



In The Name of God

Simple and Efficient Pattern Matching Algorithms for Biological Sequences

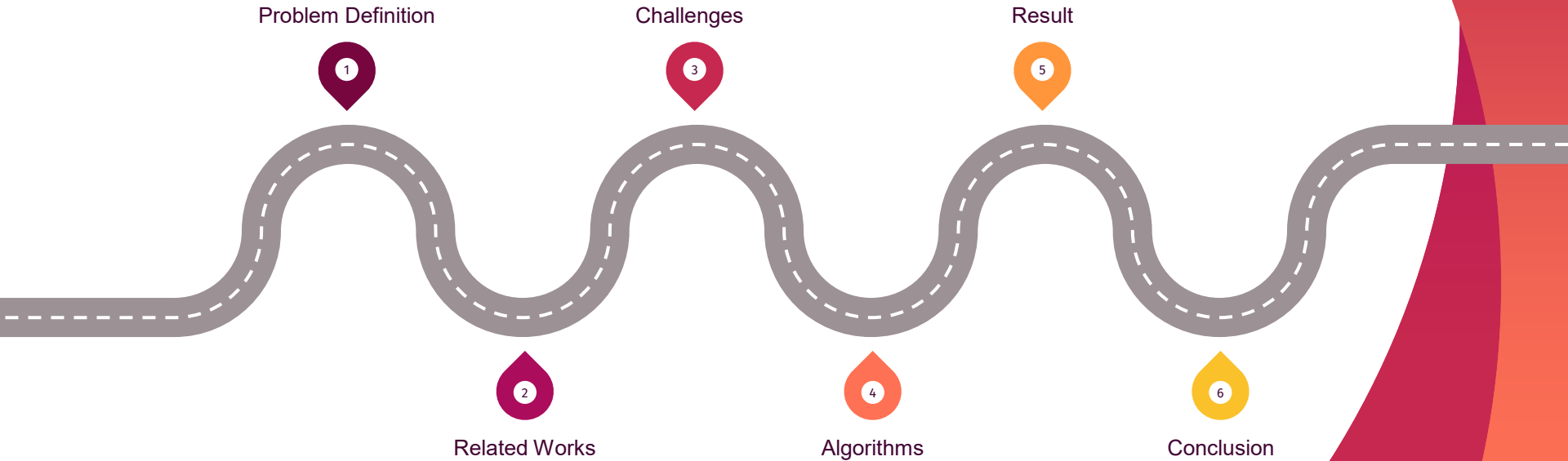
Presenters:

Nazanin Fereydoonizade – Sarvin Saravi

Supervisor:

Dr. Shakibian

Road Map



What is Pattern Matching?

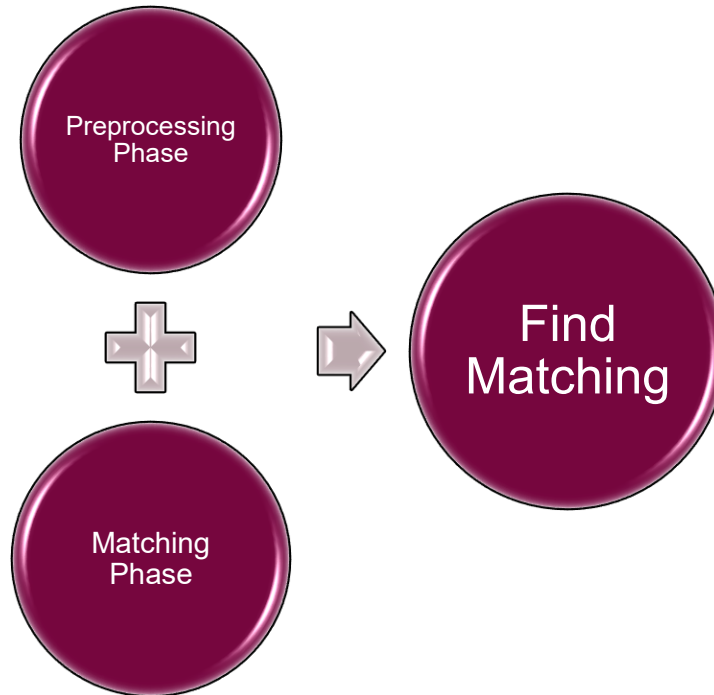
Scan a sequence or a database to detect the locations of a pattern

Applications

1. Image & Signal Processing
2. Information Retrieval
3. Search Engine
4. Text Processing
5. Parsers
6. NLP



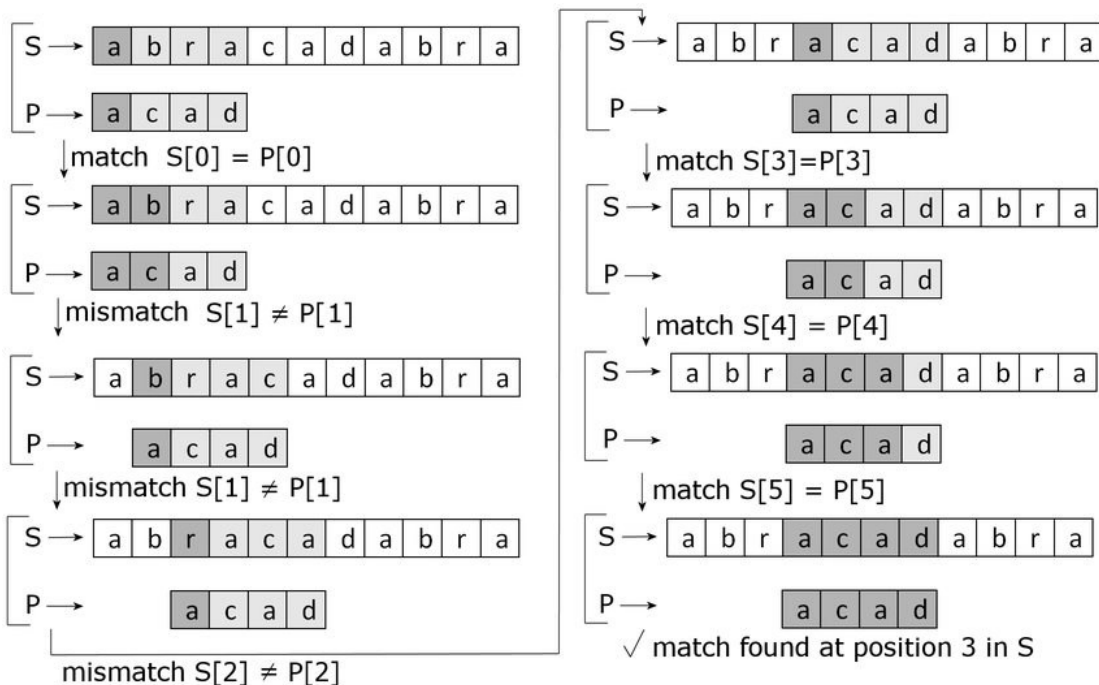
Pattern Matching Phases



Related Works

1. Brute Force ☐ Time problem, $O(mn)$
2. DFA-bases ☐ Not Scalable, $O(n)$
3. KMP ☐ Not Suitable for Small Strings, $O(n)$
4. Boyer Moore ☐ Depend on String size, $O(m+n)$ /
5. DCPM $O(mn)$

1. Brute Force



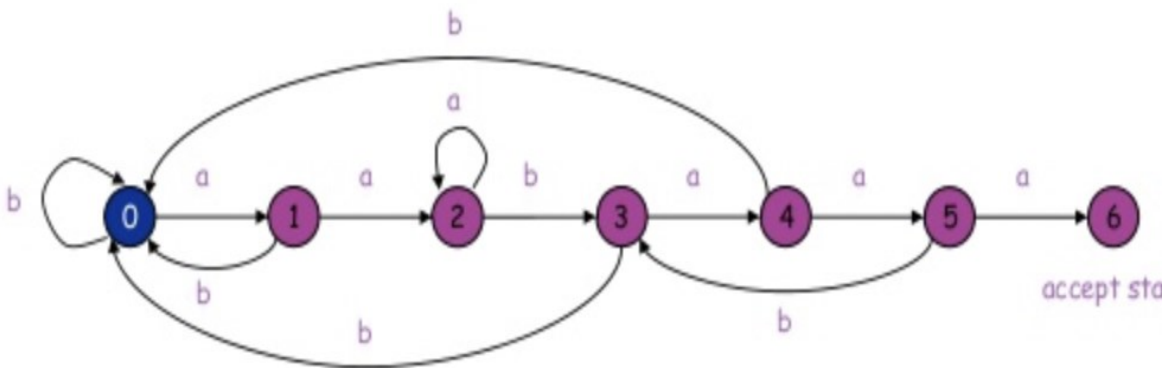
2. DFA-based

- Build DFA from Pattern
- Run DFA on Text

example

Text: aaabaabaaab

Pattern: aabaaa



Challenges

1. Library
2. Optimal File Size
3. Execution Time
4. Performance Analysis Using Different Encoding Techniques

Algorithm.1: First-Last Pattern Matching (Preprocessing Phase)

Algorithm 1 Preprocessing Phase of FLPM Algorithm

Input: Text t and pattern p stored in the arrays of $t[0..n-1]$ and $p[0..m-1]$, respectively, over finite alphabet Σ .

Output: The number of windows identified in this phase and their start indexes.

```
1.  $count \leftarrow 0$ 
2.  $num\_window \leftarrow 0$ 
3. WHILE  $count \leq n - m$  DO
4.     IF  $t[count] = p[0]$  THEN
5.         IF  $t[count + m - 1] = p[m - 1]$  THEN
6.              $window\_index[num\_window] \leftarrow count$ 
7.              $num\_window \leftarrow num\_window + 1$ 
8.         END-IF
9.     END-IF
10.     $count \leftarrow count + 1$ 
11. END-WHILE
```

Algorithm.1: First-Last Pattern Matching (Matching Phase)

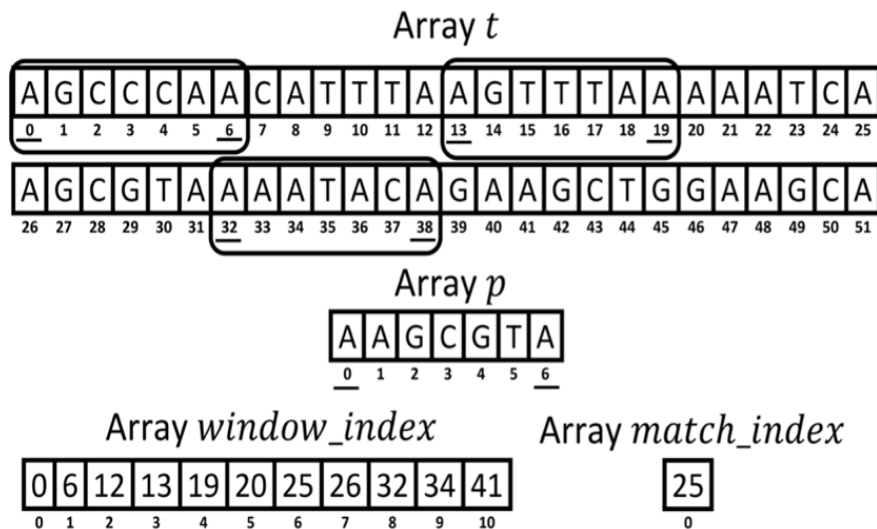
Algorithm 2 Matching Phase of FLPM Algorithm

Input: The number of windows identified in the preprocessing phase and their start indexes.

Output: The start index for all occurrences of pattern p in text t .

```
1.  $count \leftarrow 0$ 
2.  $num\_match \leftarrow$ 
3. WHILE  $count < num\_window$  DO
4.    $s \leftarrow window\_index[count]$ 
5.    $c \leftarrow 1$ 
6.   WHILE  $c \leq m - 2$  DO
7.     IF  $p[c] \neq t[s + c]$  THEN
8.       BREAK /*Exit the current loop*/
9.     END-IF
10.     $c \leftarrow c + 1$ 
11.  END-WHILE
12.  IF  $c = m - 1$  THEN
13.     $match\_index[num\_match] \leftarrow s$ 
14.     $num\_match \leftarrow num\_match + 1$ 
15.  END-IF
16.   $count \leftarrow count + 1$ 
17. END-WHILE
```

Algorithm.1: First-Last Pattern Matching (Example)



Algorithm.2: Processor-Aware Pattern Matching (Preprocessing Phase)

Algorithm 3 PAPM Algorithm

Input: Text t and pattern p stored in the arrays of $t[0..n-1]$ and $p[0..m-1]$, respectively, over finite alphabet Σ .

Output: The start index for all occurrences of pattern p in text t .

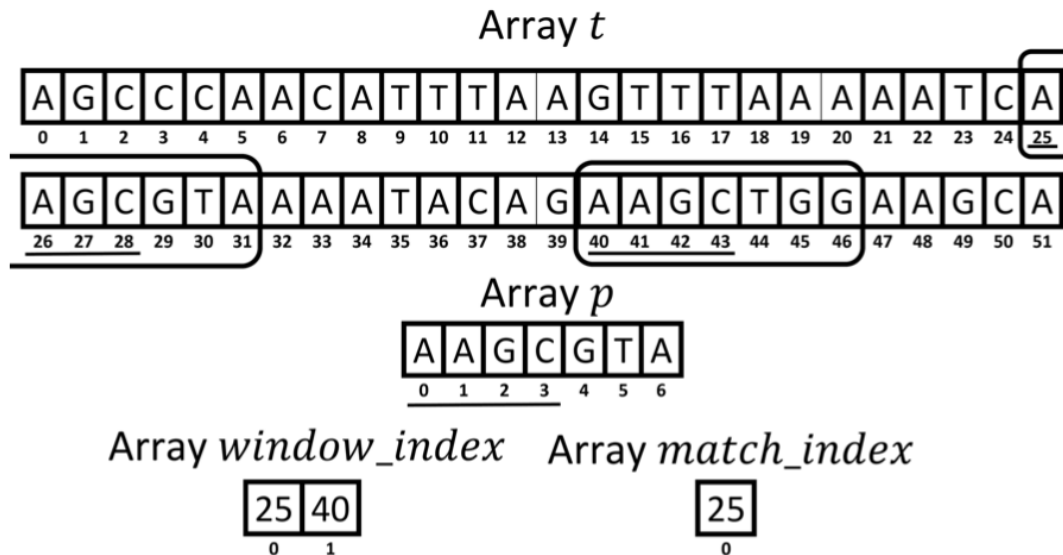
/ PREPROCESSING PHASE */*

1. $count \leftarrow 0$
2. $num_window \leftarrow 0$
3. **WHILE** $count \leq n - m$ **DO**
4. **IF** $t[count..count + word_len - 1]$
5. $= p[0..word_len - 1]$ **THEN**
6. $window_index[num_window] \leftarrow count$
7. $num_window \leftarrow num_window + 1$
8. **END-IF**
9. $count \leftarrow count + 1$
10. **END-WHILE**

Algorithm.2: Processor-Aware Pattern Matching (Matching Phase)

```
11.  $k \leftarrow m \bmod \text{word\_len}$ 
12. IF  $k = 0$  THEN
13.    $\text{start\_index} \leftarrow \text{word\_len}$ 
14. ELSE
15.    $\text{start\_index} \leftarrow k$ 
16. END-IF
17.  $\text{count} \leftarrow 0$ 
18.  $\text{num\_match} \leftarrow 0$ 
19. WHILE  $\text{count} < \text{num\_window}$  DO
20.    $s \leftarrow \text{window\_index}[\text{count}]$ 
21.    $c \leftarrow \text{start\_index}$ 
22.   WHILE  $c \leq m - \text{word\_len}$  DO
23.     IF  $p[c..c + \text{word\_len} - 1]$ 
24.        $\neq t[s + c..s + c + \text{word\_len} - 1]$  THEN
25.       BREAK
26.     END-IF
27.      $c \leftarrow c + \text{word\_len}$ 
28.   END-WHILE
29.   IF  $c = m$  THEN
30.      $\text{match\_index}[\text{num\_match}] \leftarrow s$ 
31.      $\text{num\_match} \leftarrow \text{num\_match} + 1$ 
32.   END-IF
33.    $\text{count} \leftarrow \text{count} + 1$ 
34. END-WHILE
```

Algorithm.2: Processor-Aware Pattern Matching (Example)




Algorithm.3: Least frequency pattern matching algorithm

- LFPM
- an enhancement of PAPM
- specialized for DNA
- many patterns vs. few patterns

Least frequency pattern matching algorithm

The first step before the preprocessing phase

- $\Sigma = \{ A, C, G, T \}$  HRG
- `word_len = b/8` for b-bit computer
- all possible words = $4^{\text{word_len}}$
- Create `freq_table`

Least frequency pattern matching algorithm

The first step before the preprocessing phase

TABLE 1. THE freq_table.

	0	1	2	3	4	5	...	$word_len - 2$	$word_len - 1$	$word_len$
0	A	A	A	A	A	A	...	A	A	n_0
1	A	A	A	A	A	A	...	A	C	n_1
2	A	A	A	A	A	A	...	A	G	n_2
3	A	A	A	A	A	A	...	A	T	n_3
.							...			
.							...			
.							...			
$4^{word_len} - 2$	T	T	T	T	T	T	...	T	G	$n_{4^{word_len}-2}$
$4^{word_len} - 1$	T	T	T	T	T	T	...	T	T	$n_{4^{word_len}-1}$

Algorithm 4 Filling the Table of Frequency**Input:**The array of reference and the length of this array.**Output:**The filling table of word frequency.

1. $freq_table[0..4^{word_len} - 1][0..word_len - 1] \leftarrow$
2. All possible words with $word_len$ length over $\Sigma = \{ACGT\}$
3. $freq_table[0..4^{word_len} - 1][word_len] \leftarrow 0$
4. $count \leftarrow 0$
5. **WHILE** $count \leq REF_len - word_len$ **DO**
6. $row_count \leftarrow 0, col_count \leftarrow 0$
7. $w[0..word_len - 1] \leftarrow REF$
 $[count..count + word_len - 1]$
8. **WHILE** $col_count < word_len$ **DO**
9. **SWITCH** $w[col_count]$
10. **CASE** 'A': nothing
11. **CASE** 'C': Increase row_count by
 $1 \times 4^{word_len - col_count - 1}$
12. **CASE** 'G': Increase row_count by
 $2 \times 4^{word_len - col_count - 1}$
13. **CASE** 'T': Increase row_count by
 $3 \times 4^{word_len - col_count - 1}$
14. **END-SWITCH**
15. $col_count \leftarrow col_count + 1$
16. **END-WHILE**
17. Increase $freq_table[row_count][word_len]$ by
 one
18. $count \leftarrow count + 1$
19. **END-WHILE**

Least frequency pattern matching algorithm

The preprocessing & Matching phase

- words in the pattern ➡ least frequent word ➡ start index
- similar to PAPM

Least frequency pattern matching algorithm

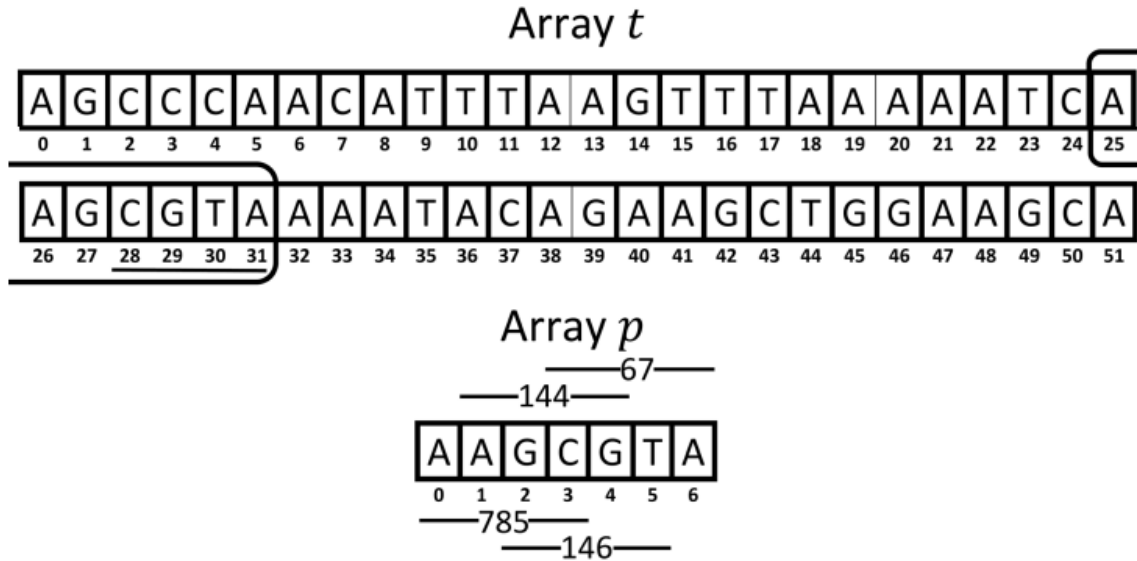


FIGURE 3. An example for the operation of LFPM algorithm.

Algorithm 5 LFPM Algorithm

Input: The filled *freq_table* on reference data. Text *t* and pattern *p* stored in the arrays of $t[0..n-1]$ and $p[0..m-1]$, respectively, over finite alphabet $\Sigma = \{A, C, G, T\}$.

Output: The first indexes for all occurrences of pattern *p* in text *t*.

/*PREPROCESSING PHASE*/

/*Step 1: finding the least frequent word in the pattern*/

```

1. count  $\leftarrow$  0
2. min_value  $\leftarrow$   $\infty$ 
3. min_index  $\leftarrow$  -1
4. WHILE count  $\leq$  m - word_len DO
5.   row_count  $\leftarrow$  0, col_count  $\leftarrow$  0
6.   w[0..word_len - 1]  $\leftarrow$ 
     p[count..count + word_len - 1]
7.   WHILE col_count < word_len DO
8.     SWITCH w[col_count]
9.     CASE 'A': nothing
10.    CASE 'C': Increase row_count by
         $1 \times 4^{\text{word\_len} - \text{col\_count} - 1}$ 
11.    CASE 'G': Increase row_count by
         $2 \times 4^{\text{word\_len} - \text{col\_count} - 1}$ 
12.    CASE 'T': Increase row_count by
         $3 \times 4^{\text{word\_len} - \text{col\_count} - 1}$ 

```

```

13.   END- SWITCH
14.   col_count  $\leftarrow$  col_count + 1
15. END-WHILE
16. IF freq_table[row_count][word_len] < min_value
    THEN
17.   min_value = freq_table[row_count][word_len]
18.   min_index  $\leftarrow$  count
19. END-IF
20. count  $\leftarrow$  count + 1
21. END-WHILE

```

Time complexity :

$$\theta(\text{word_len} \times (m - \text{word_len})) \approx \theta(m)$$

*/*preprocessing phase*/*

*/*Step 2: finding the windows*/*

```
22. count  $\leftarrow$  min_index
23. num_window  $\leftarrow$  0
24. WHILE count  $\leq n - (m - \textit{min\_index})$  DO
25.     IF t [count..count + word_len - 1]
26.         = p [min_index..min_index + word_len - 1]
27.         THEN
28.             window_index[num_window]  $\leftarrow$ 
                count - min_index
29.             num_window  $\leftarrow$  num_window + 1
30.         END-IF
31.     count  $\leftarrow$  count + 1
31. END-WHILE
```

/*MATCHING PHASE*/

```
32. count  $\leftarrow$  0
33. num_match  $\leftarrow$  0
34. k  $\leftarrow$  mmodword_len
35. WHILE count < num_window DO
36.     s  $\leftarrow$  window_index[count]
37.     c  $\leftarrow$  0
38.     w  $\leftarrow$  word_len
39.     WHILE c  $\leq$  m - 1 DO
40.         IF c > m - word_len THEN
41.             w  $\leftarrow$  k
42.         END-IF
43.         IF p[c..c + w - 1]
44.              $\neq$  t[s + c..s + c + w - 1] THEN
45.             BREAK
46.         END-IF
47.         c  $\leftarrow$  c + w
48.     END-WHILE
49.     IF c = m THEN
50.         match_index[num_match]  $\leftarrow$  s
51.         num_match  $\leftarrow$  num_match + 1
52.     END-IF
53.     count  $\leftarrow$  count + 1
54. END-WHILE
```


Results

- ❑ Brute Force (BF)
- ❑ Boyer-Moore (BM)
- ❑ Divide and Conquer
Pattern Matching
(DCPM)

- ❑ FLPM
- ❑ PAPM
- ❑ LFPM

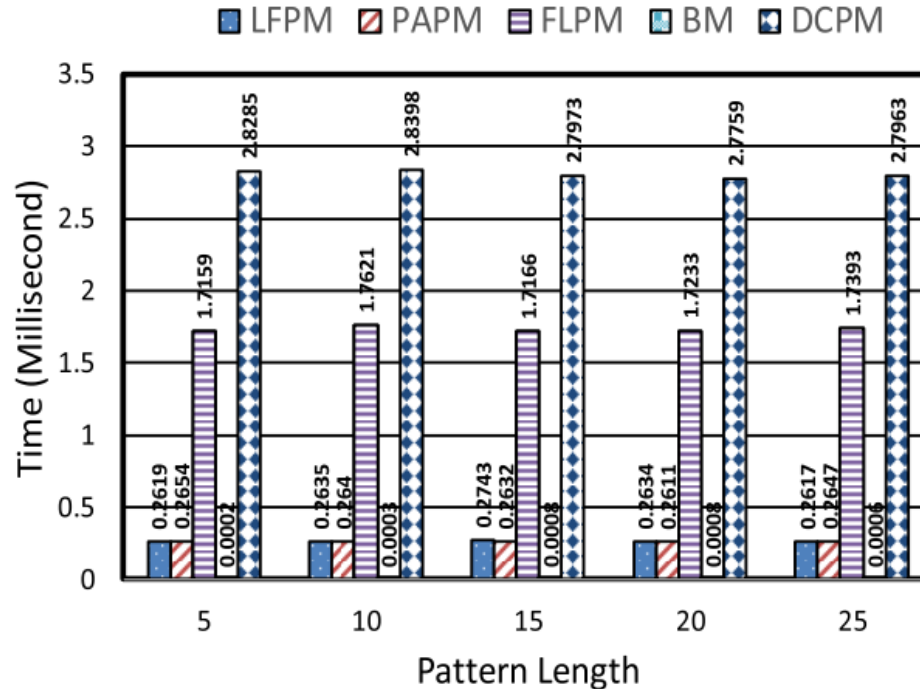
The specifications
of the computing
environment



Intell®core™2 Duo CPU T6600 (a 2.2 GHz clock)
A 2GB Memory
Acer (Aspire 5738)
Windows 7 ultimate 32 bit

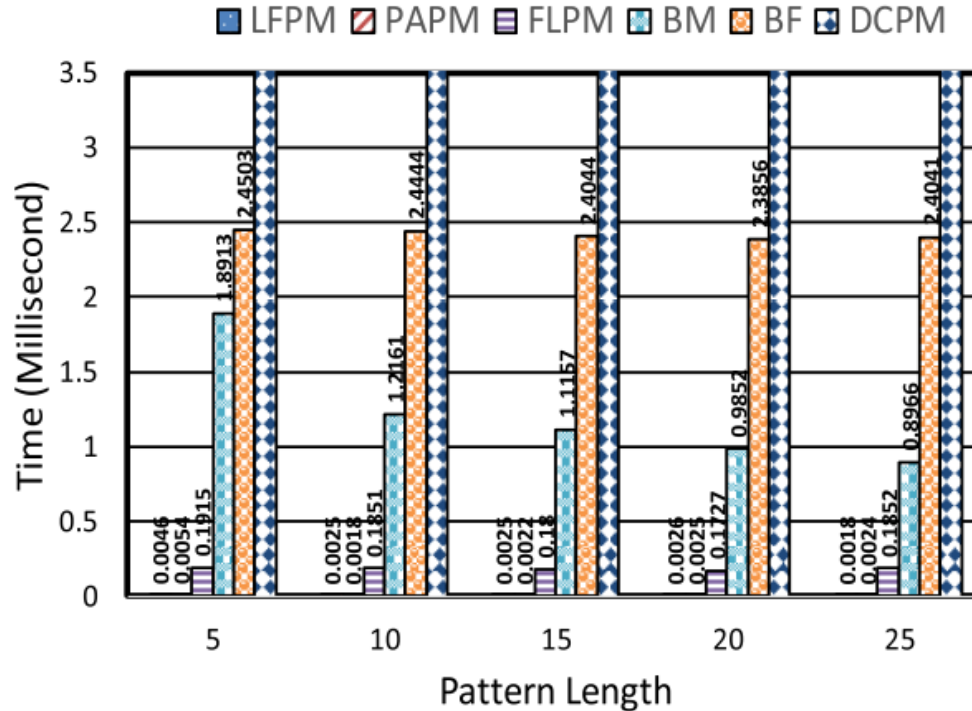
Results

A. THE TIME OF PREPROCESSING PHASE



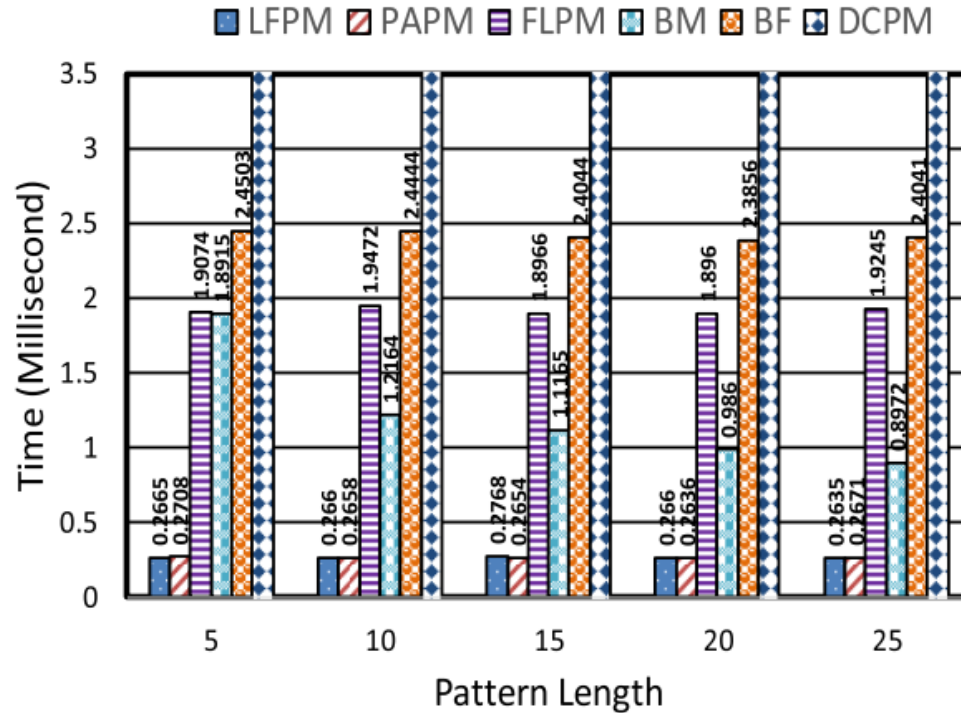
Results

B. THE TIME OF MATCHING PHASE

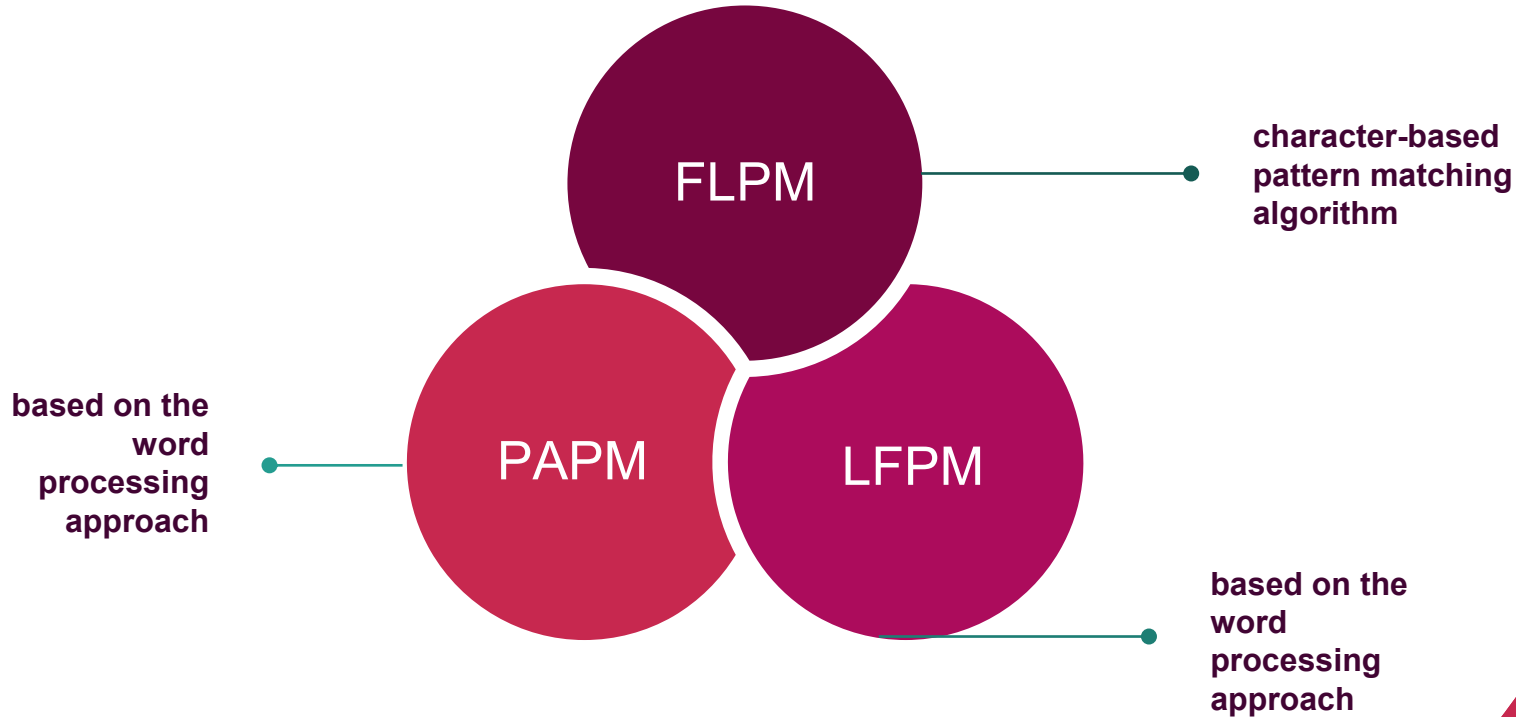


Results

C. TOTAL TIME OF PATTERNS



Conclusion



Future Research

- A parallel version of the presented algorithms
- The algorithms supporting approximate matching

Reference

- Neamatollahi, Peyman, Montassir Hadi, and Mahmoud Naghibzadeh. "Simple and efficient pattern matching algorithms for biological sequences." IEEE Access 8 (2020): 23838-23846.
- SAQIB HAKAK1, Amirrudin Kamsin1, PALAIAHNAKOTE SHIVAKUMARA1, GULSHIN "Exact String Matching Algorithms: Survey, Issues, and Future Research Directions"
- https://en.wikipedia.org/wiki/Boyer–Moore_string-search_algorithm
- https://en.wikipedia.org/wiki/Knuth–Morris–Pratt_algorithm
- https://en.wikipedia.org/wiki/Brute-force_attack

Thanks!

Any questions?

- Nazanin Fereydoonizade
- Sarvin Saravi