

LSDS Assignment 4

Command for running Problem A : `srun homework_4 1 human.txt human_reads_trimmed.fa`

Command for running Problem B : `srun homework_4 2 human.txt human_reads_trimmed.fa`

A) Needleman Wunch(Random and Completely Random) :

- a. For each of your searches (10K, 100K, and 1M), how many 'hits' with up to 2 mismatches did you find?

For Random :

- a) For 10k : Number of hits are 9864
- b) For 100k : Number of hits are 99003
- c) For 1M : Number of hits are 990003

For Completely Random :

- a) For 10k : Number of hits are 8657
- b) For 100k : Number of hits are 87005
- c) For 1M : Number of hits are 870005

- b. For each of your searches (10K, 100K, and 1M), how long did the search take?

For Random :

- a) For 10k : Total time taken is 83765.3 seconds
- b) For 100k : Total time taken is 85938.48 seconds
- c) For 1M : Total time taken is 8617538.48 seconds

For Completely Random :

- d) For 10k : Total time taken is 276091 seconds
- e) For 100k : Total time taken is 2768202.26 seconds
- f) For 1M : Total time taken is 27689202.26 seconds

```
https://ondemand.hpc.nau.edu/pun/sys/dashboard/files/fs//scratch/au282/HW_4/part1_trimmed_output.txt - Google Chrome
ondemand.hpc.nau.edu/pun/sys/dashboard/files/fs//scratch/au282/HW_4/part1_trimmed_output.txt

File size of genome data is 3057186663
File size of query dataset is 27772294
=====
Generating random indices using random function
The best score is 32
Found for string TAACCCTAACCTAAC in query buffer
Found for string TAACCCTAACCTAAC in genome buffer
Total number of hits observed in 500 n-mers are 496
Total time taken for searching all hits for 500 n-mers are 1849.82 seconds
=====
The best score is 29
Found for string GGAAGAAAGCTTTCTG in query buffer
Found for string GGAAGAAAGCTTTCTG in genome buffer
Total number of hits observed in 1000 n-mers are 965
Total time taken for searching all hits for 1000 n-mers are 10882.8 seconds
=====
The best score is 26
Found for string TGCCCGGACCTGGCGG in query buffer
Found for string CGCCAGACCTGGCGG in genome buffer
Total number of hits observed in 10000 n-mers are 9864
Total time taken for searching all hits for 10000 n-mers are 83765.3 seconds
=====
Generating completely random indices using custom function
The best score is 28
Found for string ACCGTCCTGCTGGCG in query buffer
Found for string ACCGTCCTGCTGGCG which was generated completely random
Total number of hits observed in 500 n-mers are 437
Total time taken for searching all hits for 500 n-mers are 13046.7 seconds
=====
The best score is 26
Found for string TATGTTCTATATCTAG in query buffer
Found for string TATGTTCTATCGAG which was generated completely random
Total number of hits observed in 1000 n-mers are 853
Total time taken for searching all hits for 1000 n-mers are 28222.9 seconds
=====
The best score is 26
Found for string TAGCCTCCATCCATTA in query buffer
Found for string TAGCCTCCGTCCTTA which was generated completely random
Total number of hits observed in 10000 n-mers are 8657
Total time taken for searching all hits for 10000 n-mers are 276091 seconds
=====
```

c. How long would the search take for the entire subject dataset?

For Random -> human genome(approx. 3 billion):

Number of hits are 2659752401

Total time taken is 835.646 years(approx. 8 centuries)

For Completely Random -> human genome(approx. 3 billion):

Number of hits are 2659752401

Total time taken is 2684.34 years(approx. 20 centuries)

B) BLAST(Random and Completely Random):

- a. For each of your searches (10K, 100K, and 1M), how many 'hits' with up to 2 mismatches did you find?

For Random :

- a) For 10k : Number of hits are 21325
- b) For 100k : Number of hits are 209867
- c) For 1M : Number of hits are 1437699

For Completely Random :

- g) For 10k : Number of hits are 19692
- h) For 100k : Number of hits are 258651
- i) For 1M : Number of hits are 2871541

- b. For each of your searches (10K, 100K, and 1M), how long did the search take?

For Random :

- a) For 10k : Total time taken is 113.82 seconds
- b) For 100k : Total time taken is 246.79 seconds
- c) For 1M : Total time taken is 316.57 seconds

For Completely Random :

- j) For 10k : Total time taken is 120.24 seconds
- k) For 100k : Total time taken is 176.4 seconds
- l) For 1M : Total time taken is 373.45 seconds

https://ondemand.hpc.nau.edu/pun/sys/dashboard/files/fs//scratch/au282/HW_4/output_2.txt - Google Chrome

ondemand.hpc.nau.edu/pun/sys/dashboard/files/fs//scratch/au282/HW_4/output_2.txt

File size of genome data is 3057186663
File size of query dataset is 27772294

=====

Using BLAST WAY

Generating random indices using random function

The best score is 32

Found for string TAACCCTAACCTAAC in query buffer

Found for string TAACCCTAACCTAAC in genome buffer

Total number of hits observed in 10000 n-mers are 21325

Total time taken for searching all hits for 10000 n-mers are 113.82 seconds

=====

The best score is 28

Found for string GAAGATGAAGAATAAC in query buffer

Found for string GAGATGTAAGAATAAC in genome buffer

Total number of hits observed in 100000 n-mers are 209867

Total time taken for searching all hits for 100000 n-mers are 246.79 seconds

=====

The best score is 28

Found for string TTGATTGAGGGCTCAA in query buffer

Found for string TTGTTGAGGGCTGCAA in genome buffer

Total number of hits observed in 1000000 n-mers are 1437699

Total time taken for searching all hits for 1000000 n-mers are 316.57 seconds

=====

Generating completely random indices using custom function

The best score is 28

Found for string TGAGCAAGTTCTAGGC in query buffer

Found for string TGAGCAAGTCTCTAGC which was generated completely random

Total number of hits observed in 10000 n-mers are 19692

Total time taken for searching all hits for 10000 n-mers are 120.24 seconds

=====

The best score is 26

Found for string CGAGACCTCTACCACT in query buffer

Found for string CAAGACGTCTACCACT which was generated completely random

Total number of hits observed in 100000 n-mers are 258651

Total time taken for searching all hits for 100000 n-mers are 176.4 seconds

=====

The best score is 26

Found for string CCCAGCCAACGTTTGG in query buffer

Found for string CCCAGCCAACGTTTCT which was generated completely random

Total number of hits observed in 1000000 n-mers are 2871541

Total time taken for searching all hits for 1000000 n-mers are 373.45 seconds

=====

- c. How long would the search take for the entire subject dataset?

For Random -> human genome(approx. 3 billion):

Number of hits are 4371783947

Total time taken is 7.24 days

For Completely Random -> human genome(approx. 3 billion):

Number of hits are 8804688475

Total time taken is 9.05 days

- d. How does that compare with the benchmarks from problem 1, part B?

For completely random using NW it took us approximately 20 centuries but when we used BLAST for the same it took approximately 9.05 days which is a drastic change. This is because BLAST is a heuristic approach and we ignore the k-mers which never resulted to a hit sequence but we wasted our time in NW still evaluating for it.

Time complexity for NW : $O(N) * O(G) * O(n * n)$

Time complexity for BLAST : $O(N) * O(n) * O(n)$

Estimation calculations are mentioned below for both NW and BLAST:

VM Wunsch (Random)

	x_1		x_2	100k	1M	3B
$x(\text{kmers})$	500	1000	10000	100000	1000000	3057186663
$y(\text{hits})$	496	965	9864	99003	990003	3026614799
	y_1		y_2			

Slope formula

for 100k

$$y = 0.99x + 2.95$$

$$= 0.99(100000) + 2.95$$

$$= 99000 + 3$$

$$= 99003 \text{ hits}$$

for 1M

$$= 0.99(1000000) + 3$$

$$= 990000 + 3$$

$$= 990003 \text{ hits}$$

for 3B

$$= 0.99(3057186663) + 3$$

$$= 3026614799 \text{ hits}$$

$x(\text{kmers})$	$\overset{x_1}{500}$	1000	$\overset{x_2}{10000}$	100000	1000000	3057186663
$y(\text{time})$	13046.7 $\overset{y_1}{(y_1)}$	28222.9	276091 $\overset{y_2}{(y_2)}$	2768202.26	27689202.26	≈ 20 centuries

Slope formula $a =$ For 100k

$$\begin{aligned}
 y &= 27.69x - 797.74 \\
 &= 27.69(100000) - 797.74 \\
 &= 2769000 - 797.74 \\
 &= 2768202.26 \text{ seconds}
 \end{aligned}$$

km Wunsch
Completely Random

For 1M

$$\begin{aligned}
 &= 27.69(1000000) - 797.74 \\
 &= 27690000 - 797.74 \\
 &= 27689202.26 \text{ seconds}
 \end{aligned}$$

For 3B

$$\begin{aligned}
 &= 27.69(3057186663) - 797.74 \\
 &= 2684034 \text{ years} \\
 &\approx 20 \text{ centuries}
 \end{aligned}$$

$\frac{2000}{100}$

$x(\text{kmers})$	^(a1) 500	1000	^(a2) 10000	^{100k} 100000	^{1M} 1000000	^{3B} 3057186663
$y(\text{hits})$	^(y1) 437	853	^(y2) 8657	87005	870005	2659752401

Slope formula ^{for 100k}

$$\begin{aligned}
 y &= 0.87x + 4.37 \\
 &= 0.87(100000) + 5 \\
 &= 87000 + 5 \\
 &= 87005 \text{ hits}
 \end{aligned}$$

^{for 1M}

$$\begin{aligned}
 &= 0.87(1000000) + 5 \\
 &= 870000 + 5 \\
 &= 870005 \text{ hits}
 \end{aligned}$$

^{for 3B}

$$\begin{aligned}
 &= 0.87(3057186663) + 5 \\
 &= 2659752401 \text{ hits}
 \end{aligned}$$

km Wunsch
Random

copy
pic + jo

x (kmers)	(1) 10k 10000	100k 100000	(2) 1M 1000000	3B 3057186663
y (hits)	21325 (1)	209867	1437699 (2)	4371783947

Slope formula:

$$y = 1.43x + 7018.19 \quad (\text{For 3B})$$

$$= 1.43(3057186663) + 7018.19$$

$$\approx 4371783947 \text{ hits}$$

BLAST

Random

y (kmers)	$\overset{y_1}{10000}$	100000	$\overset{y_2}{1000000}$	3057186663
x (time)	$\overset{x_1}{13.82}$	246.79	316.57	7.24 days

Slope formula :

$$x = \frac{y + 545767.20}{4882.86}$$

$$= \frac{3057186663 + 545767.2}{4882.86}$$

$$\approx 7.24 \text{ days}$$

BLAST
RANDOM

x (kmers)	(1) 10k 10000	100k 100000	(2) 1M 1000000	3B 3057186663
y (hits)	19692 (81)	258654	2871541 (82)	8804688475

Slope formula:

$$y = 2.88x - 9114.56, \text{ (For 3B)}$$

$$= 2.88(3057186663) - 9114.56$$

$$\approx 8804688475 \text{ hits}$$

BLAST completely
Random

y (kmers)	^(y1) 10000	100000	^(y2) 1000000	3057186663
x (time)	^(x1) 120.24	126.4	^(x2) 373.45	~ 9.05 days

Slope formula :

$$x = \frac{y + 460114.13}{3909.80}$$

$$= \frac{3057186663 + 460114.13}{3909.80}$$

$$\approx 9.05 \text{ days}$$

BLAST completely
RANDOM

$\frac{km \text{ Wunsch}}{2 (kmex)}$	$\frac{completing \text{ seconds}}{2}$	1000	10000	100000	1000000	3057186663
y(time)	849.82	10882.8	83265.3	85938.48	8617538.48	835.646 years
	y_1		y_2			

Slope formula = For 100k

$$\begin{aligned}
 y &= 8.62x - 2461.52 \\
 &= 8.62(100000) - 2461.52 \\
 &= 862000 - 2461.52 \\
 &= 859538.48 \text{ seconds}
 \end{aligned}$$

For 1m

$$\begin{aligned}
 &= 8.62(1000000) - 2461.52 \\
 &= 8620000 - 2461.52 \\
 &= 8617538.48 \text{ seconds}
 \end{aligned}$$

For 3B

$$\begin{aligned}
 &= 8.62(3057186663) \\
 &\quad - 2461.52 \\
 &= 26352946573.54 \text{ seconds} \\
 &\approx 835.646 \text{ years} \\
 &\approx 8 \text{ centuries}
 \end{aligned}$$