

## 1 Installation of our package

At first we need to compile the codes written in C/C++. Here naive attention score is used by default. Link `"/Ccode/main/CustomFunction.c"` to `"/Ccode/custom/CustomFunction.Attention.c"` to use alternative(frequency-based) attention score.

Recommended command: `ln -sf /Ccode/custom/CustomFunction.Attention.c /Ccode/main/CustomFunction.c`

Then run `"/compile.bash"` under `"/Ccode/main"`. Make sure your g++ supports c++11.

Links in folder `"/Ccode/bin"` to bash scripts in `"/Ccode/main"` also need to be established:

run `"/lnbash2bin.sh"` under `"/Ccode/main"`.

## 2 Demo

Here we illustrate how to use our programs by running some examples.

To be convenient, the epigenome here only have 2 chromosomes.

### 2.1 Align two epigenetic state sequences

Command: `/Ccode/bin/TwoRegions.out /Ccode/example/seq1 /Ccode/example/seq2 1.0`

1.0 is the  $\epsilon$  used by matching function and gap function described in section "Alignment algorithms".

### 2.2 Generate fake epigenome

Enter directory `/Ccode/example`.

Command: `/Ccode/bin/FakeGenomeGenerator.sh epigenomefiles epigenome/fakeepigenomefiles fakeepigenome/ /Ccode/bin/`

The script `FakeGenomeGenerator.sh` will try to remove folder `"fakeepigenome"` first, it may give you warning if this folder doesn't exist at all.

`"/Ccode/bin/"` is the path to the bash script.

### 2.3 Cut an epigenome in to small regions

Enter directory `/Ccode/example`.

Command:

`/Ccode/bin/CutFolder_Init.sh Para_Cut.example`

`/Ccode/bin/cut_sseq.out Para_Cut.example`

`"Para_Cut.example"` is a text file. There's annotation inside `"/Ccode/example"`

. The small regions are selected as the following:

$$chr_i^{[0,window)}, chr_i^{[step,step+window)}, chr_i^{[2*step,2*step+window)}, \dots$$

## 2.4 Query database search (with baseline database)

As mentioned in section "Ccode, DatabaseSearch.out", both query and database are stored as "sequence database" defined in section "Data representation, sequence database". Our code iteratively select the sequences in the first database, for each selected sequence, we implement our alignment algorithm on the second database and record the top  $k$  hits to the output file corresponding to this selected sequence.

We want to restrict the top hits such that they are apart from each other. Let the length of a query sequence to be  $l$ . Indexed in the database, a hit starts from  $s_h$ , ends at  $t_h$ , then from  $s_h - a * l$  to  $s_h + b * l$  there should be no other hit.

We also need you to provide a baseline database, each sequence in query database are also aligned to the baseline database. However, for whatever baseline database you provide, it doesn't alter the alignment result.

Command:

```
mkdir epigenomechr1align
/Ccode/bin/DatabaseSearch.out Paths_Search.example Para_Search.example
1.0
```

1.0 is the  $\epsilon$  used by matching function and gap function described in section "Alignment algorithms". You can find annotations of "Paths\_Search.example" and "Para\_Search.example" in "/Ccode/example/"

The summary is shown at "stdout". If there were  $n$  sequences in the query database, the summary would be an  $n * (k + 2)$  matrix. Form column 1 to column end are: length of this query sequence(compressed form), best alignment score to baseline set, top  $k$  alignment scores to the dataset.

For each query sequence, an output file is created to record its top hits. Explanation of these files is in "/Ccode/example/Record\_Alignments.annotation"

## 2.5 Query database search (without baseline database)

Basically this is identical to the previous example, but baseline database is no longer needed here. The top  $k$  alignments (for technical reasons, the program output more than  $k$  alignments, only the top  $k$  are correct) for each query sequence are same to the alignments found by the previous example.

Command:

```
mkdir epigenomechr1alignnative
/Ccode/bin/DatabaseSearchNative.out Paths_SearchNative.example Para_SearchNative.example
1.0
```

## 2.6 Horizontal alignment across an epigenome (with baseline)

In this example, we cut the epigenome into small regions first, for each small region, the top hits in this epigenome are searched afterwards.

This goal could be achieved by running the previous example chromosome after chromosome. Here we provide an integrated script.

Command:

```
mkdir horizontalalignment
/Ccode/bin/GenomeSearch_Path.sh Para_GenomeSearch.example /Ccode/bin/
```

"/Ccode/bin/" is the path to the bash script.

## 2.7 Horizontal alignment across an epigenome (without baseline)

This example is almost identical to the previous example, but baseline database is no longer needed here. The top  $k$  alignments for each segments are same to the alignments found by the previous example.

Command:

```
mkdir horizontalalignmentnative  
/Ccode/bin/GenomeSearchNative.Path.sh Para_GenomeSearchNative.example  
/Ccode/bin/  
"/Ccode/bin/" is the path to the bash script.
```

## 2.8 Annotate the outputs of horizontal alignment

In this example, we use "Algn2AnoBatch.sh" to annotate the outputs of horizontal alignments. For each output file of horizontal alignment, "Algn2AnoBatch.sh" does 3 things: sort the hits by alignment score, coordinating the hits to hg19, and mapping the chromosome index to chromosome name.

After the annotation, copy the given positions of query sequences and hits into UCSC website for detailed information:

<http://genome.ucsc.edu/cgi-bin/hgGateway?redirect=manual&source=genome.ucsc.edu>

Command:

```
/Ccode/bin/Algn2AnoBatch.sh Para_GenomeSearch.example Para_Annotation.example
```

or

```
/Ccode/bin/Algn2AnoNativeBatch.sh Para_GenomeSearchNative.example Para_Annotation.example
```

For each output file of horizontal alignment, the annotated file is under the same folder with the suffix ".anno".

# 3 Data representation

## 3.1 sequence file

A sequence file is a text file of  $l$  lines, each line have only one integer, representing the character(chromatin state) of this site. Usually the file name ends with ".seq".

## 3.2 compressed sequence file

A compressed sequence file is a compressed representation of a sequence file. It is a  $k$  line text file, in which each line has two integers: "(character) (repeated time)".

Usually the file name ends with ".sseq".

## 3.3 sequence database

A sequence database is composed of  $n$  compressed sequence files and one index file. The compressed sequence files should locate in the same folder.

The index file stores the paths to the compressed sequence files, each file name takes one line.

A sequence database could be a set of genome regions, an epigenome, a mixture of several epigenomes, or whatever you can come up with.

Usually an index file has no suffix.

## 4 Alignment algorithms

### 4.1 Smith-Waterman algorithm for chromatin state sequence alignment

#### 4.1.1 Definition of alignment problem

In this subsection, a chromatin state sequence is represented by a compressed format:

$$Seq := \{n, S, L\}$$

$n$  is the length of array  $S$  and  $L$ ,  $S$  is an array of chromatin states, and  $L$  an integer array. Inside  $L$ ,  $L_i$  is the repeat time of chromatin state  $S_i$ . For example:  $\{3, [a, b, c], [1, 2, 3]\} = [abbccc]$ .

The  $i$ 'th element of a sequence  $Seq$ ,  $[S_i, L_i]$  is a chromatin segment. Our goal here is to align the segments between two sequences.

Given two chromatin state sequences  $Seq_1 = \{n_1, S^{(1)}, L^{(1)}\}$ ,  $Seq_2 = \{n_2, S^{(2)}, L^{(2)}\}$ , a match  $f$  is a string. On each position of  $f$ ,  $f_i$  has three choices: 'm', 'g<sub>1</sub>', or 'g<sub>2</sub>'. 'm' is a match, 'g<sub>1</sub>' is a deletion in  $Seq_1$  and 'g<sub>2</sub>' is a deletion in  $Seq_2$ . The number of 'm' plus the number of 'g<sub>1</sub>' should equal to  $n_1$ . Symmetrically the number of 'm' plus the number of 'g<sub>2</sub>' should equal to  $n_2$ .

Let  $u^1$  and  $u^2$  be the index of  $f$  on  $Seq_1$  and  $Seq_2$ .  $u_i^1$  could either be an integer or  $\emptyset$ , indicate which segment of  $Seq_1$   $f_i$  is pointing to, or  $f_i$  is pointing to a gap.

Now we are able to define the matching score  $H$  given  $Seq_1$ ,  $Seq_2$  and  $f$ .

$$H(Seq_1, Seq_2, f) = \sum h(f_i, u_i^1, u_i^2, Seq_1, Seq_2)$$

If  $f_i$  equals to 'm',  $h(f_i, u_i^1, u_i^2, Seq_1, Seq_2) = MF(u_i^1, u_i^2, Seq_1, Seq_2)$ ,  $MF$  is the matching function, otherwise,  $f_i$  equals to 'g<sub>s</sub>', then  $h(f_i, u_i^1, u_i^2, Seq_1, Seq_2) = GF(u_i^s, Seq_s)$ ,  $GF$  is the gap function.

In alignment problem, we want to find the best match  $f^*$ , which maximize the alignment score.

Alignment is achieved by dynamic programming. This algorithm iteratively maintain and update a matrix  $M$ :

$$\begin{aligned}
M(i, 0) &= 0, \text{ for } 0 \leq i \leq n_1; \\
M(0, j) &= 0, \text{ for } 0 \leq j \leq n_2; \\
M(i, j) &= \max \begin{cases} M(i-1, j-1) + MF(i, j, Seq_1, Seq_2) & \text{(Mis)match} \\ M(i-1, j) + GF(i, Seq_1) & \text{Deletion in } Seq_1 \\ M(i, j-1) + GF(j, Seq_2) & \text{Deletion in } Seq_2 \end{cases} \\
&\text{for } 1 \leq i \leq n_1, 1 \leq j \leq n_2,
\end{aligned} \tag{1}$$

$M_{ij}$  is the maximal alignment score of the two subsequences  $Seq_1^{[1,i]}$  and  $Seq_2^{[1,j]}$ .

#### 4.1.2 Local alignment

In Practice, we prefer local alignment rather than global alignment described above. To achieve the goal of local alignment, we need to add a small modification to the dynamic programming.

$$\begin{aligned}
M(i, 0) &= 0, \text{ for } 0 \leq i \leq n_1; \\
M(0, j) &= 0, \text{ for } 0 \leq j \leq n_2; \\
M(i, j) &= \max \begin{cases} 0 & \text{Restart the alignment from here} \\ M(i-1, j-1) + MF(i, j, Seq_1, Seq_2) & \text{(Mis)match} \\ M(i-1, j) + GF(i, Seq_1) & \text{Deletion in } Seq_1 \\ M(i, j-1) + GF(j, Seq_2) & \text{Deletion in } Seq_2 \end{cases} \\
&\text{for } 1 \leq i \leq n_1, 1 \leq j \leq n_2,
\end{aligned} \tag{2}$$

## 4.2 Attention score

Before defining the Matching function and Gap function, we want to define the attention score first. Based on an observation: given a sequence  $Seq$ , some segments in  $Seq$  could be more important than other segments, we want to score the importance of each segment in  $Seq$ . The attention score is an array of float, denoted by  $AS(Seq)$ . Furthermore, the  $i$ 'th element of  $AS(Seq)$  is denoted by  $AS(Seq, i)$ .

In our package, we provided three kind of attention scores:

**Naive attention score:**

$$AS^n(Seq, i) = 1$$

**Minus log frequency(of chromatin states):**

$$AS^{lf}(Seq, i) = AS^{lf}(\{n, S, L\}, i) = -\log(\text{frequency}(S_i))$$

Frequency of the chromatin states are computed outside and passed to the alignment functions.

**Unpredictability:**

$$AS^{up,w}(Seq, i) = GramScore^w(Seq, i) * LocalRankScore^w(Seq, i) * LengthScore(Seq, i)$$

In this definition, we aim to quantify how much the neighbourhood area  $Seq^{[i-w, i+w]}$  predicts  $S_i$ . The unpredictable segments are thought to be more important than predictable segments.

In the subsequence  $Seq^{[i-w, i+w]}$ , if neither of pattern  $S_{i-1}S_i$  nor  $S_iS_{i+1}$  occurs elsewhere,  $GramScore^w(Seq, i) = 1$ , if both patterns occurs elsewhere,  $GramScore^w(Seq, i) = 4$ , otherwise  $GramScore^w(Seq, i) = 2$ .

For  $Seq^{[i-w, i+w]}$ , we sort the chromatin states occurred in this region by their frequency,  $LocalRankScore^w(Seq, i) = 1 + \frac{rank(S_i)-1}{|\{States \in Seq^{[i-w, i+w]}\}| - 1}$

If  $L_i = 1$ ,  $LengthScore(Seq, i) = 0.5$ , otherwise  $LengthScore(Seq, i) = 1$ .

### 4.3 Matching function

$$MF(Seq1, Seq2, i, j) = \begin{cases} AS(Seq1, i) + AS(Seq2, j) & S1_i = S2_j \\ -\epsilon * (AS(Seq1, i) + AS(Seq2, j)) & S1_i \neq S2_j \end{cases}$$

### 4.4 Gap function

$$GF(Seq, i) = -\epsilon * AS(Seq, i)$$