# Motif Finding

**CSE463 - Assignment**

*Prepared by:*
Lara Khan - 1905062
Majisha Jahan Disha - 1905089
Kazi Istiak Uddin Toriqe -1905104
Arnob shaha Ankon -1905108
Arko sikder -1905109

*Course Teacher:*
Dr. Muhammad Ali Nayeem
Assistant Professor, CSE, BUET

March 14, 2024

# Contents

# 1 Data

## 1.1 Biomarker

For Finding motif in dna sequence we have used some data sets. you will find all the data we have used.

# 2 Methods

To find motifs which is statistically significant we used two well known technique

## 2.1 Randomized Motif Search

### 2.1.1 Description

Randomized motif search methods employ statistical techniques to assess the significance of identified motifs. By comparing the observed motif occurrences in the input sequences to a randomized background model, these methods can determine whether the identified motifs are statistically enriched and unlikely to occur by chance alone. Despite the complexity of motif search algorithms, randomized motif search methods are designed to be computationally efficient. They leverage efficient data structures, optimization techniques, and parallelization to handle large-scale datasets and expedite motif discovery processes.

### 2.1.2 Pseudocode

---
**Algorithm 1** Randomized Motif Search

---
**Require:** $DNA$: Set of input DNA sequences, $k$: Motif length, $t$: Number of sequences
**Ensure:** Best motifs found
 1: $motifs \leftarrow$ randomlytMotifs$(DNA, k, t)$
 2: $bestMotifs \leftarrow$ motifs
 3: **while** True **do**
 4:     $profile \leftarrow$ Profile$(motifs)$
 5:     **for all** $sequence$ in $DNA$ **do**
 6:        $motif \leftarrow$ mostProbableKmer$(sequence, k, profile)$
 7:        add $motif$ to $motifs$
 8:     **end for**
 9:     **if** $score(motifs) < score(bestMotifs)$ **then**
10:        $bestMotifs \leftarrow motifs$
11:     **else**
12:        **return** $bestMotifs$
13:     **end if**
14: **end while**

---

## 2.2   Gibbs Sampler

### 2.2.1   Description

The Gibbs sampler motif search is a popular algorithm used in bioinformatics to identify overrepresented sequence motifs within a set of DNA or protein sequences. The Gibbs sampler is a stochastic algorithm that iteratively samples motif occurrences from the input sequences. It starts with an initial set of motifs and iteratively updates them to maximize the likelihood of observing the input sequences given the motifs.The Gibbs sampler uses a probabilistic model to represent motifs, typically in the form of a position weight matrix (PWM) or profile matrix. This model describes the probability of observing each nucleotide or amino acid at each position within the motif. Gibbs sampler motif search is a powerful algorithm for discovering overrepresented motifs in biological sequences. It leverages stochastic sampling techniques and probabilistic models to identify motifs that are statistically significant and biologically relevant

### 2.2.2   Pseudocode

---
**Algorithm 2** Gibbs Sampler Motif Search

---
**Require:** $DNA$: Set of input DNA sequences, $k$: Length of motif, $t$: Number of sequences, $N$ : Number of iterations
**Ensure:** Best motifs found
  1: $motifs \leftarrow$ randomlytMotifs($DNA$, $k$, $t$)
  2: $bestMotifs \leftarrow$ motifs
  3: **for** $j = 1$ **to** $N$ **do**
  4:     $i \leftarrow$ RandomRange(t)
  5:     motifs.popAtIndex(i)
  6:     $profile \leftarrow$ Profile($motifs$)
  7:     $motif \leftarrow$ mostProbableKmer($sequence[i]$, $k$, $profile$)
  8:     motifs.insertAtIndex(i,motif)
  9:     **if** $score(motifs) < score(bestMotifs)$ **then**
 10:        $bestMotifs \leftarrow motifs$
 11:     **end if**
 12: **end for**
 13: **return**   Best motifs found

---

# 3   Software

We used two Softwares : RSAT and Homer

## 3.1   Commands to run

### 3.1.1   RSAT

RSAT is an online tool.There's some filter format for visualizing data

- Format => Dataset input format

- Matrix length => The value of K

- Sites per Sequence

- Markov-Order

- Background Model

These are the important parameters of RSAT tool.

### 3.1.2   Homer

Homer is a software based on Perl language. There's an details installation guidelines available in here
Pre-requisites:

- C++/C compiler,Perl,GNU make utility

We are documenting the process how we did

- $sudo apt update

- $sudo apt install perl

- $sudo apt install build-essential

- $wget http://homer.ucsd.edu/homer/configureHomer.pl

- $perl configureHomer.pl -install

- $export PATH=/path/to/homer/bin:$PATH

Now for testing all things are running correctly give the following command

- findMotifs.pl

## 3.2  Scripts to run

### 3.2.1  Homer

We analyzed hm03.txt dataset using an background model to find more statistically significant motifs.The command :

- $findMotifs.pl x.fa fasta output/-fastaBg



We didnot use any background fasta file.For details see here

# 4  Results

## 4.1  Experiment configuration

for each dataset we have run three different programs using 3 different techniques **Varied K from 8 to 24**

### 4.1.1  Program$_1$

- We selected N = 1000 and numberSeed = 50 for progarm1

- We used Minimum hamming distance scoring technique for Randomized Motif Search and entropy technique for Giibs sampler which can be found in git repository from a commit The Entropy Function is as follow :

$$H(Nucleotide) = -\sum_{i=1}^{n} P(nuocleotide_i) \log_2 P(nucleotide_i)$$

where $Nucleotide$ is a discrete random variable, $P(nucleotide_i)$ is the probability of the $i$-th outcome, and $n$ is the total number of possible outcomes. And Nucleotide=['A','C','G','T']

### 4.1.2  Progarm$_2$

- We selected N = 1000 and numberSeed = 50 for progarm2

- We used averege Information gain technique for finding motifs on both methods

$$gain = \frac{maxEntropy * k - \sum_{i=1}^{k} H(Nucleotide)}{k}$$

where k denotes the motifs length

### 4.1.3  Program3

- We selected N = 100 and numberSeed = 50 for progarm3

- We tested with avg Minimum mismatch technique for both methods

## 4.2  Comparison

### 4.2.1  Progarm$_3$



Figure 1: for dataset1



Figure 2: for dataset2



Figure 3: for dataset3



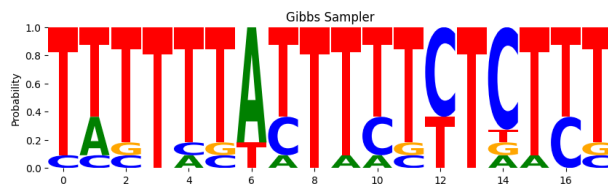Figure 4: Best Motif Logo of K=8 dataset3



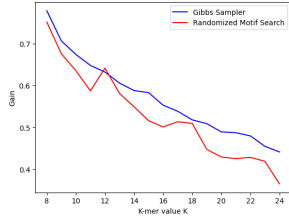Figure 5: Best Motif Logo of K=16 dataset3
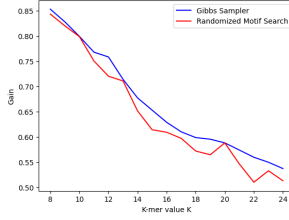
### 4.2.2 Progarm$_2$



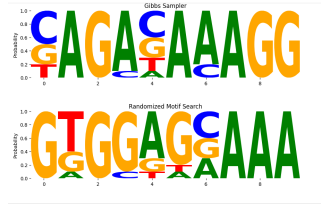Figure 6: for dataset1



Figure 7: for dataset2



Figure 8: for dataset3



Figure 9: Best Motif Logo of K=10 dataset1

### 4.2.3 Homer and RSAT Output for Dataset1
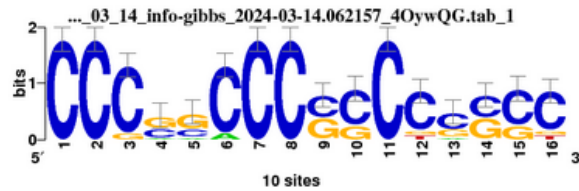


Figure 10: Homer Result for hm03 dataset



Figure 11: RSAT Result for hm03 dataset, K=16,numSeed=10,iteration=1000

# 5 Conclusion

We see from Figure1, Figure2, Figure3 **Randomized Motif Search** works significantly faster than **Gibbs Sampling**. When we vary the value of K and tends to increase it ,the running time for Gibbs sampler increases. But we can see from Figure6, Figure7 and Figure8 that Gibbs sampler gives much more statistically significant output as the Information-gain is larger and entropy is lower for Gibbs Sampler method.From those plots we also observe that for larger K the gain difference begins to increase than it was before.In the Homer tool testing the results of different K's are ranked.This is because they used some sophisticated scoring technique using mean,deviation,expected values [here] and they used a background fasta file for finding motifs which results statistically more accurate.We tested the RSAT tool output entropy and it was a little bit larger than our implementation result. Homer and RSAT both uses Gibbs Sampler by default.

# 6 References

- RSAT Article
- RSAT GibbsSampler
- Homer Motifs
- Homer Article
- Datasets