

LB/TH/08/2023

DCS 03/51

ACCELERATING K-MER COUNTING FOR GENOMIC ANALYSIS

University of Moratuwa



TH5106

LIBRARY
UNIVERSITY OF MORATUWA, SRI LANKA
MORATUWA

Gunavaran Brihadiswaran

(208032P)

Thesis submitted in partial fulfillment of the requirements for the
degree of Master of Science

Department of Computer Science and Engineering

University of Moratuwa

Sri Lanka

June 2021

TH5106

+

CD ROM

TH5106

DECLARATION

I declare that this is my own work and this thesis does not incorporate without acknowledgement any material previously submitted for a Degree or Diploma in any other University or institute of higher learning and to the best of my knowledge and belief it does not contain any material previously published or written by another person except where the acknowledgement is made in the text.

Also, I hereby grant to the University of Moratuwa the non-exclusive right to reproduce and distribute my thesis, in whole or in part in print, electronic or other media. I retain the right to use this content in whole or part in future works (such as articles or books).

Signature: *UOM Verified Signature*

Date: 10.10.2021

The above candidate has carried out research for the Master's thesis under my supervision.

Name of the supervisor: Prof. V. S. D. Jayasena

Signature of the supervisor *UOM Verified Signature*

Date: 13/10/2021



Abstract

k -mer counting is the process of counting k length substrings in a sequence. It is an important step in many bioinformatics applications including genome assembly, sequence error correction, and sequence alignment. Even though generating k -mer histograms seems simple and straightforward, processing large datasets efficiently with limited resources, especially memory, is very challenging. As the advancements in next-generation sequencing technologies have resulted in a tremendous growth of genomic data, it is inevitable for k -mer counters to be faster and more efficient. A lot of work has been done in the past decade to optimize k -mer counting.

Frigate, a fast and efficient tool capable of counting and querying k -mers is presented. Its in-memory design utilizes multithreaded, lock-free data structures to improve performance. Thread synchronization is handled using the compare-and-swap technique. The parallel processing pipeline of Frigate is the result of careful performance engineering and design. Frigate was developed with the emphasis on values of k less than 20, aiming to maximize performance by employing different algorithms for different ranges of k values.

The performance of Frigate was compared with six state-of-the-art k -mer counters: Jellyfish, DSK, Gerbil, CHTKC, KMC2, and KMC3, using two real-world datasets. The experiments were carried out for k values of 10, 15, and 17 using a different number of threads in the range [1, 32]. The results show that Frigate achieves a comparable performance or up to 2-3x speedup compared to its competitors, especially for large datasets. The k -mer counters were analyzed based on the running time, amount of memory used, and scalability. The correctness of Frigate was evaluated by comparing the k -mer frequency histogram with those of other k -mer counters.

Frigate is written in C and freely available at <https://github.com/Gunavaran/frigate> under MIT license.

Keywords: K -mer counting, Genome analysis, Performance engineering, Parallel computing

ACKNOWLEDGEMENT

I would like to express my sincere gratitude to my supervisor Prof. Sanath Jayasena for giving me the opportunity to be involved in this research and for his invaluable support and guidance throughout this research. I am always thankful for his efforts and kindness when we were struggling to acquire research grants.

I am also thankful to Systems Engineer Mr. Randika Pathirana and all the staff of the Department of Computer Science and Engineering at the University of Moratuwa who helped me a lot in getting access to the university resources.

Last but not least, I would like to thank my parents and friends for their immense support.

TABLE OF CONTENTS

DECLARATION	ii
Abstract	iii
ACKNOWLEDGEMENT	iv
TABLE OF CONTENTS	v
LIST OF FIGURES	vii
LIST OF TABLES	viii
LIST OF ABBREVIATIONS	ix
1 INTRODUCTION	1
1.1 DNA and DNA Sequencing	1
1.2 <i>K</i> -mers and <i>K</i> -mer Counting	2
1.3 The Problem	3
1.4 Motivation	4
1.5 Objectives	4
1.6 Contribution.....	4
1.7 Thesis Outline.....	5
2 LITERATURE REVIEW	6
2.1 In-memory Approaches	6
2.2 Disk-based Approaches	8
2.3 Approximate <i>K</i> -mer Counters	14
2.4 Distributed <i>K</i> -mer Counting.....	15
2.5 <i>K</i> -mer Counters Utilizing Modern Hardware.....	15

3	METHODOLOGY	17
3.1	FASTQ Files.....	17
3.2	Reader Thread and <i>K</i> -mer Queues	18
3.3	Binary Encoding.....	19
3.4	Counter Array.....	20
3.5	Lock-free Counting	20
3.6	Canonical <i>K</i> -mer Counting.....	22
3.7	Storing <i>K</i> -mer Counts on the Disk	26
3.8	Querying <i>K</i> -mers	28
3.9	Frigate Tool and its Functionalities.....	29
4	RESULTS	32
4.1	Experimental Setup	32
4.2	Execution Time	34
4.3	Performance of Writing Output.....	39
4.4	Result Validation.....	40
5	CONCLUSION AND FUTURE WORK	42
	REFERENCES.....	44
	APPENDICES	47

LIST OF FIGURES

Figure 1.1: DNA double helix structure Source: See Appendix II	1
Figure 1.2: All possible 4-mers in a DNA sequence.....	2
Figure 2.1: CHTKC hash table structure Source: see Appendix II.....	8
Figure 2.2: Minimum substrings	9
Figure 2.3: Minimum substrings partitioning.....	10
Figure 2.4: KMC2 parallel processing workflow Source: See Appendix II	12
Figure 2.5: Workflow of Gerbil's distribution phase Source: See Appendix II	12
Figure 2.6: Workflow of Gerbil's counting phase Source: See Appendix II	13
Figure 3.1: Parallel processing pipeline of Frigate	17
Figure 3.2: Two sample entries in a FASTQ file.....	18
Figure 3.3: Structure of an entry in a k -mer queue	18
Figure 3.4: Binary encoding of a 12-mer	19
Figure 3.5: Structure of counter array	20
Figure 3.6: Structure of DNA with the complementary pairs Source: See Appendix II ..	22
Figure 3.7: Steps involved in finding the canonical sequence of TACAGT	23
Figure 3.8: Encoding the first k -mer of a read and its reverse complement	24
Figure 3.9: Encoding consecutive k -mers through bitwise operations	25
Figure 3.10: Encoding the reverse complements of consecutive k -mers through bitwise operations	26
Figure 3.11: Two approaches for dividing the work among four writer threads.	27
Figure 3.12: Number of non-zero counts encountered by eight writer threads after 15-mer counting for <i>F. vesca</i> (left) and <i>H. sapiens</i> (right) datasets.	27
Figure 3.13: A sample output of histogram functionality	30
Figure 3.14: A sample output of dump command.....	31
Figure 4.1: Running times for HS dataset ($k=15$)	37
Figure 4.2: Running times for HS dataset ($k=17$)	38
Figure 4.3: Time taken by different number of threads to write the k -mer counts to the disk after 17-mer counting of FV dataset	40

LIST OF TABLES

Table 2.1: Summary of k -mer counting tools.....	14
Table 3.1: Available options for count functionality	29
Table 4.1: Test machine configuration.....	32
Table 4.2: Datasets.....	33
Table 4.3: K -mer counting execution time (in sec) for <i>F. vesca</i> ($k = 10$).....	34
Table 4.4: K -mer counting execution time (in sec) for <i>H. sapiens</i> ($k = 10$)	34
Table 4.5: K -mer counting execution time (in sec) for <i>F. vesca</i> ($k = 15$).....	35
Table 4.6: K -mer counting execution time (in sec) for <i>H. sapiens</i> ($k = 15$)	35
Table 4.7: K -mer counting execution time (in sec) for <i>F. vesca</i> ($k = 17$).....	35
Table 4.8: K -mer counting execution time (in sec) for <i>H. sapiens</i> ($k = 17$)	36
Table 4.9: Comparison of k -mer counting execution times (in sec) for <i>H. sapiens</i> ($k = 15$) with 2GB memory.....	39
Table 4.10: K -mer frequency histogram for FV dataset ($k=15$).....	41

LIST OF ABBREVIATIONS

DNA	Deoxyribonucleic Acid
NGS	Next Generation Sequencing
GPU	Graphics Processing Unit
I/O	Input/Output
SSD	Solid State Drive
DRAM	Dynamic Random Access Memory
FPGA	Field-Programmable Gate Array
NVM	Non-Volatile memory