



RUTGERS

UNIVERSITY | NEW BRUNSWICK

ECE 573 - Data Structures and Algorithms | Spring 2018

(Project# 11)

DNA Sequence Detector

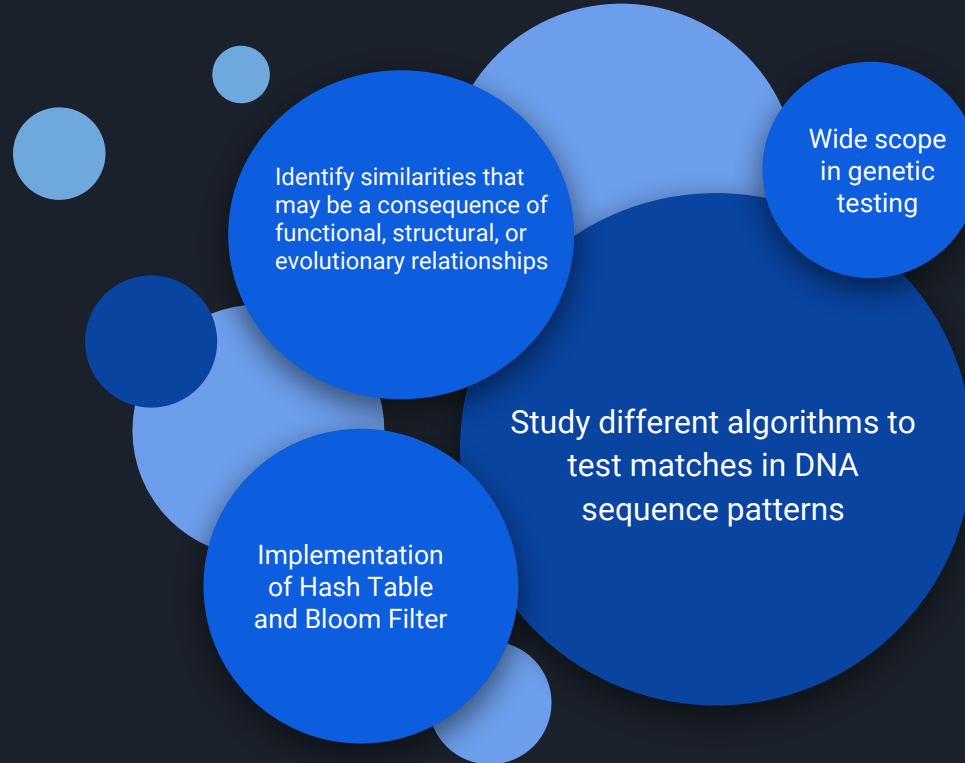
Divyaprakash Dhurandhar
(dd839)

Tina Drew
(tld95)

Anirudh Kulkarni
(ak1512)

Juhi Tripathi
(jt866)

Problem Statement - Motivation



Brute Force : Algorithm Design

- Given a text $txt[0..n-1]$ and a pattern $pat[0..m-1]$, the algorithm uses a function $search(char\ pat[], char\ txt[])$ that prints all occurrences of $pat[]$ in $txt[]$, assuming that $n > m$.

Text : A A B A A C A A D A A B A A B A

Pattern : A A B A

A A B A

A A B A

A A B A A C A A D A A B A A B A
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15
A A B A

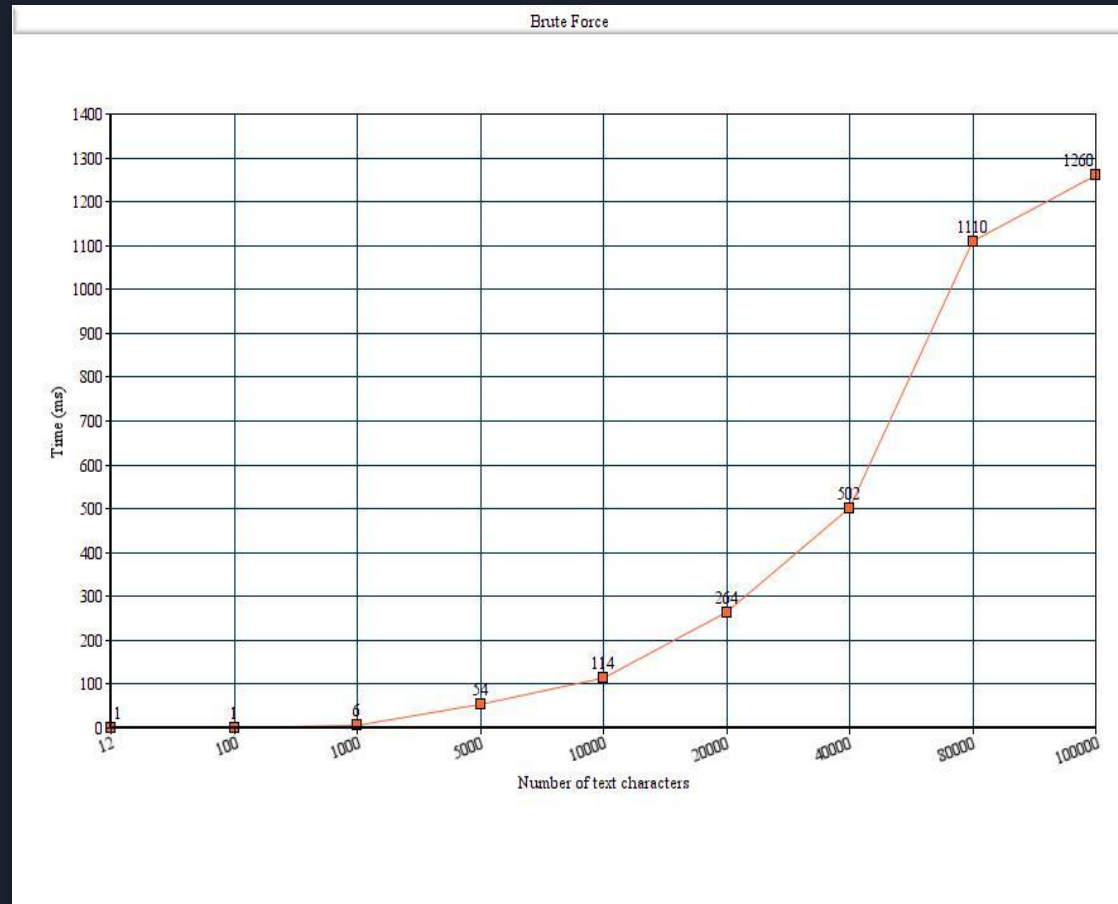
Pattern Found at 1, 9 and 12



Time Complexity

- **Best Case**
 - The best case occurs when the first character of the pattern is not present in the text at all.
 - The number of comparisons in best case is $O(n)$.
- **Worst case**
 - When all the characters in the pattern and the text are the same or the last character is different.
 - Time complexity in worst case is - $O(m*(n-m+1))$

- Created randomly generated DNA sequences.
- Searched for a pattern “CTG” in the sequence.
- The graph resembles $O(N^2)$



Goldman Algorithm

Create Separate Chaining Hash Table

Comparison

- Hash Table of optimum number of buckets 'V'
- Constant hash function : *"weight of string % number of buckets"*
- Push the string to the hashed index

- Compute hash function : *"weight of string % number of buckets"*
- Lookup in the hashtable at the hash mapping

Goldman Algorithm

GATTCC : Corpus String

CATTCG : Pattern String

$k = 3$ and $V = 2$

Possible k-length
sequences :-

GAT
ATT
TTC
TCC



Possible k-length
sequences :-

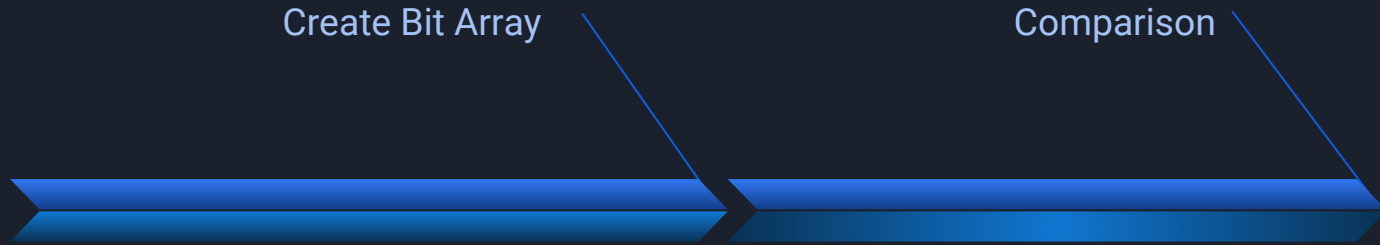
CAT
ATT
TTC
TCG

Lookup index in
hash table :

0
1
1
0

Matches Found : ATT, TTC

Bloom Filter Implementation



- Bit array of optimum size 'V', initialized with 0
- Constant hash function : *"weight of string % number of buckets"*
- Set the hash index to 1 for every string

- Compute hash function : *"weight of string % number of buckets"*
- Lookup for the hashed index bit, in the array

0 - not present

1 - present

Bloom Filter Implementation

GATTCC : Corpus String

| Possible k-length sequences :- | Corresponding index in array :- |
|--------------------------------|---------------------------------|
| GAT | 0 |
| ATT | 1 |
| TTC | 1 |
| TCC | 0 |

CATTCG : Pattern String

| Possible k-length sequences :- | Lookup index in array :- |
|--------------------------------|--------------------------|
| CAT | 0 |
| ATT | 1 |
| TTC | 1 |
| TCG | 0 |

$k = 3$ and $V = 2$



- GAT - sets index 0 to 1
- ATT - sets index 1 to 1
- TTC - index 1 is already 1
- TCC - index 0 is already 1

Matches Found : CAT, ATT, TTC, TCG
2 False positives! - Accounting for optimum value of 'V'



Complexity of Algorithms

BRUTE FORCE

$O(N^2)$

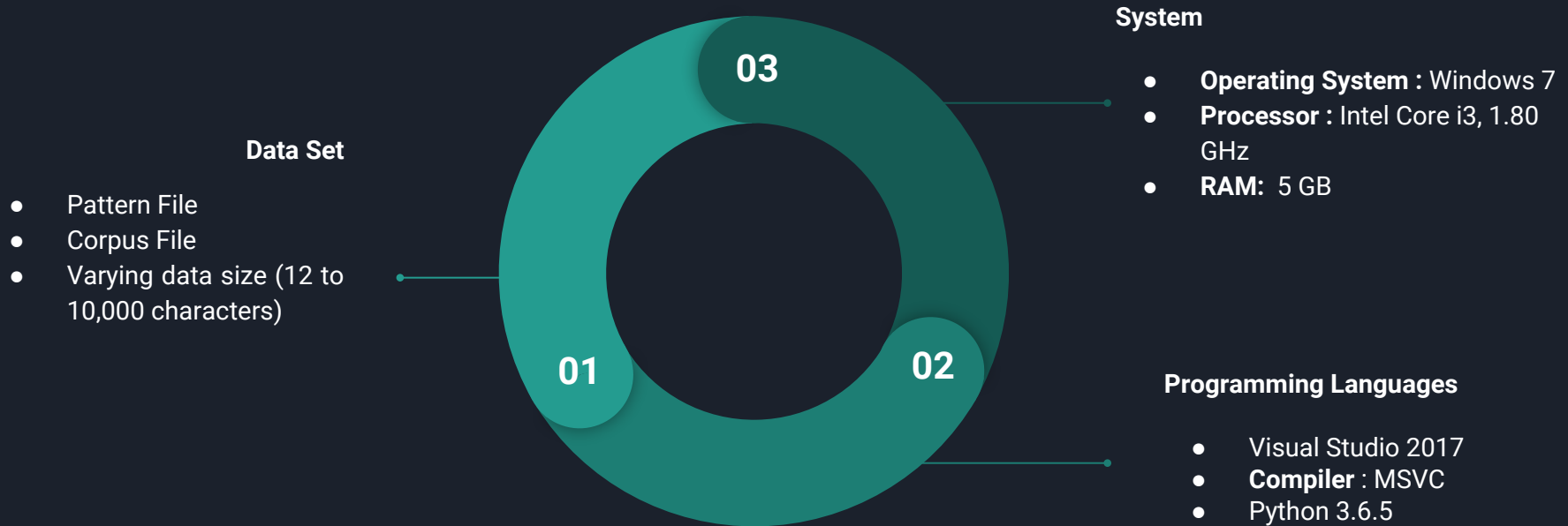
GOLDMAN METHOD

$O(N)$

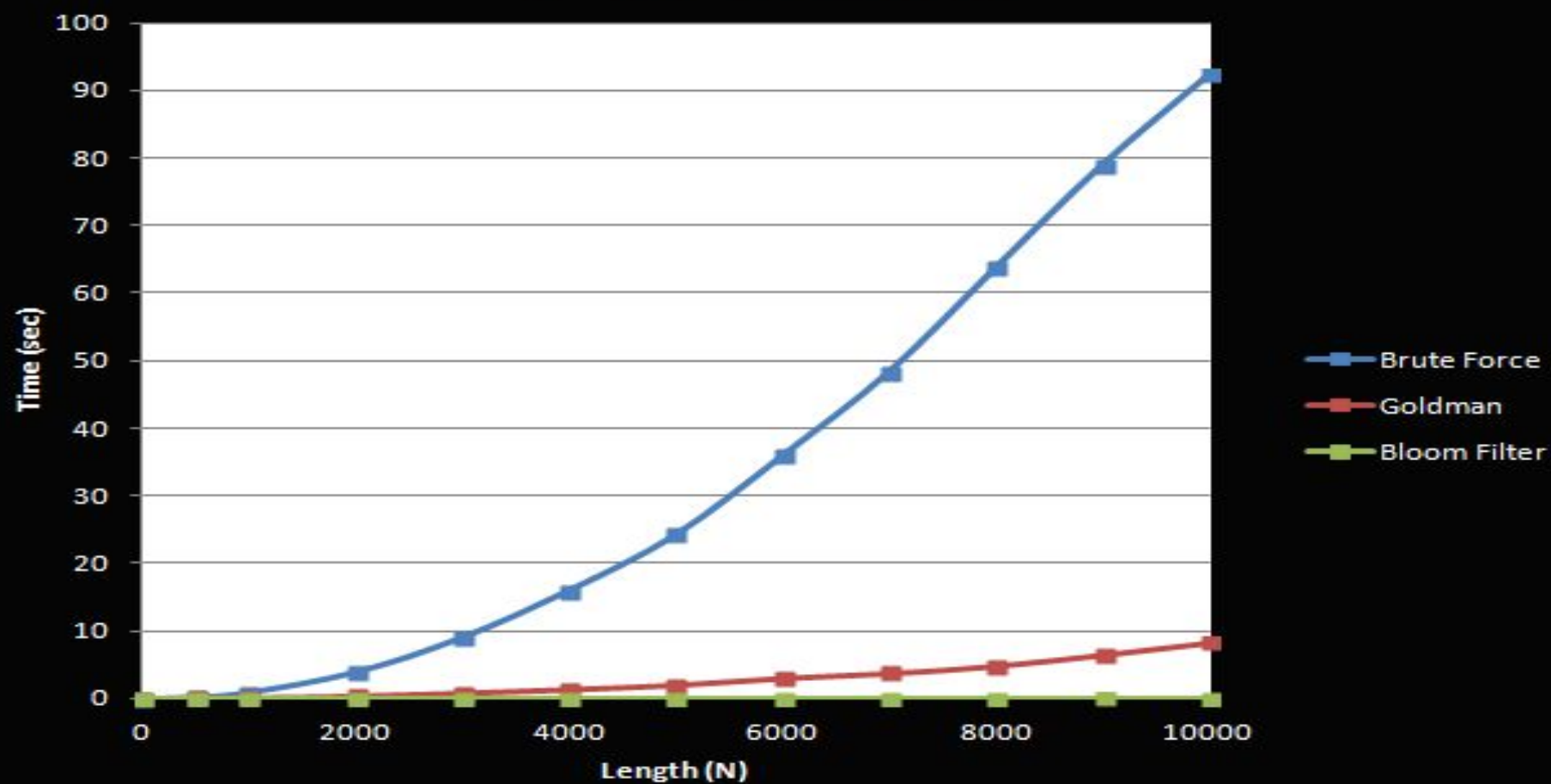
GOLDMAN - BLOOM METHOD

$O(1)$

Configuration and Details



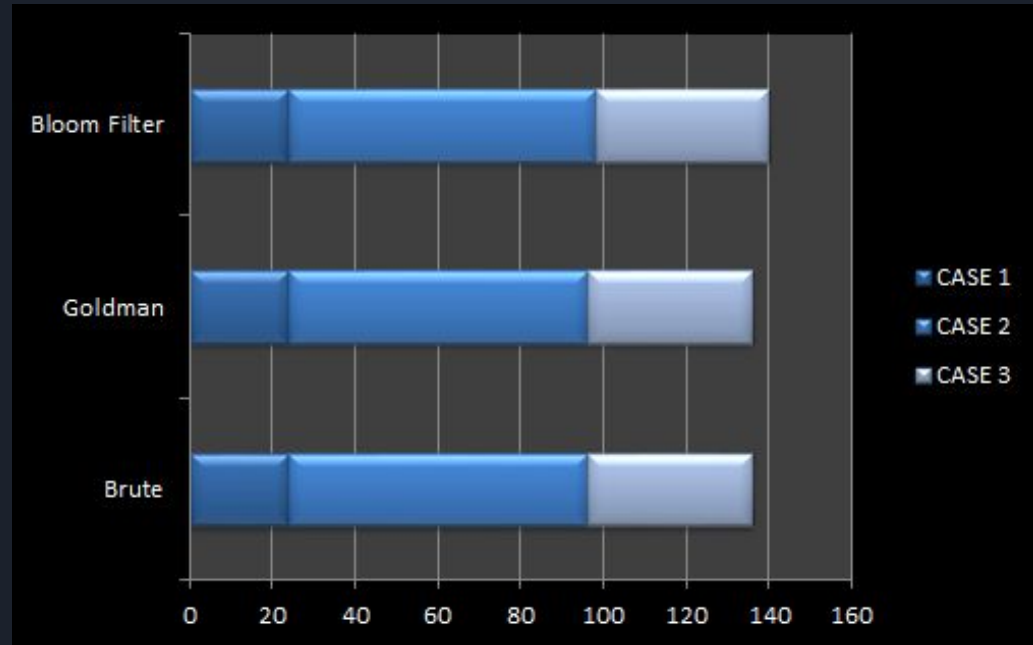
Results (Run Time)



Results (Accuracy)

Number of Matches

| | CASE 1 | CASE 2 | CASE 3 |
|--------------|--------|--------|--------|
| Brute | 24 | 72 | 40 |
| Goldman | 24 | 72 | 40 |
| Bloom Filter | 24 | 74 | 42 |



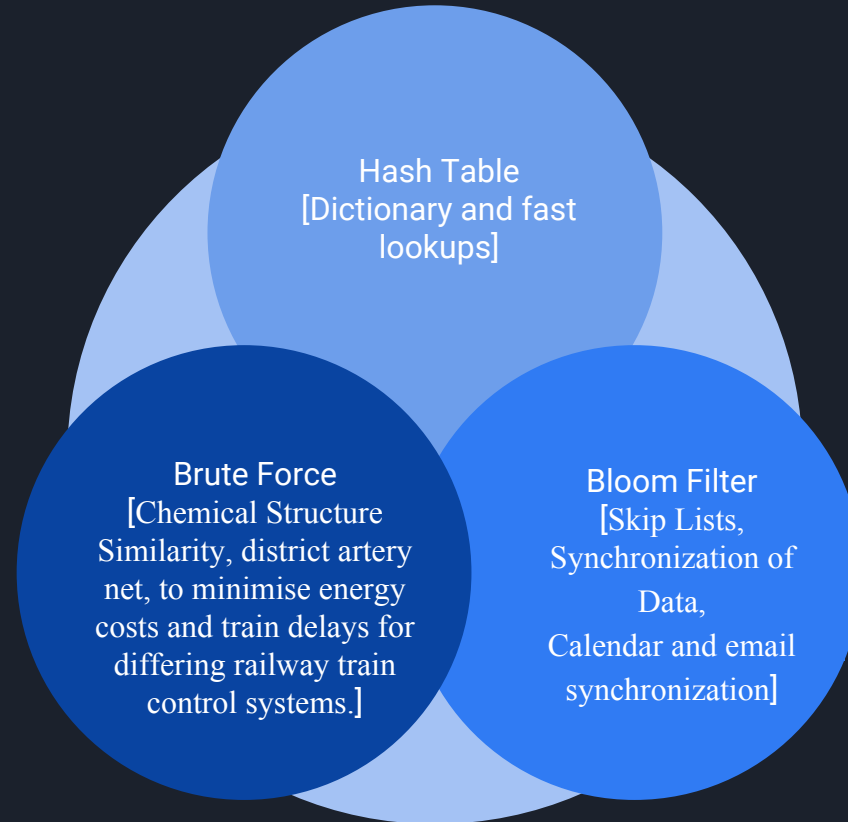


Conclusion

| | | |
|----|--------------------------------|--|
| 01 | BRUTE FORCE | <ul style="list-style-type: none">• Slowest Algorithm to implement Large Sequence |
| 02 | GOLDMAN ALGORITHM | <ul style="list-style-type: none">• Faster technique with separate chaining hash table |
| 03 | GOLDMAN - BLOOM IMPLEMENTATION | <ul style="list-style-type: none">• Reduce time and space complexity but lacks accuracy. |



Applications



... Questions?

THANK YOU FOR YOUR
ATTENTION!

