



OPTIMIZING VARIANT CALLING PERFORMANCE FOR HOTSPOT CASES BASED ON ION TORRENT SEQUENCING TECHNOLOGY

Basma Nasser Abd Elsalam Abd Elfattah

A Thesis Submitted to the
Faculty of Engineering at Cairo University
in Partial Fulfillment of the
Requirements for the Degree of
MASTER OF SCIENCE

in

Biomedical Engineering & Systems

FACULTY OF ENGINEERING, CAIRO UNIVERSITY GIZA, EGYPT





OPTIMIZING VARIANT CALLING PERFORMANCE FOR HOTSPOT CASES BASED ON ION TORRENT SEQUENCING TECHNOLOGY

by

Basma Nasser Abd-Elsalam Abd-Elfattah

A Thesis Submitted to the
Faculty of Engineering at Cairo University
in Partial Fulfillment of the
Requirements for the Degree of
MASTER OF SCIENCE

in

Biomedical Engineering & Systems

Under the Supervision of

Prof.Dr. Ayman Eldeib	Dr. Mohamed I. Abouelhoda
Professor of Biomedical Engineering	Associate Professor
Systems & Biomedical Engineering	Systems & Biomedical Engineering
Faculty of Engineering, Cairo University	Faculty of Engineering, Cairo University

Acknowledgments

I would like to thank Prof. Ayman Eldeib for providing me with the vision and giving me time to hear me. Great thanks to Dr. Mohamed I. Abouelhoda, for his guidance and patience. Without his persistent help, this thesis would not have been possible.

Thanks for all my colleagues and friends who help and support me. Special thanks for my family for their support over the last couple of years.

Dedication

Table of Contents

A	cknow	ledgments	i			
Li	st of T	ables	vi			
Li	st of F	ligures	vii			
Al	bbrevi	ations	viii			
Al	bstract	t	ix			
1	Intro	oduction	1			
	1.1	Problem Definition	1			
	1.2	Thesis Objective	2			
	1.3	Thesis Organization	2			
2	Biological Background					
	2.1	Genetic Background	3 3			
	2.2	DNA sequencing	4			
	2.3	1000 Genome projects	5			
	2.4	HapMap project	5			
	2.5	BRCA genes	5			
	2.6	Genetic variation	6			
		2.6.1 Single-nucleotide polymorphism	6			
		2.6.2 Insertion/Deletion polymorphism	6			
	2.7	Data Sequences Format	7			
		2.7.1 Fastq	7			
		2.7.2 FASTA	8			
		2.7.3 BAM/SAM File Format	8			
		2.7.4 VCF Format	10			
		2.7.5 Ion torrent Hotspot File	13			
		2.7.6 Target File	14			
	2.8	DNA Sequencing Generations	14			
	2.9	Next Generation Sequencing	15			
		2.9.1 NGS Strength	16			
		2.9.2 NGS Limitations	16			
		2.9.3 NGS Workflow Overview	16			
	2.10	Next Generation Sequencing Technologies	19			

		2.10.1	Roche 454 sequencing		 20
		2.10.2	Illumina Sequencing		 20
		2.10.3	Ion torrent: Proton / PGM sequencing		 21
		2.10.4	SOLiD Sequencing		 22
3	Rece		ant Calling Pipeline		23
	3.1		ew of variant caller pipeline		
	3.2	Read N	Mapping/Alignment tools		 24
		3.2.1	Burrows-Wheeler Aligner (BWA)		 25
		3.2.2	TMAP		 25
	3.3	Variant	t Calling		 26
		3.3.1	GATK		 26
			GATK Features		 26
			GATK workflow		
			GATK Releases		
		3.3.2	Torrent Variant Caller		
		0.0.2	TVC Workflow		
			Torrent Variant Detection Algorithms		
	3.4	GATK	and TVC Summary		
	J. 4	UATK	and I ve Summary	· · · · · · · · · · · ·	 31
4	Con	tributio	n and Methods		39
	4.1	Sequen	ncing implementation		 41
	4.2	-	l Computing		
		4.2.1	Overview		
		4.2.2	Advantages of Parallel Computing		
	4.3		lism Implementation Types		
	4.4		I Implementation		
	4.5		ds		
	4.5	4.5.1			
		4.5.1	Datasets		
			Pipelines		
		4.5.3	Aligning Tools		
			BWA Mem		
			TMAP		
			Alignment File Processing		
		4.5.4	Variant Calling Tools		
			HaplotypeCaller		48
			Torrent Variant Caller (TVC)		 48
			Hotspot Preparation		 49
			Modified Hotspot (mHotspot) Preparation .		 49
			Modified TVC (mTVC)		 49
		4.5.5	Selected Parameters		 50
		4.5.6	Parallel Algorithm		 51
		4.5.7	Hardware Specs		51
			-		
5			Discussion		52
	5.1	Sequen	nce Data Sets and Variant Calling Pipelines .		 52
	5.2	Variant	t Callers Comparison		 52
	5.3	Shared	Variants Validation		 53
	5 4	Variant	ts selection and filters		56

	5.5	Performance Results	56
6	Con	clusion and Future Work	58
	6.1	Conclusion	58
	6.2	Future work	59
R	eferen	ices	60
$\mathbf{A}_{]}$	ppend	lix A	62
$\mathbf{A}_{]}$	ppend	lix B	64
A	ppend	lix C	67

List of Tables

2.1	Summary of next generation sequencing technologies	22
4.1	Summary of Data Used	47
5.1	Variant callers pipeline results	53
5.2	Results of NA12878 using TVC with different hotspot files	55
5.3	Results of BRCA2 using TVC with different hotspot files	55

List of Figures

2.1	DNA Helix strand	4
2.2	Single-nucleotide polymorphism example.[7]	6
2.3	Insertion/Deletion polymorphism example	7
2.4	A chart demonstrating how the speed of DNA sequencing innovations	
	has expanded since the early strategies in the 1980s	15
2.5	NGS workflow from nucleic acid extraction to variant annotation	17
2.6	NGS library is prepared by fragmenting a gDNA sample and ligating	
	specialized adaptors to both fragment ends	18
2.7	library is loaded into a flow cell and the fragments hybridize to the	
	flow cell surface. Each bound fragment amplified into a clonal cluster	
	through bridge amplification	18
2.8	Reads are aligned to a reference sequence with read aligner software.	
	After alignment, differences between the reference genome and the	
	newly sequenced reads can be identified	19
3.1	Unified steps of the variant calling pipelines	24
3.2	GATK workflow.	27
3.3	Preprocessing Steps	29
3.4	Variant discovery steps	31
3.5	Evaluation steps	32
3.6	Flow diagram of torrent variant caller	35
3.7	Torrent Variant algorithm flow	37
4.1	Updated variant caller pipeline. (mTVC)	40
4.2	Serial implementation diagram	42
4.3	sequential computing, algorithm is constructed and implemented as a	
	serial stream of instructions. These instructions are executed on a CPU	
	on one machine. Only one instruction executes at a time, after that in-	
	struction is finished, the next one is executed	42
4.4	Parallel computing - breaking the problem into independent parts so that	
	each processing element can execute its part of the algorithm simultane-	
	ously with the others	43
4.5	Parallel implementation diagram	46
5.1	Intersection of variants from running TVC on NA12878 sample with	
	hapmap 3.3 and mTVC using the same parameters	
5.2	Execution time for NA12878 and BRAC2	57

Abbreviations

BAM Binary Sequence Alignment/Map.

BRCA BReast CAncer genes.

DNA Deoxyribonucleic Acid.

ENCODE Encyclopedia of DNA Elements.

GATK Genome Analysis Toolkit.

GIAB Genome in a bottle

Hapmap Haplotype Map.

HC HaplotypeCaller

InDel Insertions and DELetion.

mTVC Modified torrent variant caller script.

NGS Next Generation Sequencing.

PCR Polymerase chain reaction

PE Paired-end.

PGM Personal genome machine.

SAM Sequence Alignment/Map.

SNP Single nucleotide polymorphism.

TMAP Torrent Mapping Alignment Program.

TVC Torrent Variant Caller.

VCF Variant call Format.

VQSR Variant Quality Score Recalibration

Abstract

In last few years, Next generation sequencing (NGS) revolutionized in DNA sequencing research. The latest major technologies released are Illumina and Ion Torrent. In this study, we analyze the performance of calling variants using both platforms. In order to compare between both platforms, we used two sequenced data sets from ion community, which contain flow spaces required by Ion Torrent. We are concerned with the execution time of both platforms. Ion torrent detects genome variants faster than Illumina but with low accuracy. Moreover, we found that Ion Torrent called slightly more variants but with higher false positive rate. The hotspot option provided by ion torrent variant caller provides more accurate variant calling positions as the filtration restricted to specific positions in the genome, but this leads to time consumption. Here, we enhanced the execution time to attain the accurate positive variants using torrent variant caller by using the NOCALL variants as hotspot regions to torrent variant caller (TVC). We developed two dependent packages the first one to create the new hotspot file (mHotspot) required by TVC and the second one to rerun the TVC with mHotspot file generated to get the exact positions of variants.

Chapter 1

Introduction

Sequencing technologies are evolving rapidly; in 2011, a lot of sequencing technologies were developed. The advances of sequencing technologies make it possible to detect enormous number of potential disease-causing variants; also, next generation sequencing (NGS) data has been used for diagnostic purposes. This is due to the improvements made to the bioinformatics tools used to analyze the massive data produced by NGS instruments.

NGS technologies use same workflow for searching for mutations in a patient. Align the raw data to human reference genome, and identify single nucleotide variants and INDELs that cause the phenotype of interest. We have to decide the best tools to use for each stage of the pipeline. There are various tools for sample sequencing, but the most important parts are to select the proper aligner and variant caller to achieve the optimal results for SNPs and INDEL variants.

Different pipelines with different combinations between aligners and variant callers identify Gene variants (SNPs and INDELs). The most widely used aligners are BWA. GATK and torrent variant caller (TVC) are popular tools for variant calling [1].

In this thesis, we describe the pipelines and typical tools for read alignment and variant callers. In addition, we will elaborate in details the results of using ion torrent workflow for detecting variants from two different samples (NA12878 and BRCA2). In addition, we will describe new workflow for TVC in order to enhance time of execution using our plugin mTVC.

1.1 Problem Definition

Several studies have conducted comparison between different variant callers. From these studies, quail et al. BMC Genomics 2012 [2] mentioned that ion torrent has poor performance with only 65% of genome is covered with high quality reads compared to other platforms. It was mentioned that the overall rate of SNP Calling was slightly higher in ion torrent than Illumina, $\tilde{8}2\%$ of SNPs being called correctly. But far less false positives ion torrent from Illumina for SNPs and InDels. They noted that the inbuilt automatic variant calling in ion torrent calling only 1.4% of variants and to detect small InDels using proton platform we need to develop more accurate variant calling technique [4]. From torrent variant caller [5], using hotspot and target file, TVC evaluates the variants in specific locations and regions in order to achieve high accuracy and increase the true

positive variants. However, using hotspot file takes days to process and generate variants report. Although, TVC takes minutes to generate the variants without using hotspot file.

1.2 Thesis Objective

The main purpose of the thesis is to optimize the processing time of Torrent Variant Caller (TVC) to call variants when using hotspot file. We achieved this by using the NOCALL variants and reuse it as hotspot positions to TVC. Additionally, we present solution to automate the generation of new hotspot file and enhance the run time of torrent variant caller to detect variants.

1.3 Thesis Organization

In chapter 2, we will introduce briefly the biological terms related to the problem and DNA sequencing generations and technologies. As well as, we defined an overview of NGS. A detailed review of variant calling pipelines presented in chapter 3. Chapter 4 review parallel computing implementation types, advantages, and the contribution also elaborate the steps and tools used to achieve. The results are discussed in chapter 5. Then finally, chapter 6 concludes and discusses our future work. Appendix A and B summarize the common steps to run variant calling pipelines. Appendix C shows the steps to run the developed package.

Chapter 2

Biological Background

In this chapter, we introduce the basic concepts of biology and bioinformatics that are needed in the thesis. We will show the main topics of bioinformatics such as sequencing technologies, biological sequence databases and sequences format.

2.1 Genetic Background

In 1857, Gregor Mendel performed an experiment with plants that increased the interest in the study of genetics. After that, Friedrich Miescher in 1869 discovered a substance he called "nuclein" or "nucleic acid" that exists only in the chromosome later he isolated sample of material known now as "DNA". James Watson and Francis Crick began to examine the DNA's structure guided by X-ray diffraction photos of DNA fibers taken by Maurice Wilkins. They created double helix model with little bars called nucleotides connecting the two strands.

Our human body is build of billions of cells. Each cell varies according to the organ function. DNA (deoxyribonucleic acid) considered as the blueprint for our life, which exists in the nucleus of the cell. It is present in all living organisms from the smallest bacterium to the largest whale. DNA determines the genetic characteristics and the behavior of organism including the diseases that may develop.

The DNA structure includes four main nucleotide bases, Adenine (A), Cytosine (C), Guanine (G) and Thymine (T), attached to sugar and phosphate string. Each DNA bases bound to each other using hydrogen bonds to form units called base pairs, A with T and C with G. Each base attached to a sugar and a phosphate molecule as shown in figure 2.1. Human DNA contains around 3 billion bases. These bases sequenced according to the organ and information that needs to be transmitted.

The nucleus contains different number of chromosomes for each organism. The set of chromosomes make up a genome. The human genome arranged into 46 chromosomes. The sequence of pieces of DNA called genes that act as instructions to make molecules called proteins. According to Human Genome Project estimation, humans have between 20,000 and 25,000 genes. Every human has pairs of each gene, inherited from each parent.

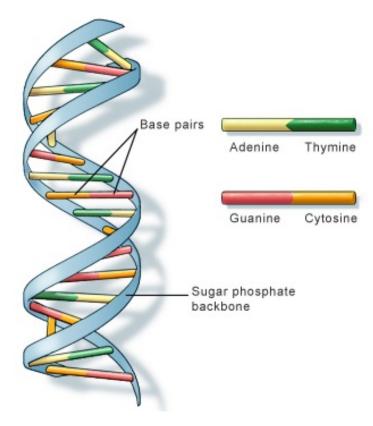


Figure 2.1: DNA Helix strand.

2.2 DNA sequencing

DNA sequencing is a laboratory method to determine the order of the nucleotides 4 bases (Adenine (A),thymine (T), cytosine (C), guanine (G)) in a piece of DNA strand. we cannot sequence a genome from start to end. We have to break the DNA strand into many smaller fragments, then sequencing the pieces and label each base individually with different colors. The detailed flow described in the next chapter. Previously, to sequence one or two genes we may spend several years. DNA sequencing technologies have been developed rapidly since the completion of Human Genome Project. We can sequence the entire genome in few hours and with less money.

DNA sequencing can tell the scientists the genetic information in particular DNA segment.

In addition, sequencing can be used in molecular biology to allow researchers to identify the changes in gene that may cause disease and identify potential drug targets. DNA sequencing can be used to study how different organisms related to each other and how they evolved. Moreover, we can determine the risk of genetic diseases.