

# Subject: Large Scale Data Structures

Assignment: 04

By Saurabh Kakade

Sk2354@nau.edu

## 1. Part 01: A (OUTPUT)

```
[sk2354@ondemand /scratch/sk2354/A04/Part01]$ srun --mem=05GB -t 00:60:00 ./homework A ./hw3_dataset.fa ./test_genome.fasta

The number of arguments passed: 4
The first argument is: /scratch/sk2354/A04/Part01/./homework
The second argument is: A
The third argument is: ./hw3_dataset.fa
The fourth argument is: ./test_genome.fasta

Function in C++ that implements the Smith-Waterman alignment between two genomic sequences.:

Number Of Read Query Lines (reads): 30

Total 70-character fragments: 427

*****

Sequence 01 and Sequence 02:

AGGGTTCAGGAAAGAGTTTGAAGGATGCTGAACATCTTGAATAGGAGAC
ATTAAAGGTTTATACCTTCCAGGTAACAAACCAACCACTTTCGATCTCTGTAGATCTGTTCTCTAAA

Max score in the matrix is: 18

Alignment of sequence 01 and sequence 02: (traceback_array)

TGAAGGATGCTGAACATC
| x||x|x|xxx x||x|x
T_TAAAGGTTTA_TACCTT

*****

*****

Sequence 01 and Sequence 02:

AGGGTTCAGGAAAGAGTTTGAAGGATGCTGAACATCTTGAATAGGAGAC
ACAAATAGACATATAAAGATAGACAAACATACAAACAGACAAAGATAAATAAAAAATAAAAATAAAAAC

Max score in the matrix is: 20

Alignment of sequence 01 and sequence 02: (traceback_array)

AAAGAG_TTTGAAGGAT_GCTGAACATCTTGAATAG
x||xxx xxxxxx|xxxx xxxx||xxx |xx||xxx
CAAATAGACATATAAAGATAGACAAACA_TACAAACA
```

\*\*\*\*\*

Sequence 01 and Sequence 02:

AGGGTTCAGGAAAGAGTTTGAAGGATGCTGAACATCTTGAATAGGAGAC  
CGCACGCGCCACCCGCTCCCACACCCACTCCCTCTCACACACGCACTCGCGCCACCCTCTCGCTCGCG

Max score in the matrix is: 11

Alignment of sequence 01 and sequence 02: (traceback\_array)

GCTGAACATCTTG  
xxxxx| |xxxx|x  
CGCCACCCCTCTC

\*\*\*\*\*

\*\*\*\*\*

Sequence 01 and Sequence 02:

AGGGTTCAGGAAAGAGTTTGAAGGATGCTGAACATCTTGAATAGGAGAC  
GTGTGGGAGAGCGAGGGCGCGCGTGAGTGGGTGGGTGTGCGAGTGCGAGAGAGCGGGTGTGCGGGGGAGT

Max score in the matrix is: 23

Alignment of sequence 01 and sequence 02: (traceback\_array)

TTGAAGGATGCTGAACATCTTGAATAGGAG  
|xx|xxx| |xxxx|xx|xxxxxxxx|xx| |xx  
TGTGGGAGAGCGAGGGCGCGCGTGAGTGGGT

\*\*\*\*\*

\*\*\*\*\*

Sequence 01 and Sequence 02:

AGGGTTCAGGAAAGAGTTTGAAGGATGCTGAACATCTTGAATAGGAGAC  
TGTCTTTGTGTTTATGTCTATGTATATCTTTCTGTATATGTGTCTATTTTCTATGTTTATCTGTGTTTC

Max score in the matrix is: 19

Alignment of sequence 01 and sequence 02: (traceback\_array)

GAGTTTGAAGGATGCTGAACATCTTGAATAGGAG  
xxxxx| |xxxxxxxxxxxx|xx|xx|x|xxxxxxxxxxxx  
TGTGTTTATGTCTATGTATATCTTTCTGTATATGT

\*\*\*\*\*

Sequence 01 and Sequence 02:

```
AGGGTTCAGGAAAGAGTTTGAAGGATGCTGAACATCTTGAATAGGAGAC
TATTTATGTTTTATCTGTGTCTGTCTCTGTATTTCTTTATATATGTTTCTATTTTTTGTATCTTTTA
```

Max score in the matrix is: 19

Alignment of sequence 01 and sequence 02: (traceback\_array)

```
TTTGAAGGATGCTGAACATCTTGA_ATAGG
x|xx|xxx|xx xxxxxx|x|x xxxxx
CTGTGTCTGTC_TCTGTATTTCTTTATATAT
```

\*\*\*\*\*

\*\*\*\*\*

Sequence 01 and Sequence 02:

```
AGGGTTCAGGAAAGAGTTTGAAGGATGCTGAACATCTTGAATAGGAGAC
ATACAAAAGAAAAAAAAACATAGAGAAAAACAAACATAAAAAACAAATAAAGACAAAGATAGAGATAG
```

Max score in the matrix is: 20

Alignment of sequence 01 and sequence 02: (traceback\_array)

```
GGATGCTGAACATCTTGAATAGGAGA
xxxxxxxxxx|x|xxxxx|xxx xxxx
TACAAAAGAAAAAAAAACAT_AGAG
```

\*\*\*\*\*

\*\*\*\*\*

Sequence 01 and Sequence 02:

```
AGGGTTCAGGAAAGAGTTTGAAGGATGCTGAACATCTTGAATAGGAGAC
GGGCGAGTGAGCGAGCGTGC GGCGTGAGTGGGTGCGGGAGTGAGAGCGAGAGCGTGTGCGTGGGTGGGG
```

Max score in the matrix is: 22

Alignment of sequence 01 and sequence 02: (traceback\_array)

```
GGGTTCAGGAAAGAGTTTGAAGGATGCTG
x| |x|xx| |x|xxxxxxxxxx| |xx|xxx xx
TGGGTGCGGGAGTGAGAGCGAGAGCGT_GT
```

## 2. Part 01 – B

- For 1k random genomic sequence – Time required: 02 sec (by jobstats monsoon command)

```
[sk2354@ondemand /scratch/sk2354/A04/Part01]$ srun --mem=05GB -t 00:60:00 ./homework B ./hw3_dataset.fa ./test_genome.fasta

The number of arguments passed: 4
The first argument is: /scratch/sk2354/A04/Part01/./homework
The second argument is: B
The third argument is: ./hw3_dataset.fa
The fourth argument is: ./test_genome.fasta

Generate 1K, 10K, 100K, and 1M (million) completely random genomic sequences (50nt) to use as targets for alignment and use S
and record time to completion (in seconds / minutes).:

Total 70-character fragments: 427

*****

Sequence 01 and Sequence 02:

CAGAAACTCGAACGGTGCGCTTCTTGAGTCGACAACACACACGCTACCCT
ATTAAAGGTTTATACCTTCCAGGTAACAAACCAACCACTTCGATCTCTGTAGATCTGTTCTCTAAA

Max score in the matrix is: 20

Alignment of sequence 01 and sequence 02: (traceback_array)

AGAAAC_TCGAACGGTGCGCT_TCTTG_AG_TC
x |||| xxx|x x|xxxx xxx|x xx xx
C_AAACCAACCA_CTTTCGATCTCTGTAGAT

*****

Sequence 01 and Sequence 02:

CAGAAACTCGAACGGTGCGCTTCTTGAGTCGACAACACACGCTACCCT
ACAAATAGACATATAAGATAGACAAACATACAAACAGACAAAGATAAATAAAATATAAAATAAAAC

Max score in the matrix is: 24

Alignment of sequence 01 and sequence 02: (traceback_array)

_GGTGCGCTTCTTGAG_TCGAC_AACACACA
xxxx|xxxx|x|x xxxx ||xxxxx
CAAATAGACATATAAGATAGACAAACATAC
```

- b. For 10k random genomic sequence – Time required: 01 sec (by jobstats monsoon command)

```
=====
[sk2354@ondemand /scratch/sk2354/A04/Part01]$ make
g++ -c homework.cpp homework.h homework_BLAST.h
g++ -o homework homework.o
[sk2354@ondemand /scratch/sk2354/A04/Part01]$ srun --mem=10GB -t 00:60:00 ./homework B ./hw3_dataset.fa ./test_genome.fasta

The number of arguments passed: 4
The first argument is: /scratch/sk2354/A04/Part01/./homework
The second argument is: B
The third argument is: ./hw3_dataset.fa
The fourth argument is: ./test_genome.fasta

Generate 1K, 10K, 100K, and 1M (million) completely random genomic sequences (50nt) to use as targets for alignment and use SARS-
and record time to completion (in seconds / minutes).:

Total 70-character fragments: 427

*****

Sequence 01 and Sequence 02:

ACCACTTGCGGTAATCTTTTACGATACGTGTTATATGTCACACCTGATC
ATTAAAGGTTTATACCTTCCCAAGTAACAAACCAACCACTTCGATCTCTGTAGATCTGTTCTCTAAA

Max score in the matrix is: 31

Alignment of sequence 01 and sequence 02: (traceback_array)

CC_ACTTGCGGTAATCTTTTACGATACGTGTT_TATA
x| |xx|xxx|x xxx| |xx xxx x xxx| xxx
ACCAACTTTCGA_TCTCTTGT_AGA_T_CTGTTCTCT

*****

*****

Sequence 01 and Sequence 02:

ACCACTTGCGGTAATCTTTTACGATACGTGTTATATGTCACACCTGATC
ACAAATAGACATATAAGATAGACAAACATACAAACAGACAAAGATAAATAAAATATAAAATAAAAC

Max score in the matrix is: 16

Alignment of sequence 01 and sequence 02: (traceback_array)

ACGATACGTGTTTATAT_GTCACACCT
|xxxxxxxxxxx|x|xxx xxxxx|x|x
```

- c. For 100k random genomic sequence – Time required: 01 sec (by jobstats monsoon command)

```
[sk2354@ondemand /scratch/sk2354/A04/Part01 ]$ make
g++ -c homework.cpp homework.h homework_BLAST.h
g++ -o homework homework.o
[sk2354@ondemand /scratch/sk2354/A04/Part01 ]$ srun --mem=10GB -t 00:60:00 ./homework B ./hw3_dataset.fa ./test_genome.fasta

The number of arguments passed: 4
The first argument is: /scratch/sk2354/A04/Part01/./homework
The second argument is: B
The third argument is: ./hw3_dataset.fa
The fourth argument is: ./test_genome.fasta

Generate 1K, 10K, 100K, and 1M (million) completely random genomic sequences (50nt) to use as targets for alignment and use SARS-C
and record time to completion (in seconds / minutes).:

Total 70-character fragments: 427

*****

Sequence 01 and Sequence 02:

CTACGTGTGGTCTCTGGCTATTATTATAAAGGTACCGCTAACTCAC
ATTAAAGTTTATACCTTCCCAGGTAACAAACCAACCACTTCGATCTCTGTAGATCTGTTCTCTAAA

Max score in the matrix is: 25

Alignment of sequence 01 and sequence 02: (traceback_array)

TAAAGGT
|x| |x|x
TTAAAGG

*****

*****

Sequence 01 and Sequence 02:

CTACGTGTGGTCTCTGGCTATTATTATAAAGGTACCGCTAACTCAC
ACAAATAGACATATAAGATAGACAAACATACAAACAGACAAAGATAAATAAAATATAAAATAAAAAC

Max score in the matrix is: 19

Alignment of sequence 01 and sequence 02: (traceback_array)

TATTATTATAAAGGTACCGCTAAC
xx xxxxxxxx|x|xx x|xx|x|x
```

- d. For 1M random genomic sequence – Time required: 03 sec (by jobstats monsoon command)

```
[sk2354@ondemand /scratch/sk2354/A04/Part01]$ make
g++ -c homework.cpp homework.h homework_BLAST.h
g++ -o homework homework.o
[sk2354@ondemand /scratch/sk2354/A04/Part01]$ srun --mem=10GB -t 00:60:00 ./homework B ./hw3_dataset.fa ./test_genome.fasta

The number of arguments passed: 4
The first argument is: /scratch/sk2354/A04/Part01/./homework
The second argument is: B
The third argument is: ./hw3_dataset.fa
The fourth argument is: ./test_genome.fasta

Generate 1K, 10K, 100K, and 1M (million) completely random genomic sequences (50nt) to use as targets for alignment and use SAR
and record time to completion (in seconds / minutes).:

Total 70-character fragments: 427

*****

Sequence 01 and Sequence 02:

TATCCTCGAAAAAGGAGTTTCATACGTGAATACTCGGTAGTAGAGCAAAA
ATTAAAGGTTTATACCTTCCAGGTAACAACCAACCACTTTCGATCTCTGTAGATCTGTTCTCTAAA

Max score in the matrix is: 25

Alignment of sequence 01 and sequence 02: (traceback_array)

 _GGAGTTTCATACGTGAATACTCGGTAGTAGAGCAA
xx|xxx|xxxxxxx |x|x|x|xxxxx|x|x|
TTAAAGGTTTATACC__TTCAGGTAACAACCA

*****

*****

Sequence 01 and Sequence 02:

TATCCTCGAAAAAGGAGTTTCATACGTGAATACTCGGTAGTAGAGCAAAA
ACAAATAGACATATAAAGATAGACAAACATACAAACAGACAAAGATAAATAAAATATAAAATAAAAAC

Max score in the matrix is: 22

Alignment of sequence 01 and sequence 02: (traceback_array)

TAGTAGAGCAAAA
x|xxxxx xx|x
```

### 3. Part 02: A

```
[sk2354@ondemand /scratch/sk2354/A04/Part01]$ srun --mem=05GB -t 00:60:00 ./homework C ./hw3_dataset.fa ./test_genome.fasta

The number of arguments passed: 4
The first argument is: /scratch/sk2354/A04/Part01/./homework
The second argument is: C
The third argument is: ./hw3_dataset.fa
The fourth argument is: ./test_genome.fasta

:Having a BLAST - Implement a seed-based Smith Waterman:

Number of lines in read query dataset: 30
Number of lines in read genome dataset: 429
Size of HASH table :4194271
Collisions :236
Number of UNIQUE sequences: 25390
Sequence 01 and Sequence 02:

ACAAGTGTGCC
ACAAGTGTGCCCAAGTGTGCCTACGTGCTAGCGCTAACATAGGTTGTAAC

Max score in the matrix is: 20

Alignment of sequence 01 and sequence 02: (traceback_array)

AAGTGTGC
x|xxxxxx
CAAGTGTG

*****
Sequence 01 and Sequence 02:

CAACACAACAA
CAACACAACAAAACACAACAAAGGAGGTATGCACTGTTATCCGATTAC

Max score in the matrix is: 20

Alignment of sequence 01 and sequence 02: (traceback_array)

ACACAACA
|xxxx|xx
AACACAAC

*****
Sequence 01 and Sequence 02:
```



```

*****
Sequence 01 and Sequence 02:

CAACAGAATCT
CAACAGAATCTAACGTTAGATTTCCTAATATTACAACTTGTGCCCTTT

Max score in the matrix is: 20

Alignment of sequence 01 and sequence 02: (traceback_array)

ACAGAATC
|xxxx|xx
AACAGAAT

*****
Sequence 01 and Sequence 02:

AACATGGCAAG
AACATGGCAAGACATGGCAAGGAAGACCTTAAATTCCTCGAGGACAATT

Max score in the matrix is: 20

Alignment of sequence 01 and sequence 02: (traceback_array)

CATGGCAA
xxxx|xx|
ACATGGCA

*****
Sequence 01 and Sequence 02:

ACATGGCAAGG
ACATGGCAAGGCATGGCAAGGAAGACCTTAAATTCCTCGAGGACAATTA

Max score in the matrix is: 20

Alignment of sequence 01 and sequence 02: (traceback_array)

ATGGCAAG
xxx|xx|x
CATGGCAA

*****
Sequence 01 and Sequence 02:

CATGGCAAGGA
CATGGCAAGGAATGGCAAGGAAGACCTTAAATTCCTCGAGGACAATTAA

Max score in the matrix is: 20

```

#### 4. Part 02: B

- a. For 1k random 50mers sequence – Time required: 02 sec (by jobstats monsoon command)

```
[sk2354@ondemand /scratch/sk2354/A04/Part01]$ make
g++ -c homework.cpp homework.h homework_BLAST.h
g++ -o homework homework.o
[sk2354@ondemand /scratch/sk2354/A04/Part01]$ srun --mem=10GB -t 00:60:00 ./homework D ./hw3_dataset.fa ./test_genome.fasta

The number of arguments passed: 4
The first argument is: /scratch/sk2354/A04/Part01/./homework
The second argument is: D
The third argument is: ./hw3_dataset.fa
The fourth argument is: ./test_genome.fasta

Test your code on a set of 1K, 10K, 100K, and 1M (million) completely random 50-mers, aligning them to SARS-COV2 genome. How 1

Number of lines in read genome dataset: 429
Size of HASH table :4194271
Collisions :236
Number of UNIQUE sequences: 25390
Sequence 01 and Sequence 02:

TGGAGGAGATTCTATGAGATTCAGGGGAGTCCCAACGAACACCGCTG
ATTAAAGGTTTTTAAAGGTTTATACCTTCCAGGTAACAAACCAACCAAC

Max score in the matrix is: 23

Alignment of sequence 01 and sequence 02: (traceback_array)

AGATTCTATGAGATT_CAGGGGAGTCCCAACGAAC
|xxx| |xx xxxx| |xxx| |xxxxxxxx|x|xxx|x
AAGGT_TTA_TACCTTCCAGGTAACAAACCAACAA

*****
Sequence 01 and Sequence 02:

TGGAGGAGATTCTATGAGATTCAGGGGAGTCCCAACGAACACCGCTG
TTTCGATCTTTTCGATCTCTGTAGATCTCGAACTTTAAATCTGTGTG

Max score in the matrix is: 22

Alignment of sequence 01 and sequence 02: (traceback_array)

GAT_TCTATGAGAT_TCAGGGGAGTCCCAAC
xxx xxx |xxxxx xx xxx|xxxxx|x|x|x
CGATCTC_TTTCGATCT_CTTGTAGATCTCGAA

*****
```

- b. For 10k random 50mers sequence – Time required: 01 sec (by jobstats monsoon command)

```
[sk2354@ondemand /scratch/sk2354/A04/Part01]$ make
g++ -c homework.cpp homework.h homework_BLAST.h
g++ -o homework homework.o
[sk2354@ondemand /scratch/sk2354/A04/Part01]$ srun --mem=10GB -t 00:60:00 ./homework D ./hw3_dataset.fa ./test_genome.fasta

The number of arguments passed: 4
The first argument is: /scratch/sk2354/A04/Part01/./homework
The second argument is: D
The third argument is: ./hw3_dataset.fa
The fourth argument is: ./test_genome.fasta

Test your code on a set of 1K, 10K, 100K, and 1M (million) completely random 50-mers, aligning them to SARS-COV2 genome. How long?

Number of lines in read genome dataset: 429
Size of HASH table :4194271
Collisions :236
Number of UNIQUE sequences: 25390
Sequence 01 and Sequence 02:

GGGAACATCCTGCACTTATAGAATATCAAGTGATTCGTATATGGTTCTA
ATTAAAGGTTTTTAAAGGTTTATACCTTCCCAGGTAACAAACCAACCAAC

Max score in the matrix is: 24

Alignment of sequence 01 and sequence 02: (traceback_array)

ATATGGTTCT
x|xx|x|x|
TTAAAGGTTT

*****
Sequence 01 and Sequence 02:

GGGAACATCCTGCACTTATAGAATATCAAGTGATTCGTATATGGTTCTA
TTTCGATCTCTTTCGATCTCTGTAGATCTCGAACTTTAAATCTGTGTG

Max score in the matrix is: 29

Alignment of sequence 01 and sequence 02: (traceback_array)

TCGTATATGGTTC
|xx xxxxxx|x
TTC_GATCTCTTT

*****
Sequence 01 and Sequence 02:

GGGAACATCCTGCACTTATAGAATATCAAGTGATTCGTATATGGTTCTA
```

- c. For 100k random 50mers sequence – Time required: 01 sec (by jobstats monsoon command)

```
[sk2354@ondemand /scratch/sk2354/A04/Part01]$ make
g++ -c homework.cpp homework.h homework_BLAST.h
g++ -o homework homework.o
[sk2354@ondemand /scratch/sk2354/A04/Part01]$ srun --mem=10GB -t 00:60:00 ./homework D ./hw3_dataset.fa ./test_genome.fasta

The number of arguments passed: 4
The first argument is: /scratch/sk2354/A04/Part01/./homework
The second argument is: D
The third argument is: ./hw3_dataset.fa
The fourth argument is: ./test_genome.fasta

Test your code on a set of 1K, 10K, 100K, and 1M (million) completely random 50-mers, aligning them to SARS-COV2 genome. How lo

Number of lines in read genome dataset: 429
Size of HASH table :4194271
Collisions :236
Number of UNIQUE sequences: 25390
Sequence 01 and Sequence 02:

AAACGTGCTCGACAGCTATAAGTGTATGCGCCTGGCCATCCCCGGTTCAG
ATTAAAGGTTTTTAAAGGTTTATACCTTCCCAGGTAACAAACCAACCAAC

Max score in the matrix is: 19

Alignment of sequence 01 and sequence 02: (traceback_array)

TATAAG__TGTATGCGCCT_GGCCATCCCCGGT
|x||| |x|xxx|x| xxx|x|||x|x
TTAAAGGTTTTTAAAGGTTTATACCTTCCCAGG

*****
Sequence 01 and Sequence 02:

AAACGTGCTCGACAGCTATAAGTGTATGCGCCTGGCCATCCCCGGTTCAG
TTTCGATCTCTTTCGATCTCTGTAGATCTCGAACTTTAAATCTGTGTG

Max score in the matrix is: 25

Alignment of sequence 01 and sequence 02: (traceback_array)

GCCATCCC
xx|xxx|x
TTCGATCT

*****
Sequence 01 and Sequence 02:
```

- d. For 1M random 50mers sequence – Time required: 03 sec (by jobstats monsoon command)

```
=====
[sk2354@ondemand /scratch/sk2354/A04/Part01]$ make
g++ -c homework.cpp homework.h homework_BLAST.h
g++ -o homework homework.o
[sk2354@ondemand /scratch/sk2354/A04/Part01]$ srun --mem=10GB -t 00:60:00 ./homework D ./hw3_dataset.fa ./test_genome.fasta

The number of arguments passed: 4
The first argument is: /scratch/sk2354/A04/Part01/./homework
The second argument is: D
The third argument is: ./hw3_dataset.fa
The fourth argument is: ./test_genome.fasta

Test your code on a set of 1K, 10K, 100K, and 1M (million) completely random 50-mers, aligning them to SARS-COV2 genome. How long d

Number of lines in read genome dataset: 429
Size of HASH table :4194271
Collisions :236
Number of UNIQUE sequences: 25390
Sequence 01 and Sequence 02:

TAGTCTTCACCTACCGATCAGGTGTCAATTTGTCGAACGAGTTTCACCAT
ATTAAAGGTTTTAAAGGTTTATACCTTCCCAGGTAACAAACCAACCAAC

Max score in the matrix is: 24

Alignment of sequence 01 and sequence 02: (traceback_array)

GTCAATTTGTCGAACGAGTTT_CACC
 |x||xxxx|xxx|xx |x|| |xxx|
 _TTAAAGGTTTTAAA_GGTTTATAC

*****
Sequence 01 and Sequence 02:

TAGTCTTCACCTACCGATCAGGTGTCAATTTGTCGAACGAGTTTCACCAT
TTTCGATCTCTTTCGATCTCTTGATATCGAACTTTAAATCTGTGTG

Max score in the matrix is: 23

Alignment of sequence 01 and sequence 02: (traceback_array)

TC_A_C_CTACCGATCAGGTGTCAATTTGTCGAAC
 |x x x xxxxxxxxxxxxxx|xx xxxxxxxxxxx|x
 TTTCGATCTCTTTCGATCTCTTG_TAGATCTCGAA

*****
Sequence 01 and Sequence 02:
```

Comparison: Both results from Part 01 B and Part 02 B are same in timings as mentioned above individually.