# A NEW AUTOMATA BASED APPROXIMATE STRING MATCHING APPROACH AND WEB INTERFACE FOR BIOINFORMATICS ALGORITHMS

*Burak KOCA*[*]
*Gıyasettin ÖZCAN*[**]

**Abstract:** In this study, we present a new web interface for major bioinformatics algorithms and introduce a novel approximate string matching algorithm. Our web interface executes major algorithms on the field for the use of computational biologists, students or any other interested researchers. In the web interface, algorithms come under three sections: Sequence alignment, pattern matching and motif finding. In each section, we introduce algorithms in order to find best fitting one for specific dataset and problem. The interface introduces execution time, memory usage and context specific results of algorithms such as alignment score. The interface utilizes emerging open source languages and tools. In order to develop light and user-friendly interface, all parts of the interface coded with Python language. On the other hand, Django is used for web interface. Second contribution of the study is novel A-BOM algorithm, which is designed for approximate pattern matching problem. The algorithm is approximate matching variation of Backward Oracle Matching. We compare our algorithm with popular approximate string matching algorithms. Results denote that A-BOM introduces %30 to %80 short runtime improvement when compared to current approximate pattern matching algorithms on long patterns.

**Keywords:** Bioinformatics, A-BOM, Interface, Approximate Pattern Matching

**Başlıca Biyoinformatik Algoritmaları için Web Ara yüzü ve Yeni Otomat Tabanlı Yaklaşık Desen Eşleştirme Yaklaşımı**

**Öz:** Bu çalışmada temel biyoinformatik algoritmaları için yeni bir web ara yüzü ve özgün bir yaklaşık desen eşleştirme algoritması sunmaktayız. Web ara yüzümüz biyologlar, öğrenciler ve ilgili araştırmacılar için bu alandaki temel algoritmaları çalıştırmaktadır. Web ara yüzünde algoritmalar üç bölüm altında toplanmaktadır: Dizilim hizalama, desen eşleştirme ve motif bulma. Her bir bölümde, özgül veri seti ve problemlere en iyi uyan algoritmanın bulunabilmesi için sonuçlarını karşılaştırabilecekleri algoritmalar sunulmaktadır. Web ara yüzü çalışma süreleri, hafıza kullanımı ve hizalama skoru gibi konuya özel sonuçları sunmaktadır. Ara yüz yeni geliştirilen açık kaynak kodlu dilleri ve araçları kullanmaktadır. Hafif ve kullanıcı dostu bir ara yüz olması amacıyla ara yüzün tüm kısımları Python dili ile kodlanmıştır. Diğer yandan web ara yüzü için Django kullanılmıştır. Çalışmanın ikinci katkısı, yaklaşık desen eşleştirme için tasarlanmış yeni A-BOM algoritmasıdır. Bu algoritma Backwards Oracle Matching algoritmasının yaklaşık varyasyonudur. Algoritmamızı popüler yaklaşık desen eşleştirme algoritmaları ile kıyasladık. Sonuçlar, A-BOM algoritmasını güncel yaklaşık desen eşleştirme algoritmaları ile uzun desenler üzerinde karşılaştırdığımızda, çalışma süresinde %30 ile %80 arasında kısalma gelişimi olduğunu göstermektedir.

**Anahtar Kelimeler:** Biyoinformatik, A-BOM, Ara yüz, Yaklaşık Desen Eşleştirme

[*] Faculty of Engineering, Computer Engineering Department, Gebze Technical University, 16059, Bursa, Turkey
[**] Faculty of Engineering, Computer Engineering Department, Uludag University, 16059, Bursa, Turkey
Corresponding Author: Gıyasettin Özcan (gozcan@uludag.edu.tr)

## 1.  INTRODUCTION

Recent development of the technology have introduced big amount of data in scientific fields. For instance, biologists extract have DNA sequences of organisms, where a human genome consist of nearly 3 billion nucleotides (Pevsner, 2015). In order to the DNA store and extract its features, new computational methods and tools are needed. As a result of this fact, a new discipline, Bioinformatics, has been emerged.

Bioinformatics is an interdisciplinary study field which tries to understand biological information. For this goal, researchers of the computer science and biology introduces various tools and software that can collect, store and process biological data. Particularly, main motivation of computer scientists presenting new algorithms and software tools.

One sub field of Bioinformatics is fast and accurate sequence matching among long nucleotide sequences. The sequence matching studies are important since DNA strand of living organisms are very long.  For instance, human genome consist of nearly 3 billion nucleotides and sequence alignment among the genome sequences are computationally expensive. Therefore efficient algorithms and software tools are highly demanded.

In terms of computational aspects of biology, there exist three major sequence alignment problems. Literature denotes these problems as sequence alignment, pattern matching, motif finding.

The sequence alignment process is finding relationships between the sequences to identify similarity of species. The problem is mutations can occurs in DNA sequences and a single mutated nucleotide on middle of long sequence corrupt all alignment. This problem handled by dynamic programming and its variations like Smith – Waterman and Needleman – Wunsch (Smith and Waterman 1981) (Needleman and Wunsch, 1970).

Pattern matching is the second challenging problem in bioinformatics (Bishop, 2006). The process in the basic is detecting the exactly same of pattern presences of a given pattern in a long sequence. Since the single one sequence consist of about 3 billion nucleotides, in case of programming 3 billion characters, brute force approaches can't handle the problem in reasonable time. To solve this problem, searching approach should detect the positions which have no chance to match and skip these points for reduce volume of searching points.

Motif finding is the third and still under development problem in bioinformatics. The main idea for motif finding is detect the most repetitive sub sequences (D'haeseleer, 2006). There are many problem for the process like how many is motif length should be or how can group k length patterns. Current approaches usually offers divide and conquer technique. In the interface, an algorithm had presented for motif finding with the technique.

Sequence alignment is commonly is used by biologists to compare nucleotide sequences and to find functions of the genomes. There exists various software utilities that contain tools to do string matching methods such as sequence alignment, pattern matching and motif finding. On the other hand, advancements in web framework technologies, and programming languages enables to design better software tools. Also novel string matching algorithms give rise to new interfaces and tools.

Performance of the string matching algorithms depends on the data set and problem (Ozcan and Ünsal, 2015) Even further performance of and approximate string matching depends on the data. Most commonly used techniques are based on Dynamic Programming (Langmead and Salzberg, 2012). However, the techniques require high memory consumption.

Finally, some of the most efficient genomic analysis tools require licenses. Also the tools may have access limitations. In contrast, developing an open source tool with easy access property contributes to the educational demands. So that native language support can be introduced as well.

This study aims to present free, complete, user-friendly interface for whole bioinformatics field. The tool supports both English and Turkish languages. Therefore it can be useful for biology students who cannot read English.

The tool also introduces a novel approximate string matching algorithm. The algorithm speeds up string matching time. Due to its automata based technique, it also reduces memory consumption.

Overall the study has two contributions to the literature. First, it presents a novel approximate string matching algorithm. Second it introduces a new bioinformatics interface, which is coded with open source languages. The bioinformatics interface presents a simple and efficient interface. Together with its native language support it support academic improvement.

## 2. DEFINITIONS AND LITERATURE

Sequence alignment, pattern matching and motif finding problems have several solution approaches. In this section each major problem and their fundamental solutions will be mentioned in separated subheadings. Only de facto algorithms have explained in detail, other approaches in literature are variations of these major algorithms.

### 2.1. Sequence Alignment

Sequence alignment aims to find similarity of two sequences. Let's suppose that we have two sequences defined as

$$T_1 = t_0, t_1, \ldots, t_{n-1} \quad and \quad T_2 = t_0, t_1, \ldots, t_{n-1} \qquad \{t_i \in A, C, G, T\}$$

(1)

The sequences is not exactly the same but they are very similar to each other. In example, there are two text that are $T_1$ and $T_2$ and all characters of the texts are same except i-th character. The i-th character of $T_1$ is not same with i-th character of $T_2$. So the equation can be defined as:

$$T_{1(0)} = T_{2(0)}, T_{1(1)} = T_{2(1)}, \ldots, T_{1(i)} \neq T_{2(i)}, \ldots, T_{1(n-1)} = T_{2(n-1)}$$

(2)

To find optimum relativity between the sequences, sequences need to be realigned with gaps.

Sequence alignment is an essential problem because in real world, sequences not always remain in their original form of being created due to mutations. On the other hand, corruptions may arise during sequencing. To solve these kind of problems, there are two major approaches in literature; Smith-Waterman and Needleman-Wunsch. Both algorithms are variation of dynamic programming (Smith and Waterman, 1981).

Smith-Waterman Algorithm is a variation of dynamic programming. Dynamic programming approaches for sequence alignment have common variables like match score, mismatch score and gap score to calculate similarity score. Dynamic programming using for creating a relativity matrix from the sequences in Smith-Waterman algorithm. Each node of matrix value is maximum value of transitions from left, top and left top diagonal nodes. There are one more value which is zero for calculation maximum value additionally. Zero value gives a guarantee to there are no negative value in matrix. This particular precaution increases alignment success efficiency. The diagonal transition represent match, and other transitions means gaps. Once matrix have crated, trackback on matrix from last node to first node for calculating alignment score. Algorithm and explanations can be found on (Smith and Waterman, 1981).

Needleman – Wunsch is another derivation of dynamic programming which differs from Smith Waterman with negative values because of there is no zero condition on score function. Detailed explanation can be found on (Needleman and Wunsch, 1970)

### 2.2 Pattern Matching

In terms of exact string search, pattern matching can be defined as detecting occurrences of the pattern on a long sequences.

Let suppose we have a sequence $T$ as defined in sequence alignment. On the other hand we have another short sequence which entitled pattern, $P$, as:

$$P = p_0, p_1, \dots, p_{m-1} \qquad (3)$$

Pattern matching aims to locate the subsequences which is same with $P$ exactly or tolerance contrast in a range. In general, there are two approach to pattern matching: exact and approximate matching.

Exact pattern matching aims to find presences of exactly the same of $P$ in sequence as follows:

$$p_0 = t_i, p_1 = t_{i+1}, p_2 = t_{i+2}, \dots, p_{m-1} = t_{i+m+1} \qquad (4)$$

Current approaches using several skip algorithms to do this process efficiently. Skip algorithms boost matching process because many position skips and that means far less operations while matching. Essentially there are two main idea behind the skip algorithms; bad character and good suffix. Bad character means if there is any mismatch while matching, shift the pattern until the bad character is not in current sub sequence. The good suffix means if there is any prefix which same with suffix on mismatch point, shift pattern to align prefix with suffix. All major algorithms developed with this two approaches like KMP, Boyer – Moore, BOM etc.

Knuth – Morris – Pratt algorithm is an exact pattern matching algorithm which searches for presences of $P$ within a subsequence $T$ by using bad character approach. Before the matching process, preprocess should be done on $P$ for calculating skip count for every position of $P$. Detailed explanation can be found on (Knuth, Morris and Pratt, 1977).

The Boyer-Moore algorithm is another exact matching algorithm. As a distinct from KMP, Boyer-Moore algorithm combining good suffix and bad character approaches. While matching, if there is a mismatch between current part of $T$ and $P$, looking bad character and good suffix tables respectively for decide shift count on current position. Details can be found on (Boyer and Moore, 1977).

Another exact pattern matching algorithm is Backward Oracle Matching, BOM (Allauzen, Crochemoore and Raffinot, 1999). BOM is an automat based algorithm which is variation of BNDM algorithm. The details of BNDM algorithm can be found on (Navarro and Raffinot, 2002). The BOM algorithm based on the Boyer-Moore strategy. Thereupon try to match prefix with suffix of the pattern on mismatch position. On the other hand matching progress performs right as a necessity of good suffix approach. The algorithm using automat instead of tables unlike other Boyer-Moore approaches.

The first step of generating BOM automat is taking reverse of pattern and generate states for each character in reversed pattern and character transitions are added between the states respectively. After produced all factors of $P$, transitions for factors appends to the automat. Search algorithm and details can be found on (Allauzen, Crochemoore and Raffinot, 1999).

The second approach for pattern matching is approximate pattern matching. Approximate pattern matching differs from exact matching with mismatch tolerance. That means matching process tolerance to mismatches as long as number of mismatch is under threshold. Formula,

$$[P/t_{i\dots i+m+1}] < k \qquad (0 \leq k < m : m = pattern\ length) \qquad (5)$$

The approach make possible to find out mutated presences but this gain also cause computational weight to the matching process. For reducing this weight, approximate matching algorithms should have very efficient skip algorithms. On the other hand producing skip algorithm for approximate pattern matching algorithms harder than exact algorithms because skipped part could contain possible matches unlike exact approaches. To solve this problem usually skip algorithms does pre-processing on pattern, text or both of them.

Approximate pattern matching approaches compare pattern and text characters one by one until mismatch counter reach to the threshold or overall characters of the pattern has been compared. If mismatch counter exceeds the threshold the text shifts one character. On the other hand if does not exceed the threshold after all characters have been matched, that means there is a match on current position. In other words, naïve search using hamming distance to decide matching occurred on current position or not. If distance is under threshold there is a match or exceeding threshold is not. There is no skip mechanism in naïve search that means naïve search is a linear brute force matching algorithm but still useful small patterns and sequences due to no needs preprocess on neither text nor pattern.

An efficient approximate matching algorithm which is Burrows Wheeler transform firstly developed for data compression but nowadays there are many usage areas like pattern matching and sequence alignment (Durbin and others, 1998). The basic idea behind BWT is produce the permutations of the characters of text and positioning closely to similar contexts. That means in approximate matching, $k$ mismatched contexts can be found in $k$ neighborhoods. This process increases efficiency of approximate matching but on the other hand, pre-process on long patterns takes long execution times. Exhaustive explanations and detailed example could be found on (Burrows and Wheeler, 1994).
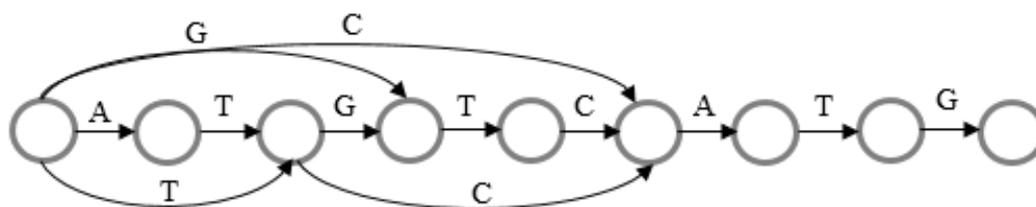
## 3. APPROXIMATE BOM

In this study, we present the approximate version of Backwards Oracle Matching algorithm. Recall that approximate pattern matching enables to find very similar presences of pattern on text. The flexible matching approach extends scope of matching but the profit comes with computation weight because of permutations of the pattern. In general approximate search algorithms are slower by nature. Especially matching takes huge execution times on long patterns. To overcome the problem, any preprocess should be done on pattern before matching.

BOM algorithm is an automat based exact pattern matching algorithm as mentioned above. The algorithm offers an automat for permutation problem. The automat provide how many shift performs any location on mismatch. The automat accelerate the matching process because shift counts for all permutations have already calculated. From this idea, the automat based approach could be apply on approximate pattern matching. The novel A-BOM algorithm is approximate variation of Backward Oracle Matching algorithm. BOM algorithm is best fit when long pattern searching case because all suffix combinations (factors) are calculated before search process and factor automaton prepared for search process. That means when any mismatch occurs on any position, search already know to how many shifts are necessary. Therefore like BOM algorithm, approximate BOM algorithm is supposed to be powerful on long pattern search.
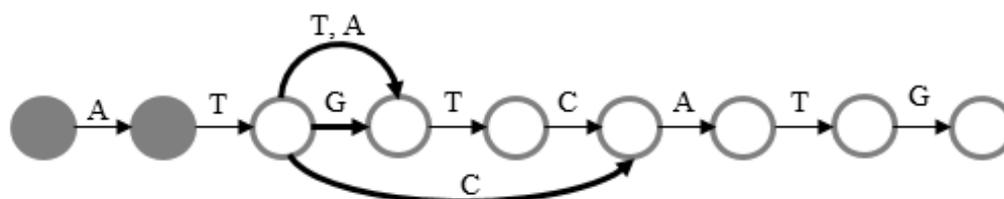
Approximate BOM algorithm using same automata logic and matching function with BOM algorithm can be found on (Allauzen and others, 1999). Approximation feature provided on calculating match score of current subsequence. Unlike BOM, the algorithm doesn't skip current position on mismatch until error counter is under threshold. When any mismatch occurs as long as error counter under threshold, matching branch out sub matching process by all transitions of current state.

Let's suppose there is a pattern like $P="GTACTGTA"$. The automat of reversed pattern shown in Figure 1.

***Figure 1:***
*Factor oracle automat of the pattern P=GTAACTGTA*

On the other hand let assume that also there is a sequence *T="GTACTTTA…"*. Let suppose that the threshold is 3. The score function performs matching from end to begin due to Boyer – Moore characteristics. When the score function come at third letter, the letter *T* is not match with the third character of pattern *G*. The approximation mechanism step in and branching starts at position 3. The root process branch out four sub matching process because of alphabet consist of four letter which *A, T, G* and *C*. The branching shown in Figure 2.



***Figure 2:***
*Branching on the mismatch at third character*

The sub processes perform matching after mismatch location and they can branch out as long as error doesn't reach up to threshold. Therefore matching score function designed as recursive. Branches go on matching with related transition of current state. There is a significant detail on the transitions. If there are no transition or the transition offers to jump over left error tolerance, branch go on matching with state of next expected character on the pattern. After all branches done of any parent process, largest matching score of branches adds parent's score and this adding process continues until the root matching process. After all branches of root process's done, score function returns the matching score to matching function. The matching function announces there is a match on current position when the matching score equals to pattern length. On the other hand if they are not equals, skips the matching location as much as subtraction of pattern length and matching score. Pseudo code of match score function explained in Algorithm 1.

## 4. EXPERIMENTAL RESULTS

In this section we introduce experimental performance comparison results of our approximate matching algorithm against Barrows Wheeler and Naive hamming distance based approximate matching algorithms. All the experiments we perform on a computer, with an Intel i5 2.30 GHz CPU with 4 GB of RAM and running Ubuntu 16.7 64-Bit. The code was written in C and compiled with CLion IDE.

1. # Score Function
2. WHILE index > 0 and current_state != length of automat
3. IF sequence[index] in current_state
4. move to next state
5. ELSE
   a. BREAK
6. END WHILE
7. IF index > 0 and  error_counter < threshold
8. FOREACH transition in alphabet
9. IF (next_state_transition – current_state + error_counter) <= threshold and
10. next_state_transition is not null
11. error counter += next_state – current state
12. append Score(next_state_transition) to temp_indexes
13. ELSE
14. append Score(next_expected_state) to temp_indexes
15. index =min(temp_index)
16. RETURN index

1

***Algorithm 1:***

*Match score function algorithm*

Test Sequence has 50K nucleotide and pattern lengths are variable. All experiments are repeated 100 times and averaged results are collected. In Table 1 execution times of algorithms are represented. Short names in Table 1 used as; A-BOM for Approximate Backward Oracle Matching, BWT for Burrows Wheeler Transform and NAIVE for Hamming Distance based approximate string matching.

The results on the table 1 presents execution time of algorithms in seconds. The observations donates A-BOM algorithm performance increasing with pattern length progressively unlike naïve matching algorithm. On the other hand, performance of BWT algorithm doesn't show significant variance on different patterns lengths.

The results denote that, A-BOM algorithm yields best performance results on long patterns. Results of Table 1 denotes that highest performance improvement occurs when the pattern length is 50 or 100. In general, observations donates the algorithm has from %30 to %80 better performance when pattern lengths over 10. Table 1 concludes that A-BOM is slower than naïve algorithm on short patterns, but still donates reasonable execution time. Increasing error rate influence unfavorably all algorithms. Our algorithm is affected than high error rates because of the branching characteristic.

In summary, A-BOM yields efficient approximate string matching for long patterns. Since pattern search on long DNA sequences is common, our algorithm can make sense for DNA sequences that contain mutations.

**Table 1. Average running time for 50K length DNA sequence**

| Pattern Length | Error Rate | A-BOM | BWT | NAIVE |
|:---:|:---:|:---:|:---:|:---:|
| 5 | 1 | 0. 084729 | 3.896761 | 0. 052436 |
| 10 | 1 | 0. 037810 | 3.867000 | 0. 046877 |
| 15 | 1 | 0. 031255 | 3.838748 | 0. 046883 |
| 25 | 1 | 0. 015626 | 3.932374 | 0. 053443 |
| 50 | 1 | 0. 002548 | 3.983852 | 0. 062507 |
| 5 | 2 | 0. 226232 | 3.821203 | 0. 062556 |
| 10 | 2 | 0. 084656 | 4.098129 | 0. 069030 |
| 15 | 2 | 0. 037353 | 3.824674 | 0. 068765 |
| 25 | 2 | 0. 022149 | 4.027949 | 0. 062505 |
| 50 | 2 | 0. 015628 | 3.898494 | 0. 069025 |
| 10 | 5 | 0. 268755 | 4.067687 | 0. 099815 |
| 25 | 5 | 0. 115907 | 3.974571 | 0. 099862 |
| 50 | 5 | 0. 052977 | 4.229776 | 0. 099817 |
| 75 | 5 | 0. 046839 | 4.139655 | 0.100351 |
| 100 | 5 | 0. 031255 | 4.007760 | 0. 100320 |

## 5. WEB INTERFACE

The web interface can accessible on "https://github.com/burakkoca/BioLab" address. In the interface; sequence alignment, pattern matching and motif finding can be easily done with user friendly graphic interface. Interface supports big amount of data. That means the interface can be used for academic and research projects. Students can use learning major solutions and try on own datasets also compare with several algorithms for the best fit solution. All algorithms which mentioned in proposal are presented in the interface. There are Smith-Waterman and Needleman- Wunsch algorithms for sequence alignment. For exact pattern matching KMP, Boyer-Moore and BOM algorithms are available and BWT, naïve search and A-BOM presented to perform approximate pattern matching. Motif finding can be done with greedy algorithm. How to use the interface introduced in separated subheadings for all solutions.

### 5.1. Sequence Alignment

The Sequence alignment algorithms can reachable "Sequence Alignment" collapsible item on left menu. Both KMP and Boyer-Moore algorithms have same interface for alignment. There are five field that has been labeled for sequences, gap, match score and mismatch score. After

all fields filled, press "Align" button for alignment. Note that, sequence length must be under ten thousand and must be consist of nucleotide letters.

When a query is given on the interface, results are returned in another page. Aligned string, gap score, gap ratio, match score, match ratio introduced in that page. In figure 3 first row presents aligned sequence. In Aligned sequence string, dashes stand for gaps and rods represent matches. Second row colored and presents match, mismatch and gap statics.

| 1. Dizilim = | C | T | - | G | C | T | - | G | A | T | A | C |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | : | | | | | | | | | | | |
| 2. Dizilim = | A | T | A | G | - | T | T | G | A | - | - | C |

| | |
|---|---|
| Eşleşme sayısı | : 6 / 12 |
| Eşleşme oranı | : %50.0 |
| Gap sayısı | : 5 / 12 |
| Gap oranı | : %41.6667 |

*Figure 3:*
*Sequence alignment results page*

## 5.2. Pattern Matching

The Pattern search algorithms can reachable "Exact Pattern Search" and "Approximate Pattern Search" collapsible items on left menu. Both exact and approximate matching pages have same interface for all algorithms but approximate pattern matching page has threshold field additively. The fields introduce with related parameter of matching. After all fields filled, press "Match" button for pattern matching. Note that, for only pattern matching sequence length could reach up to one million.

Pattern matching results are presented in consecutive web page. In other words sequence, match points, and match count presented in result page. In Figure 4, First row on the page presents treated sequence with colored presences of pattern. Second row represents matched locations and last row demonstrate match count.



*Figure 4:*
*Pattern matching results page*

## 5.3. Motif Finding

Motif finding algorithm can reachable "Motif Finding" collapsible item on left menu. There are two input for motif finding on the web page. The first input takes sequence and

99

second one for motif length. Multi sequences could be using for motif finding by separating sequences with comma. After all fields filled, press "Find Motif" button for motif finding.

Founded motifs and consensus motif presented in Results page that shown in Figure 3. Treated sequence stay on the first row of results. Second row presents founded motifs and last row show us the consensus motif.



*Figure 3:*
*Motif finding results page*

### 5.4. Comparison

Each solution group have own "Compare" tab in collapsible menus. Comparison pages have same interface with related solution page. Comparison results pages are same for all solutions. The results introduced in a table which each comparison parameter heading for each algorithm.

## 6. CONCLUSION

In this study we introduced a useful interface for all major bioinformatics problems solution algorithms. The interface differs from variations with wide scope. From educational to scientific purpose, any people who interested in bioinformatics can take the advantage of the interface because of the extensive contents opportunity. Also the interface offers to execute algorithms with large amount of data opportunity in free form. On the other hand, to the best of our knowledge there is no national interface that provide pattern matching, sequence alignment and motif finding for bioinformatics field and this study fulfill the need. We believe that the study nourish bioinformatics studies in our country and worldwide.

The second contribution of this study is a novel approximate string matching algorithm which present best performance for long patterns. Experimental results show that approximate approach of BOM speeds up approximate matching on long patterns. Our solution yields up to %80 better performance compared to Burrows Wheeler and Hamming Distance approach if pattern length is longer than 10. The results may contribute to the recent bioinformatics researches. For example A-BOM may fit better for approximate matching problems like error correction or merging read data from new generation DNA sequencing methods like Nanopores. In summary, the algorithm can be used for tolerant pattern matching with long patterns.

**REFERENCES**

**1.** Pevsner, J. (2015) *Bioinformatics and functional genomics*, John Wiley & Sons, UK

**2.** Smith, T.F. and Waterman, M.S. (1981) Identification of common molecular subsequences, *Journal of Molecular Biology*, Academic Press Incorporated, London, 40-48. doi: 10.1016/0022-2836(81)90087-5

3.  Needleman, S.B. and Wunsch, C.D. (1970) A general method applicable to the search for similarities in the amino acid sequence of two proteins, *Journal of Molecular Biology*, Academic Press Incorporated, London, 443-453.  doi: 10.1016/0022-2836(70)90057-4

4.  Bishop, C. M. (2006) *Machine learning and pattern recognition.* Information Science and Statistics. Springer, Heidelberg.

5.  D'haeseleer, P. (2006) How does DNA sequence motif discovery work?. *Nature biotechnology*, 24(8), 959-961

6.  Özcan, G., and Ünsal, O. S. (2015). Fast bitwise pattern-matching algorithm for DNA sequences on modern hardware. *Turkish Journal of Electrical Engineering & Computer Sciences*, 23(5), 1405-1417.

7.  Langmead, B., and Salzberg, S. L. (2012) Fast gapped-read alignment with Bowtie 2. *Nature methods*, 9(4), 357.

8.  Knuth, D.E., Morris, J.H and Pratt, W.R. (1977) Fast Pattern Matching in Strings, *Journal of Molecular Biology*, SIAM Journal on Computing, Philadelphia, 323-350. doi: 10.1137/0206024

9.  Boyer, R.S., Moore, J.S and Pratt, W.R. (1977) A Fast String Searching Algorithm, *Journal of Molecular Biology*, Communications of the ACM, New York, 762-772. doi: 10.1145/359842.359859

10. Alluzen, C., Crochemore, M. and Raffinot, M. (1999) Factor Oracle: A New Structure for Pattern Matching, *SOFSEM'99: Theory and Practice of Informatics*, Lecture Notes in Computer Science, Berlin, 291-306.  doi: 10.1007/3-540-47849-3_18

11. Ji, H. and Shendure, J. (2008) Next-generation DNA sequencing, *Nature biotechnology volume 26,* Nature Publishing Group, London, 1135-1145. doi: 10.1038/nbt1486

12. Burrows, W. and Wheeler, D. J. (1994) A block-sorting lossless data compression algorithm, *Technical Report 124, Digital Equipment Corporation,* Digital Equipment Corporation, California.

13. Durbin, R., Eddy, S. R., Krogh, A. and Mitchison, G. (1998) *Biological Sequence Analysis: Probabilistic Models of Proteins and Nucleic Acids*, Cambridge University Press, Cambridge.

14. Navarro, R. and Raffinot, M. (2002) *Flexible Pattern Matching in String*, The press Syndicate of The University of Cambridge, Cambridge.

15. Özcan, G. (2016) Detection of P53 Consensus Sequence: A Novel String Matching With Classes Algorithm, *Uludag University Journal of The Faculty of Engineering 21 (2),* Bursa, 269-282.