



Distributed Bioinformatics Computing System for DNA Sequence Analysis

By Mohammad Ibrahim Khan, Kaushik Deb & Chotan Sheel

Chittagong University of Engineering & Technology, Bangladesh

Abstract- This paper provides an effective design of computing technique of a distributed bioinformatics computing system for analysis of DNA sequences using OPTSDNA algorithm. This system could be used for disease detection, criminal forensic analysis, gene prediction, genetic system and protein analysis. Different types of distributed algorithms for the search and identification for DNA segments and repeat pattern in a given DNA sequence are developed. The search algorithm was developed to compute the number of DNA sequence which contains the same consecutive types of DNA segments. A distributed subsequence identifications algorithm was designed and implemented to detect the segment containing DNA sequences. Sequential and distributed implementation of these algorithms was executed with different length of search segments patterns and genetic sequences. OPTSDNA algorithm is used for storing various sizes of DNA sequence into database. DNA sequences of different lengths were tested by using this algorithm. These input DNA sequences varied in size from very small to very large. The performance of search technique distributed system is compared with sequential approach.

Keywords: *distributed bioinformatics system, DNA sequence, search segments, identify DNA sequences, reported gene sequences.*

GJCST-A Classification: *H.1.1*



Strictly as per the compliance and regulations of:



Distributed Bioinformatics Computing System for DNA Sequence Analysis

Mohammad Ibrahim Khan^α, Kaushik Deb^σ & Chotan Sheel^ρ

Abstract- This paper provides an effective design of computing technique of a distributed bioinformatics computing system for analysis of DNA sequences using OPTSDNA algorithm. This system could be used for disease detection, criminal forensic analysis, gene prediction, genetic system and protein analysis. Different types of distributed algorithms for the search and identification for DNA segments and repeat pattern in a given DNA sequence are developed. The search algorithm was developed to compute the number of DNA sequence which contains the same consecutive types of DNA segments. A distributed subsequence identifications algorithm was designed and implemented to detect the segment containing DNA sequences. Sequential and distributed implementation of these algorithms was executed with different length of search segments patterns and genetic sequences. OPTSDNA algorithm is used for storing various sizes of DNA sequence into database. DNA sequences of different lengths were tested by using this algorithm. These input DNA sequences varied in size from very small to very large. The performance of search technique distributed system is compared with sequential approach.

Keywords: distributed bioinformatics system, DNA sequence, search segments, identify DNA sequences, reported gene sequences.

I. INTRODUCTION

Distributed Computing (DC) provides a cost effective frame work with efficient execution of a solution on multiple computers connected by a network. For distributed Computing (DC), large tasks are divided into smaller problems which can then be executed on multiple computers at the same time independent of each other. The task must be broken up into independent problems to minimize inter-computers communication; otherwise distributed computing will not be effective. Over the past few years, the intermixing of computer science and the complexity of biology has lead to the prosperous field of bioinformatics [1-2] Advances in molecular biology and technology for research have facilitated the process of sequencing of

large portions of genomes in various species. Today computers have made medical research more efficient and accurate, by using parallel and distributed computers and complex biological modeling. Bioinformatics, is one of the newer areas, and has opened our eyes to a whole new world of biology [1].

The fusion of computers and biology has helped scientists learn more about species, especially humans [3-5]. With the aid of the computers, we have learned a great deal about genetics, but there still stand many unanswered questions, that are being researched today. DNA sequence analysis can be a lengthy process ranging from several hours to many days. This paper builds a distributed system that provides the solution for many bioinformatics related applications.

The overall goal of this paper is to build a Distributed Bioinformatics Computing System for genetic sequence analysis of DNA. This system is capable of searching and identifying gene patterns in a given DNA sequence. For the purpose of computing we stored a large no. of DNA sequence using OPTSDNA algorithm [13] and segments is divided two to six consecutive nucleotide [13]. The system was tested for its correctness and efficiency. Different lengths of DNA sequences were used for the consecutive and non-consecutive pattern search to compare the system's response time obtained using single and multiple computers [6]. In addition, different lengths of DNA sequences were also used for the pattern identification to compare its response time observed using a single computer and multiple computers. Several different distributed implementations of search algorithms have been reported in the literature. The characteristics of some of those distributed algorithms are listed in Table 1.

It can be observed that the most of the existing approaches require high performance parallel processors and are not implemented on loosely coupled distributed network. Moreover, most of them require specialized programming language for their implementation on these parallel processors.

The specific objective of the proposed distributed algorithm for analysis of DNA sequences are:

1. Develop an effective distributed DNA sequence analysis algorithms for pattern matching of DNA Gene sequence and sub-sequences identification.
2. Implement them on loosely coupled distributed network such as regular local area network and

Author α : Professor, Department of Computer Science & Engineering (CSE), Chittagong University of Engineering & Technology (CUET), Chittagong, Bangladesh. e-mail: muhammad_ikhancuet@yahoo.com.

Author σ : Associate Professor, Department of Computer Science & Engineering (CSE), Chittagong University of Engineering & Technology (CUET), Chittagong, Bangladesh. e-mail: deb.kaushik99@gmail.com

Author ρ : Research Scholar, Department of Computer Science & Engineering (CSE), Chittagong University of Engineering & Technology (CUET), Chittagong, Bangladesh. e-mail: chotan_cuetcse03@yahoo.com.

wide area network using standard programming language.

This paper is organized in four sections. Section 2 discusses the material and method of algorithm. Section 3 discusses the results and discussion and conclusions included in section 4.

Table 1 : Comparative Study with Existing Approaches

Reference	Algorithm Complexity	Special Purpose Computer Required	No. of Computers Required	Special Language Required	Useful on General Network
[2]	$O(n)$	Yes	Flexible	Yes	No
[9]	$O(n)$	Yes	Not Flexible	Yes	No
[10]	$O(n)$	Yes	Not Flexible	Yes	No
[11]	$O(n^2)$	Yes	Not Flexible	Yes	No
[12]	$O(n)$	Yes	Not Flexible	Yes	No

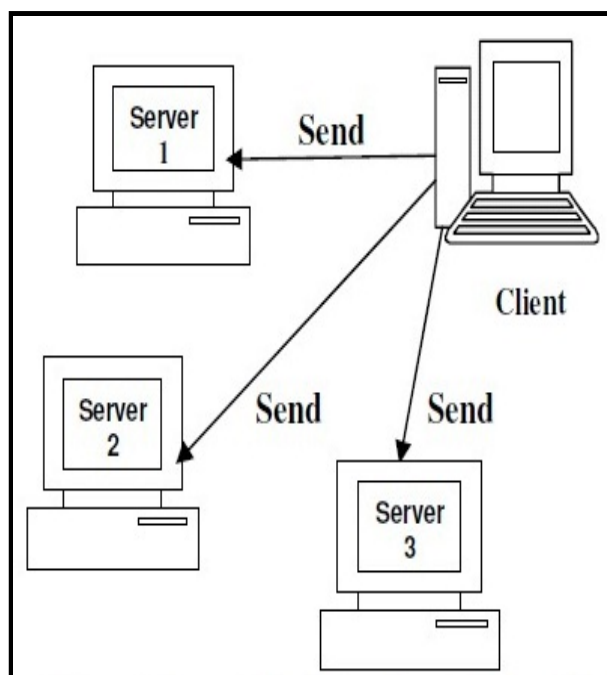


Figure 1 : Layout of the System (Sending of Data)

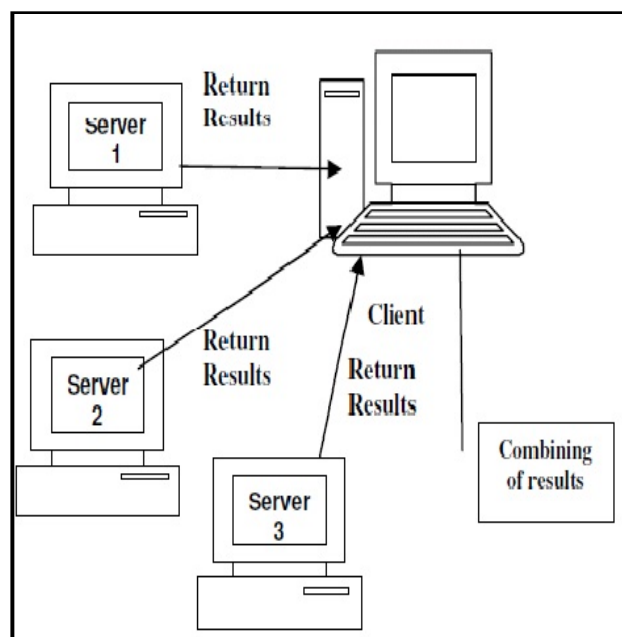


Figure 2 : Layout of System (Returning Results)

a) Application of the Proposed Distributed Algorithm

This distributed Bioinformatics system developed in this paper could be used for disease detection, criminal forensics analysis, genetics systems and protein analysis. Di-let, Triplet, Tetra-let, Pentad-let, Hexed-let repeats formally known as a Di-nucleotide, Tri-nucleotide, Tetra-nucleotide, Pent nucleotide, Hex nucleotide. Repeat occurs when two, three, four, five and six consecutive nucleotides are repeated within a specific region of DNA sequence. These repeats can occur within or between genes. These consecutive repeats are frequently located in genes that encode transcription factors and which are active in the organism development process. Extensive Di-let, Triplet, Tetra-let, Pant-lets, Hex-let repeats are found when a mutation occurs in a gene. This mutation increases the number of occurrences of a particular nucleotide which can lead to a number of neurodegenerative diseases. These diseases include, Huntington's Disease (HD), Fragile X Syndrome, Kennedy's Disease, Myotonic Dystrophy, Spinocerebellar Ataxia Type 1 (SCA1), Dentatorubral Pallidoluysian atrophy (DRPLA), and Fragile X E mental retardation (FRAAXE). In Kennedy's Disease, Huntington's disease, Spinocerebellar Ataxia Type 1, and Dentatorubral Pallidoluysian atrophy, the number of triplet repeats is quite small, in contrast to Fragile X Syndrome, Myotonic Dystrophy, and FRAAXE, where the number of consecutive repeats may be very large, producing alleles that consist of thousands of repeats. These algorithms can help to detect Di-let, Triplet, Tetra-let, Patna-led and Hex-let repeats in gene sequence, and can also search through DNA sequences to identify most frequently occurring repeats.

The proposed distributed algorithms will be able to first identify a DNA sequence Gene pattern in the DNA obtained from the crime scene and then it can search for those patterns in suspects DNA, which will be helpful for criminal investigation, Disease analysis, Gene Sequence Prediction, Human Identification etc. Criminal investigation can now be facilitated by the DNA forensic analysis. Forensic analysis is a process by which two organism's DNA is compared with each other. DNA analysis is effective in finding criminals, because two different individuals will have different DNA sequence. In DNA analysis one can look for matching gene patterns at different locations of the suspect's DNA and the DNA obtained at the crime scene. Gene pattern matching at one, two or three locations in DNA usually aren't enough to associate a suspect with a crime, but gene pattern matches at 5 or more locations in DNA are usually good enough to identify a criminal. Experts believe that DNA forensic technology is more reliable than eyewitnesses, where the odds are fifty-fifty. In DNA analysis one can look for matches based on number of repeating patterns at different locations of the suspect's genome.

II. MATERIALS AND METHOD

The proposed distributed algorithm is based on client server model. For distributed search and identification algorithms on DNA sequence, the proposed framework avoids duplicates computations on server machines. The two input items are provided by the user for pattern search and identification:

1. The DNA sequence which is stored by OPTSDNA algorithm with extend two to six consecutive nucleotides division.
2. Search string DNA subsequences or identification DNA segments (Di-nucleotides to Hex-Nucleotides Segment pattern).

Using OPTSDNA algorithm, the DNA sequence is broken up in X segments where $X = m * p$. Here m = number of storage DNA and p = length of storage nucleotide base. Number of storage DNA is also used as number of servers used in distributed algorithm implementation and length of storage nucleotide base represents the length of pattern for search or identification. In the first step each server gets one segment of data and the required search or identification pattern for carrying out its computation as shown in Figure 1. In addition, an offset value is sent to the server as well to make sure that no two servers are performing the same computation for search or identification. The individual results from each server are sent back to the clients where partial results are combined as shown in Figure 2. The complete details of client and server side interaction are shown in Figure 3. The actual pattern search for a DNA sequence with three servers is shown in Figure 4, where each server starts the match at different Gene chromosome.

Different starting point at various servers guarantees that no comparison for pattern search and identification is performed more than once on any server. The worst case complexity of this distributed search or identification algorithm is $O(L/X)$, where L is the length of DNA sequence and $X = m/p$. In case of Figure 4 value of $X = 1$ because $m = 3$ and $p = 3$. That implies that complete DNA sequence is end to all three servers and the offset for starting the search or identification.

a) Implementation of Distributed Algorithms

A Dot net based client server system was developed for this project[7-8] shown in figure A and figure B. The client and server side logic implementation is given in Figure 3 and figure 4. This framework can distribute the workload across multiple servers as specified by the user. In this paper, a client provides the user input from Graphical User Interface (GUI) and then send this input to one or more server computers as directed by the user (shown figure A and B). The processing option is developed in GUI. When a client selects a processing option such as pattern identification, appropriate input for carrying out a search or identification in a DNA sequence displayed (shown in figure A and B). The client program then sends the input data to multiple servers (as specified by the user). The code at the server executes the desired algorithm and returns its results to the client. The client then receives the results from all the servers and combines to individual results to generate a final output of the processing a shown in figure A and figure B.

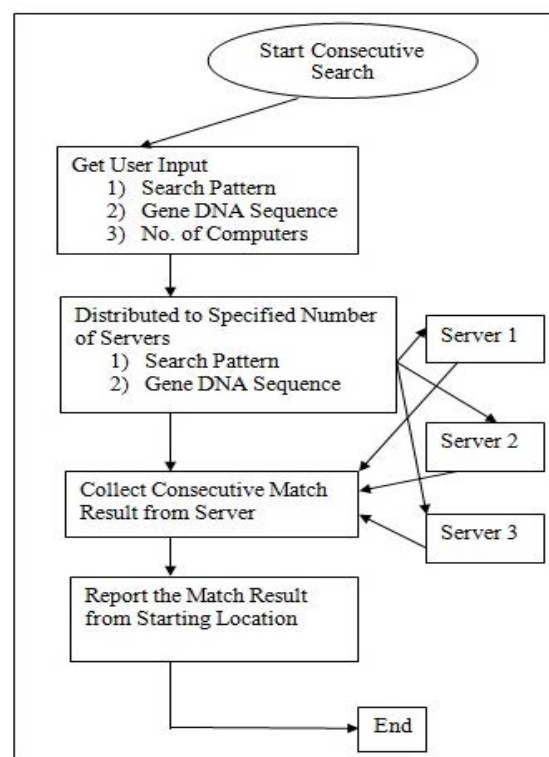


Figure 3 : Flow Diagram for Client Side Implementation

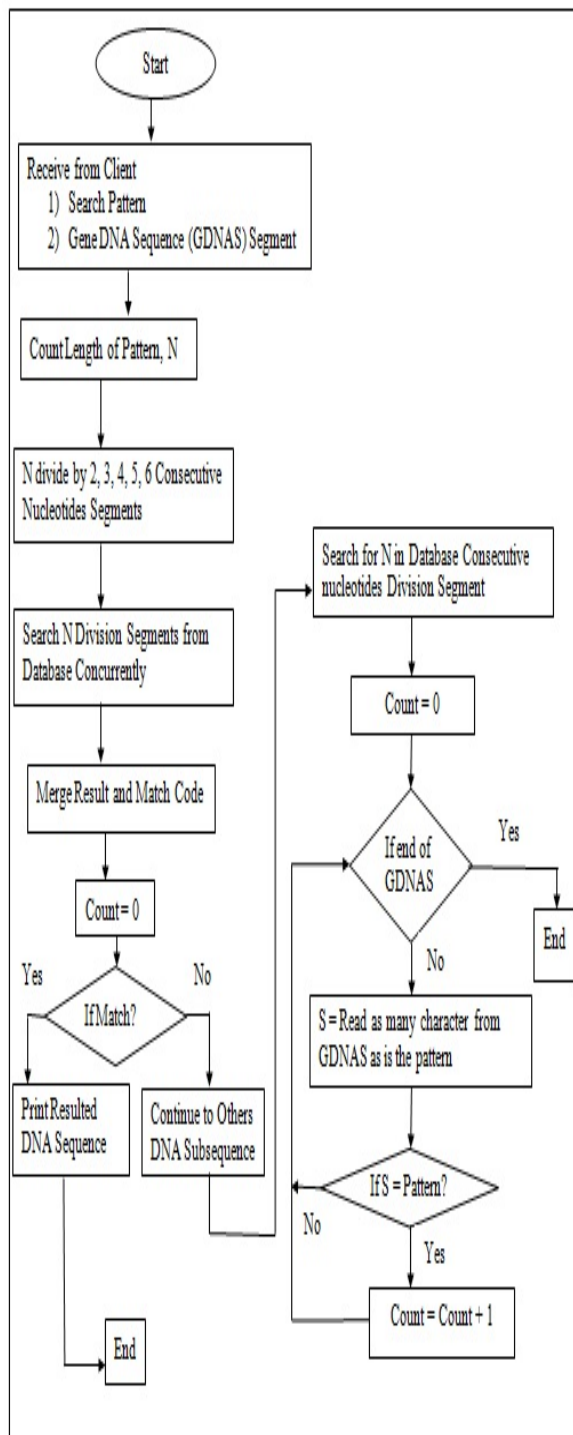


Figure 4 : Flow Diagram for Consecutive Search Pattern from Server

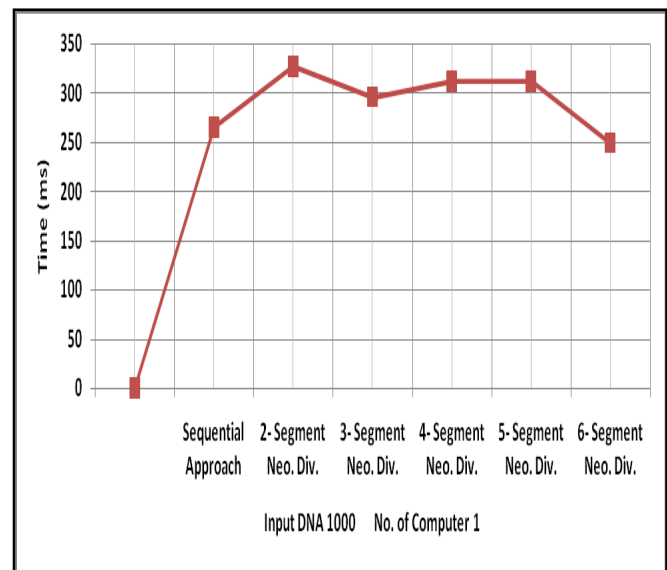


Figure 5 : Effect of Data size on using Single Computer

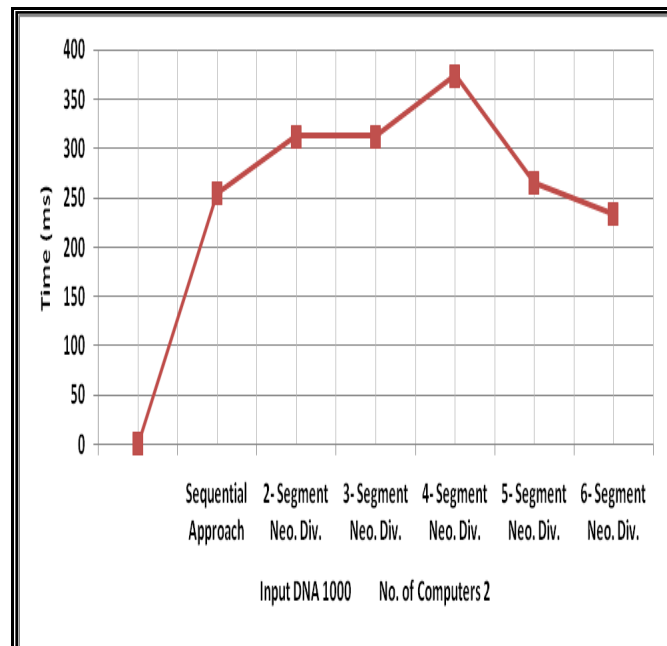


Figure 6 : Effect of Data size on using Two Computers

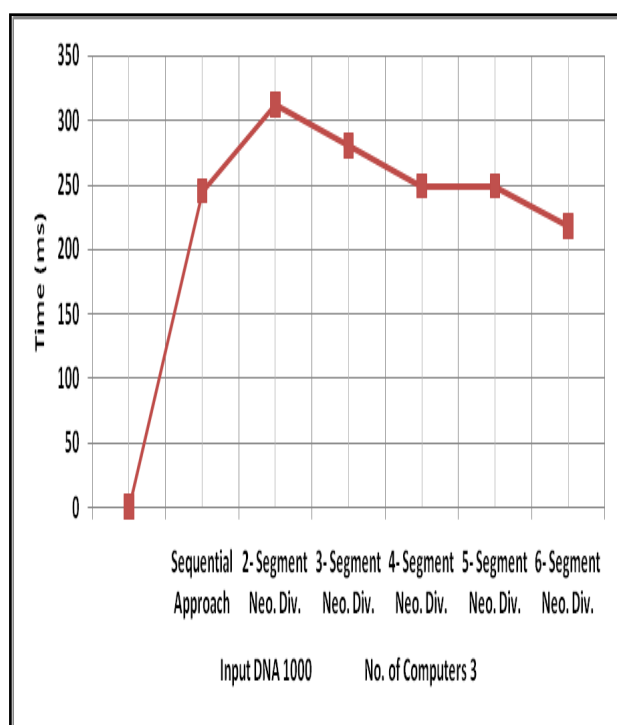


Figure 7 : Effect of Data size on using Three Computers

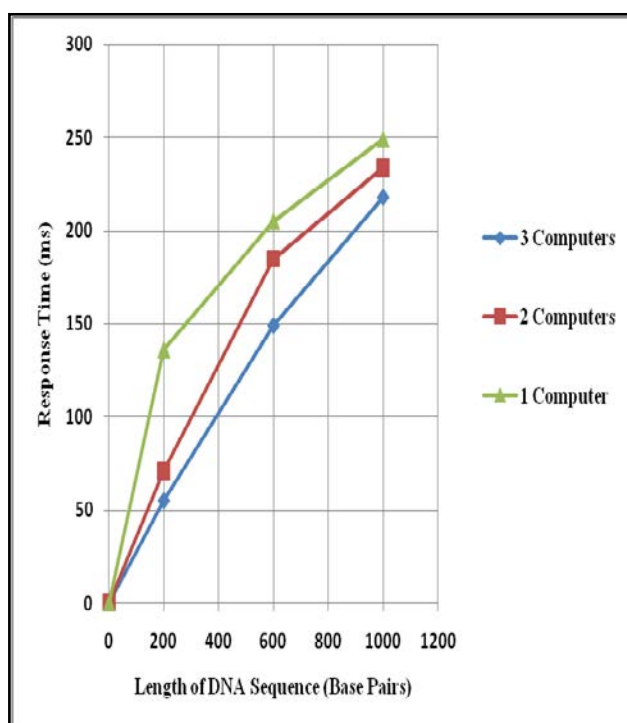


Figure 8 : Effect of Data Size on Computation Time

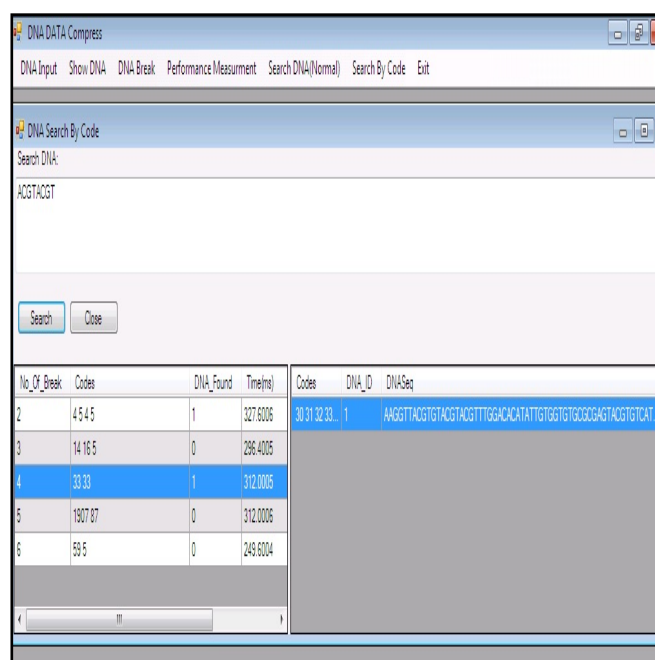


Figure A : Screen Shot for the Search Process by Generating Code

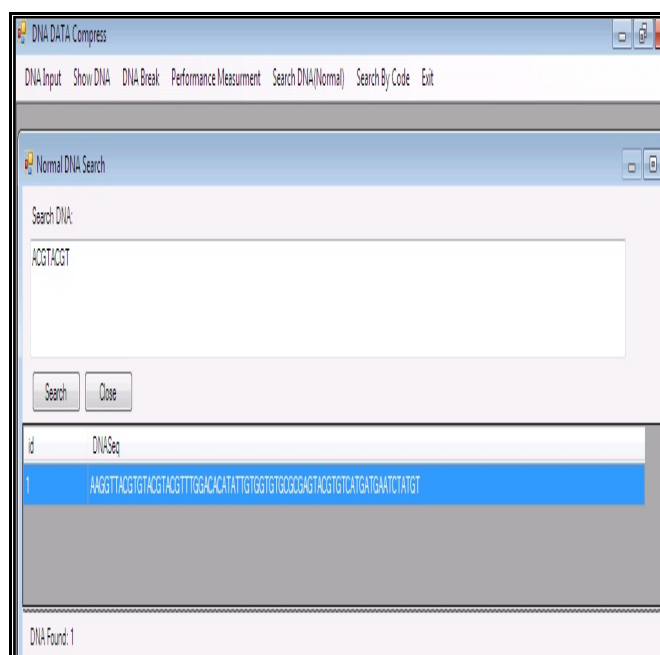


Figure B : Screen Shot for the Search Process by Sequential Approach

III. RESULTS AND DISCUSSION

Sequential and Distributed versions of the algorithms were executed with different patterns of genetic sequences. These sequences were of different sizes ranging from very small to very large. The response times for sequential and distributed versions of the programs were plotted to demonstrate the effectiveness of distributed DNA sequence analysis algorithms. Figure 5, 6, and 7 shows the response time

of consecutive pattern search execution on single machine and multiple machines. The execution time was calculated for DNA sequences of sizes 1 to 1000 sequences. It can be observed that execution time reduces significantly as number of servers increased. Moreover, the improvement in execution time is significant when DNA sequence size is 600 with 3 servers. Figure 5, 6 and 7 shows the response time of consecutive pattern identification execution on single machine and multiple machines. It can be observed that the execution time reduces significantly as number of servers increased.

Similar observation was made for sequential approach consecutive pattern identification algorithm execution shown in Figure 5, 6, and 7. Figure 8 demonstrates how the data size affected the computation time. With a single computer the response time of each gene sequence was significantly more than that of the distributed execution using two and three servers. In addition, rate of growth of execution time is almost linear with three servers as the size of DNA sequence increases.

IV. CONCLUSIONS

As shown in the previous figures, it is clear that as complexity of the algorithm increases the response time also increases. The algorithm for the Pattern Identification was the most complex one and the algorithm for the pattern search was the least complex. It can be seen in Figure 5 the response times for the Pattern Identification were much lower compared to the other two studies shown in Figures 5 and 6. This is due to the fact that more complex algorithms usually involve more steps, which increases the response time. To help get a better understanding of the effects of Distributed Systems on DNA sequences, more DNA sequences of various lengths should be tested. This would provide more data for a larger analysis. It is also recommended that the computers used in the investigation should not exceed the length of the repeat pattern that is being searched or identified, because this will not improve the response time. The complexity of our algorithm is $O(n)$. For computing DNA sequences special purpose of computer is required. Using this algorithm no. of computer required is flexible and special language is required. Our algorithm is useful on general network. So our algorithm is more efficient than previous all. In addition, this system could be interfaced with the Internet, so that all these feature of DNA analysis are accessible to everyone via Web.

REFERENCES RÉFÉRENCES REFERENCIAS

- a) *Chapter or Article Book*
 1. C. Sheel, M. I. Khan, M. I. H. Sarker, T. Alam, "Algorithm for Optimal Storage of a Distributed Bioinformatics System for Analysis of DNA Sequences", International Journal of Computational Bioinformatics And In Silico Modeling, Vol. 2(3), 2013, pp 106 - 109, 2013.
- b) *Dissertation or thesis*
 2. R. Kumar, A. Kumar, and S. Agarwal, "A Distributed bioinformatics Computing System for Analysis of DNA Sequences", IEEE Spectrum, Vol. 7, pp. 358-363, 2007.
 3. Baxevaris, Andreas, and B.F Ouellette, ed., "Bioinformatics: A Practical Guide to the Analysis of Genes and Proteins", 2nd ed. Danvers: John Wiley & Sons, Inc, 2001.
- c) *Article in Conference preceding*
 4. Durban, Richard, S. Eddy, A. Krogh, and G. Mitchison., "Biological Sequence Analysis: Probabilistic Models of Proteins and Nucleic Acids", Cambridge: Cambridge University Press, 1998.
 5. Gusfield, Dan, "Algorithms On Strings, Trees, and Sequences: Computer Science and Computational Biology", New York: The Press Syndicate of the University of Cambridge, 1997.
- d) *Electronic Source*
 6. "OMIM", The OMIM Gene Map, NCBI, 03 Nov, 2005 (<http://www.ncbi.nlm.nih.gov/Omim>).
- e) *Book*
 7. E. Petroustos, "Mastering Visual Basic.NET", 2nd Edition, 2006.
 8. R. Mistry, "Microsoft SQL Server 2008 Management and Administration", McGraw Hill Edition, 2006.
- f) *Patent*
 9. U. Vishkin, "Parallel Pattern Matching in String", Information Control, pp. 91-113, 1985.
 10. C. H. Huang et. al, "Parallel Pattern Identification in Biological Sequences on Clusters", IEEE Transactions on Nan bioscience, pp. 1-6, 2003.
 11. Strumpen, "Coupling Hundreds of Workstations for Parallel Molecular Sequence Analysis", Software Practices and experienced in Parallel Molecular Sequence Analysis, Vol. 25(3) pp. 291-304, 1995.
 12. C. Janiki, R. R. Joshi et. al, "Accelerating Comparative Genomics Using Parallel Computing", Bio-information System, Scientific and Engineering Computing in Silico Biology, Vol. 3(36), pp 123-128, 2003.
- g) *Periodical*
 13. M. I. Khan and C. Sheel, "OPTSDNA: Performance evaluation of efficient distributed bioinformatics system for DNA sequence analysis", Bioinformation - Biomedical Informatics, Vol. 9(16), pp 842 - 846, 2013.

- a) *Chapter or Article Book*
 1. C. Sheel, M. I. Khan, M. I. H. Sarker, T. Alam, "Algorithm for Optimal Storage of a Distributed Bioinformatics System for Analysis of DNA