

COMP5112 / MSBD5009 Parallel Programming (Spring 2022)

Assignment 2: Pthread Programming

Submission deadline: **23:59 pm on April 8, 2022**

I. General Notes

1. This assignment counts for 10 points.
2. This is an individual assignment. You can discuss with others and search online resources, but your submission must be your own code. Plagiarism checking will be enforced, and penalty will be given to all students involved in an incident of academic dishonesty.
3. Add your ***name, student id, and email*** at the first two lines of comments in your submission.
4. Submit your assignment through Canvas before the deadline.
5. All submission will be compiled and tested uniformly on given machines of the course (CSE Lab 2 machines for COMP5112 and Azure virtual machines for MSBD 5009).
6. Please direct any inquiries about this assignment to the designated TAs listed in the CANVAS assignment page.
7. No late submission will be accepted!

II. Problem Description

This assignment is on the Pthread parallelization of a ***super-mer*** generation algorithm given an input list of DNA fragments, called ***reads***. Each read is represented as a string of length n , and the characters in the string include 'A', 'T', 'C', and 'G'. A ***k-mer*** is a substring of length k of a read, so a read of length n contains $n - k + 1$ k-mers. A ***minimizer*** of length p of a k-mer is the lexicographically smallest length- p substring of the k-mer. Finally, a ***super-mer*** is generated by merging consecutive k-mers in the read that have the same minimizer. Two k-mers are consecutive in a read if their starting positions in the read are consecutive. Figure 1 illustrates the super-mer generation algorithm (Algorithm 1) with a simple example.

Input: $k = 9, p = 5, n = 14$		Super-mer Generation (<u>minimizers</u> are underlined)
<i>Read</i> = CAAATTACTGCATA		
$i=0$ (k-mer #1)	CAAATTACT	<i>minimizer</i> \leftarrow AAATT, <i>minimizer_beg_pos</i> \leftarrow 1, <i>supermer_beg_pos</i> \leftarrow 0
$i=1$ (k-mer #2)	AAATTACTG	no update on minimizer
$i=2$ (k-mer #3)	AATTACTGC	<i>minimizer</i> \leftarrow new_minimizer \leftarrow AATTA, <i>minimizer_beg_pos</i> \leftarrow 2
$i=2$ (super-mer #1)	CAAATTACTG	generate current super-mer, <i>supermer_beg_pos</i> \leftarrow 2
super-mer #1 is made up of k-mer #1 and #2, minimizer = AAATT		
$i=3$ (k-mer #4)	ATTACTGCA	<i>minimizer</i> \leftarrow new_minimizer \leftarrow ACTGC, <i>minimizer_beg_pos</i> \leftarrow 7
$i=3$ (super-mer #2)	AATTACTGC	generate current supermer, <i>supermer_beg_pos</i> \leftarrow 3
super-mer #2 is made up of k-mer #3 only, minimizer = AATTA		
$i=4$ (k-mer #5)	T TACTGCAT	no update
$i=5$ (k-mer #6)	TACTGCATA	no update
(super-mer #3)	ATTACTGCATA	<i>supermer_beg_pos</i> $\neq n - k$, generate the last super-mer
super-mer #3 is made up of k-mer #4 #5 #6, minimizer = ACTGC		

Figure 1: Illustration of super-mer generation (The minimizers are underlined)

III. Sequential Algorithm

Algorithm 1: Super-mer Generation (sequential)

```
Input:    read string  $S = s_0, s_1, s_2, \dots, s_{n-1}$ 
           k-mer length  $k$ , minimizer length  $p$ 

1. Procedure GenerateSupermer ( $S, k, p$ ):
2.    $minimizer$  = the minimum  $p$ -substring of  $S[0, k - 1]$ 
3.    $minimizer\_beg\_pos$  = the starting position of minimizer in  $S$ 
4.    $supermer\_beg\_pos = 0$ 
5.   for  $i$  from 1 to  $n - k$  do:
6.     if  $i > minimizer\_beg\_pos$  then:
7.        $new\_minimizer$  = the minimum  $p$ -substring of  $S[i, i + k - 1]$ 
8.       update  $minimizer\_beg\_pos$ 
9.     else:
10.      if  $S[i + k - p, i + k - 1] \leq minimizer$  then:
11.         $new\_minimizer = S[i + k - p, i + k - 1]$ 
12.        update  $minimizer\_beg\_pos$ 
13.      end if
14.    end if
15.    if  $new\_minimizer \neq minimizer$  then:
16.       $minimizer = new\_minimizer$ 
17.      save current super-mer:  $S[supermer\_beg\_pos, i + k - 1]$ 
18.       $supermer\_beg\_pos = i$ 
19.    end if
20.  end for
21.  if  $supermer\_beg\_pos \neq n - k$  then:
22.    save the last super-mer:  $S[supermer\_beg\_pos, n]$ 
23.  end if
24. end procedure
```

IV. Your Implementation Task

The code skeleton “gensupermer_pthread.cpp” is provided. In this assignment, your task is to implement a Pthread parallelization of the super-mer generation algorithm. Each thread handles a subset of the input reads and generates super-mers.

1. Implement the thread function “parallel” for each thread to generate super-mers. You can call the provided function “read2supermers”.
2. In the gensupermers function, initiate threads, call the thread function, and store results into the vector of strings all_supermers.

Notes:

1. **CSR Format:** The CSR (compressed sparse row) format stores all rows in a single data array and has an offset array to point to the beginning position of each row in the data array. For example, we store all reads in the array *read_CSR*, and the index of each read in the array *read_CSR_offset*. Figure 2 shows an example of four reads AC, ACT, AG, and GCT stored in CSR.

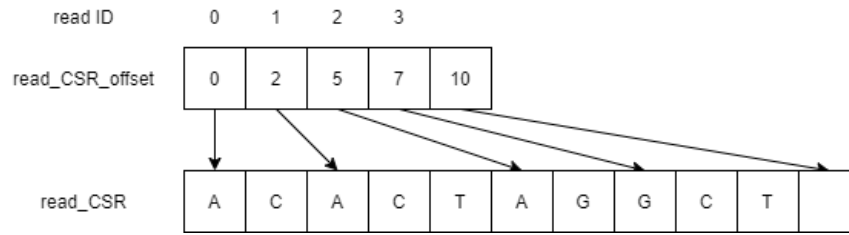


Figure 2: An Example of CSR Format

2. We provide the function: *read2supermers*. The function takes one read as input, generates all super-mers in the read, and stores them in a `vector<string>`. You can reference its usage in the sequential code.

V. The Given Package

gensupermer_thread.cpp	The Pthread program skeleton for you to fill. You need to write your own code only in this file.
gensupermer.hpp	A header file containing a function “gensupermers”.
main.cpp	The main function of the pthread version.
gensupermer_sequential.cpp	The sequential version of super-mer generation. It is for your reference and result comparison.
dataset/*.txt	Each text file is a test dataset with each line containing a read.
result/	An empty folder for super-mer output.
utilities.hpp	A header file containing utility functions such as file loading, result output, and result comparison.

VI. Compile and Run

Notes:

1. We will ***not*** apply any compiler optimization during the assessment (e.g., `-O2 -O3 ...`).
2. ***No*** external library or header file is allowed (e.g., Boost or other third-party libraries).

An example of input file and corresponding output file:

The input file contains multiple lines of reads. The following is an example containing 5 reads.

```
GATAACGAGTTCTAGAAACTGGGTCCC
AGATCAGCAAACCCTGAGAAAAA
AGAAAAATTAGCAATAATTAGCAGTGTTTCATAACA
AGAAGACAACTGGGCCCGGGGGAC
GAGGGAGAACCTGATTTCCAGAGT
```

Example input: dataset1.txt

The output will be the super-mers generated from the reads. The order of the super-mers in the output file is lexicographic. For example, the 7th and the last super-mers in the example output file are generated from the first read in the example input file. The 7th super-mer contains 7 consecutive k-mers. These 7 k-mers share the same minimizer AAAACTG.

```
AAAATTAGCAATAATTAGCAG
AAATTAGCAATAATTAGCAGT
AATTAGCAATAATTAGCAGTGTTTCATAA
AGAAAAATTAGCAATAATTAGCA
AGAAGACAACCTGGGCCCGGGGAC
AGATCAGCAAACCCTGAGAAAAA
ATAACGAGTTCTAGAAAAGTGGGTCCC
ATAATTAGCAGTGTTTCATAACA
GAGGGAGAACCTGATTTCCAGAGT
GATAACGAGTTCTAGAAAAGT
```

Example output (k=21, p=7): my_gs_output.txt

Example of compiling and running the sequential program:

```
# Compile:
g++ gensupermer_sequential.cpp -o seq_gs -std=c++11
# Run and save super-mers to text file:
# (e.g., k=21, p=7, data file is ./dataset/dataset1.txt, and result file is ./result/seq_gs_output.txt):
./seq_gs 21 7 ./dataset/dataset1.txt ./result/
# Run without saving super-mers to text file:
./seq_gs 21 7 ./dataset/dataset1.txt
```

Example of compiling and running your Pthread program:

```
# Compile:
g++ -std=c++11 -lpthread gensupermer_pthread.cpp main.cpp -o pthread_gs
# Run:
./pthread_gs 21 7 2 ./dataset/dataset1.txt <correctness_check> <result_folder
(optional)>
# <correctness_check> can be 0 or 1. If 1, the program will automatically check the
correctness by comparing all_supermers with the sequential version. <result_folder> is
optional; if not provided, the super-mers will not be saved to any file.
```

VII. Submission and Evaluation

You only need to submit your completed *gensupermer_pthread.cpp* to Canvas before the deadline. You can add your code only in the specified section. Your program should not have any extra output to the result file or the standard output.

We will evaluate your program in different parameter settings. We will vary p and k ($10 \leq 2p \leq k < 23$). The length of each read is greater than 22. Both the total length of all reads and that of the generated super-mers will be less than INT_MAX (0x7FFFFFFF). The number of threads will be set to values ≤ 8 . We will prepare multiple sets of evaluation arguments and datasets. Your grade will be marked based on both correctness and best speedup achieved.