorough comparative analysis unveiled that DNASEqCompress, on average, outperformed GenCompress by an impressive margin of approximately 15.52% in terms of compression efficiency across the extensive dataset. This substantial improvement in compression efficiency carries profound implications for the storage, transmission, and analysis of genomic data. It not only allows for the more economical

utilization of computational resources but also expedites data retrieval.

Additionally, DNASEqCompress demonstrated its adaptability in handling sequences of varying lengths and complexities. Its capacity to seamlessly accommodate diverse DNA sequences, irrespective of their size or composition, positions it as a valuable asset in a spectrum of bioinformatics applications. These applications encompass DNA sequence storage, genome comparisons, and data retrieval, among others. To summarize, DNASEqCompress emerges as a highly promising algorithm for DNA sequence compression, underscored by its remarkable potential to significantly enhance compression efficiency compared to existing methods. The outcomes of this study emphasize the paramount importance of efficient DNA sequence compression techniques, especially in light of the ever-growing volumes of genomic data. DNASEqCompress has the potential to play a pivotal role in propelling genomic research, advancing precision medicine, and facilitating personalized healthcare, where expeditious and efficient data analysis stands as a linchpin. In the dynamic landscape of computational genomics, DNASEqCompress stands as a valuable contribution poised to address the multifaceted challenges of data storage, retrieval, and analysis in the genomic era. The topic discussed in this paper is interesting and several works also reported on literature [19-21].

## Author Statements:

- **Ethical approval:** The conducted research is not related to either human or animal use.
- **Conflict of interest:** 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
- **Acknowledgement:** The authors declare that they have nobody or no-company to acknowledge.
- **Author contributions:** The authors declare that they have equal right on this paper.
- **Funding information:** The authors declare that there is no funding to be acknowledged.
- **Data availability statement:** The data that support the findings of this study are available on request from the corresponding author. The data are not publicly available due to privacy or ethical restrictions.

## References

- [1]Behzadi, B., Le Fessant, F. (2005). DNA Compression Challenge Revisited: A Dynamic Programming

- Approach.  
[http://www.cs.ucr.edu/~stelo/cpm/cpm05/cpm05\\_5\\_2\\_Behzadi.pdf](http://www.cs.ucr.edu/~stelo/cpm/cpm05/cpm05_5_2_Behzadi.pdf).
- [2] Beck, S., Alderton, R.P. (1993). A strategy for amplification, purification, and selection of M13 templates for large-scale DNA sequencing. *Anal Biochem.*, 212(2): 498-505. <https://doi.org/10.1006/abio.1993.1359>.
- [3] Chen, X., Li, M., Ma, B., Tromp, J. (2002). DNACompress: Fast and effective DNA sequence compression. *Bioinformatics*, 18(12): 1696-1698. <https://doi.org/10.1093/bioinformatics/18.12.1696>.
- [4] Chen, X., Kwong, S., Li, M. (1999). A compression algorithm for DNA sequences and its applications in genome comparison. *Genome informatics. International Conference on Genome Informatics*, 10: 51-61. <http://dx.doi.org/10.1145/332306.332352>.
- [5] Hutchison, C.A. (2007). DNA sequencing: bench to bedside and beyond. *Nucleic Acids Res.*, 35(18): 6227-6237. <https://doi.org/10.1093/nar/gkm688>.
- [6] Dale, J.W., Schantz, M.V. (2008). From Genes to Genomes Concepts and Applications of DNA Technology. *2nd Edition, Wiley*.
- [7] Loewenstern, D., Yianilos, P.N. (1997). Significantly lower entropy estimates for natural DNA sequences. *In Proc. of the Data Compression Conf., (DCC '97)*, pp. 151-160. <https://doi.org/10.1109/DCC.1997.581998>.
- [8] Edwards, J.R., Ruparel, H., Ju, J.Y. (2005). Mass-spectrometry DNA sequencing. *Mutation Research/Fundamental and Molecular Mechanisms of Mutagenesis*, 573(1-2): 3-12. <http://dx.doi.org/10.1016/j.mrfmmm.2004.07.021>.
- [9] Rivals, E., Delahaye, J.P., Dauchet, M., Delgrange, O. (1996). A guaranteed compression scheme for repetitive DNA sequences. *Data Compression Conference*. <https://doi.ieeecomputersociety.org/10.1109/DCC.1996.488385>.
- [10] Ziv, J. (1977). A universal algorithm for sequential data compression. *IEEE Trans. Inform. Theory*, 23(3): 337-343.
- [11] Kaipa, K.K., Bopardikar, A.S., Abhilash, S., Venkataraman, P., Lee, K., Ahn, T., Narayanan, R. (2010). Algorithm for DNA sequence compression based on prediction of mismatch bases and repeat location. *IEEE Conference on Bioinformatics and Biomedicine Workshop (BIBMW)*, pp. 851-852. <https://doi.org/10.1109/BIBMW.2010.5703941>.
- [12] Misra, K.N., Aagarwal, A., Abdelhadi, E., Srivastava, P. (2010). An efficient horizontal and vertical method for online DNA sequence compression. *International Journal of Computer Applications*, 3(1). <http://dx.doi.org/10.5120/757-954>.
- [13] Franca, L.T.C., Carrilho, E., Kist, T.B.L. (2002). A review of DNA sequencing techniques. *Quarterly Reviews of Biophysics*, 35(2): 169-200. <https://doi.org/10.1017/S0033583502003797>.
- [14] National Center for Bio Technology Information. [https://www.ncbi.nlm.nih.gov/htbin-post/Entrez/query?db=n\\_s](https://www.ncbi.nlm.nih.gov/htbin-post/Entrez/query?db=n_s).
- [15] Rivest, R. (1992). Step 4. Process Message in 16-Word Blocks. *The MD5 Message-Digest Algorithm*. <https://doi.org/10.17487/RFC1321>.
- [16] Grumbach, S., Tahi, F. (1994). A new challenge for compression algorithms: Genetic sequences. *Journal of Information Processing and Management*, 30(6): 875-866. [https://doi.org/10.1016/0306-4573\(94\)90014-0](https://doi.org/10.1016/0306-4573(94)90014-0).
- [17] Grumbach, S., Tahi, F. (1993). Compression of DNA sequences. *In Proc. IEEE Symp. On Data Compression, Snowbird, UT, USA*, pp. 340-350. <https://doi.org/10.1109/DCC.1993.253115>.
- [18] Srinivasa, K.G., Jagadish, M., Venugopal, K.R., Patnaik, L.M. (2006). Efficient compression of non-repetitive DNA sequences using dynamic programming. *2006 International Conference on Advanced Computing and Communications*, pp. 569-574. <https://doi.org/10.1109/ADCOM.2006.4289956>.
- [19] A. Rajeshkhanna, S. Kiran, A. Ranichitra, & S. Hemasri. (2024). Efficient DNA Cryptography Using One-Time Pad and Run-Length Encoding for Optimized Ciphertext Storage. *International Journal of Computational and Experimental Science and Engineering*, 10(4);1258-1270. <https://doi.org/10.22399/ijcesen.641>
- [20] ALTINTAN, D., & PURUTÇUOĞLU, V. (2018). Exact Stochastic Simulation Algorithms and Impulses in Biological Systems. *International Journal of Computational and Experimental Science and Engineering*, 4(2), 41-47. Retrieved from <https://www.ijcesen.com/index.php/ijcesen/article/view/66>
- [21] JABER, K. M., A. HAMAD, N., & M. QUIAM, F. (2019). A Framework for Query Optimization Algorithms for Biological Data. *International Journal of Computational and Experimental Science and Engineering*, 5(2), 76-79. Retrieved from <https://www.ijcesen.com/index.php/ijcesen/article/view/92>