# Question A

For the first 1 million reads, the total CPU time it takes to initialize (fill up) the array is **7 seconds** and the RAM usage is about **251 MB**. For the total 36 million reads, the estimated CPU time is **7*36=252 seconds** and estimated RAM usage is about **251*36=9036 MB**.

```
yg336@rain:~/inf503/homework1
The number of arguments passed: 3
The first argument is: main
The second argument is: problem1A
The third argument is: /common/contrib/classroom/inf503/hw_dataset.fa
Got 1000000 reads.
time used for 1000000 reads is: 7 s
memory released!
time used for release the memory is about: 44 ms
~
```
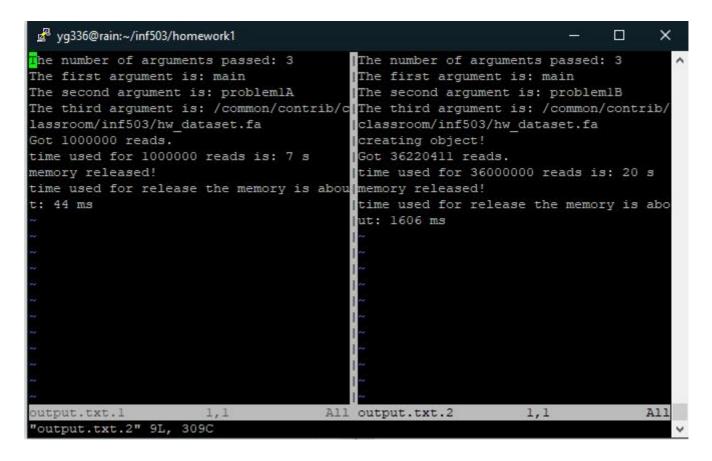
Fig. 1

| JobID | JobName | ReqMem | MaxRSS | ReqCPUS | UserCPU | Timelimit | Elapsed | State | JobEff |
|-------|---------|--------|--------|---------|---------|-----------|---------|-------|--------|
| 37172160 | lazy | 3.91G | 0.0M | 1 | 00:06.468 | 00:10:00 | 00:00:12 | COMPLETED | 2.0 |
| 37172189 | lazy | 3.91G | 3.84M | 1 | 00:07.674 | 00:10:00 | 00:00:13 | OUT_OF_MEMORY | - |
| 37172191 | lazy | 7.81G | 3.84M | 1 | 00:14.615 | 00:10:00 | 00:00:25 | OUT_OF_MEMORY | - |
| 37172221 | lazy | 3.91G | 0.0M | 1 | 00:07.633 | 00:10:00 | 00:00:14 | OUT_OF_MEMORY | - |
| 37172222 | lazy | 29.3G | 8.91G | 1 | 03:12.684 | 00:10:00 | 00:03:23 | COMPLETED | 32.12 |
| 37172250 | lazy | 9.77G | 8.91G | 1 | 03:12.608 | 00:10:00 | 00:03:23 | COMPLETED | 62.51 |
| 37172388 | lazy | 19.5G | 8.91G | 1 | 00:37.741 | 00:10:00 | 00:00:51 | COMPLETED | 27.05 |
| 37172389 | lazy | 19.5G | 8.91G | 1 | 00:37.834 | 00:10:00 | 00:00:47 | COMPLETED | 26.72 |
| 37172391 | lazy | 19.5G | 8.91G | 1 | 00:37.884 | 00:10:00 | 00:00:48 | COMPLETED | 26.8 |
| 37172392 | lazy | 19.5G | 0.0M | 1 | 00:00.048 | 00:10:00 | 00:00:02 | COMPLETED | 0.33 |
| 37172393 | lazy | 19.5G | 8.91G | 1 | 03:12.588 | 00:10:00 | 00:03:23 | COMPLETED | 39.71 |
| 37172396 | lazy | 19.5G | 3.97M | 1 | 00:06.512 | 00:10:00 | 00:00:12 | COMPLETED | 1.01 |
| 37172397 | problem | 19.5G | 0.0M | 1 | 00:06.590 | 00:10:00 | 00:00:11 | COMPLETED | 1.83 |
| 37172398 | problem1A | 19.5G | 0.0M | 1 | 00:06.505 | 00:10:00 | 00:00:11 | COMPLETED | 1.83 |
| 37172399 | problem1B | 19.5G | 3.97M | 1 | 00:16.804 | 00:10:00 | 00:00:26 | COMPLETED | 2.18 |
| 37172505 | problem1A | 19.5G | 3.97M | 1 | 00:05.665 | 00:10:00 | 00:00:18 | COMPLETED | 1.51 |
| 37172507 | problem1B | 19.5G | 0.0M | 1 | 00:10.797 | 00:10:00 | 00:00:25 | COMPLETED | 4.17 |
| 37172514 | problem1A | 19.5G | 3.97M | 1 | 00:06.534 | 00:10:00 | 00:00:20 | COMPLETED | 1.68 |
| 37172515 | problem1B | 19.5G | 3.84M | 1 | 00:16.886 | 00:10:00 | 00:00:26 | COMPLETED | 2.18 |
| 37172516 | problem1A | 19.5G | 0.0M | 1 | 00:06.388 | 00:10:00 | 00:00:11 | COMPLETED | 1.83 |
| 37172517 | problem1A | 19.5G | 0.0M | 1 | 00:06.349 | 00:10:00 | 00:00:12 | COMPLETED | 2.0 |
| 37172518 | problem1A | 19.5G | 3.97M | 1 | 00:06.599 | 00:10:00 | 00:00:12 | COMPLETED | 1.01 |
| 37172520 | problem1A | 19.5G | 0.0M | 1 | 00:06.390 | 00:10:00 | 00:00:11 | COMPLETED | 1.83 |
| 37172521 | problem1B | 19.5G | 0.0M | 1 | 00:11.726 | 00:10:00 | 00:00:30 | COMPLETED | 3.33 |
| 37172524 | problem1A | 19.5G | 251M | 1 | 00:06.506 | 00:10:00 | 00:00:38 | COMPLETED | 3.8 |
| 37172526 | problem1B | 19.5G | 8.91G | 1 | 00:16.837 | 00:10:00 | 00:00:46 | COMPLETED | 26.63 |

Fig. 2

# Question B

For the entire 36 million reads, the RAM usage is about **8.91 GB** (shown in Question A) and the CPU time is about **20 seconds**.

Fig. 3

Compared with the previously estimated time, there is a big gap. In order to understand why there's such a big mismatch, the code in **test.cpp** was used to test.

You may use the following commands to compile and run the code in a Linux system computer:

```
cd /home/yg336/inf503/homework1
g++ -o test.out test.cpp
sbatch test.sh
```

The results are shown in the following figure:



Fig. 4

From the figure we know that if we just allocate memories for the reads, the time consumption will increase linearly. However, if we write data into memories, the time consumption will not increase linearly. Comparing these two processes, we can know that **the time it takes to read data from a file and write to memory is uncertain and that will possibly lead to inaccurate estimation.**

If we compare the results from **test.cpp** to the results shown in Fig. 3, we can get another process that may influence the CPU time consumption: **creating the FASTA_readset object**.

## Question C

Fig. 5

## Question D



Fig. 6

From Fig. 6 we can know that it takes **1606 ms** to deallocate the memories. It makes sense. The time increases linearly when comparing with the time shown in Fig. 1.

## Question E

The quick sort was used to sort the segments, so the time complexity is O(nlogn). The first 20 lines are:

Fig. 7