

# Approximate string matching for high-throughput sequencing

Dissertation zur Erlangung des Grades eines  
Doktors der Naturwissenschaften (Dr. rer. nat.)  
vorgelegt von

Enrico Siragusa



am Fachbereich Mathematik und Informatik  
der Freien Universität Berlin

Berlin 2015

Gutachter:

***Prof. Dr. Knut Reinert***, Freie Universität Berlin, Deutschland

***Prof. Dr. Raffaele Giancarlo***, Università degli Studi di Palermo, Italien





## Abstract

Over the last years, high-throughput sequencing (HTS) has become an invaluable method of investigation in molecular and medical biology. HTS technologies allow to sequence cheaply and rapidly an individual's DNA sample under the form of billions of short DNA reads. The ability to assess the content of a DNA sample at base-level resolution opens to a myriad of applications, including individual genotyping and assessment of large structural variations, measurement of gene expression levels and characterization of epigenetic features. Nonetheless, the quantity and quality of data produced by HTS instruments call for computationally efficient and accurate analysis methods.

In this thesis, I present novel methods for the mapping of high-throughput sequencing DNA reads, based on state of the art approximate string matching algorithms and data structures. Read mapping is a fundamental step of any HTS data analysis pipeline in resequencing projects, where DNA reads are reassembled by aligning them back to a previously known reference genome. The ingenuity of approximate string matching methods is crucial to design efficient and accurate read mapping tools.

In the first part of this thesis, I cover practical indexed and filtering methods for exact and approximate string matching. I present state of the art algorithms and data structures, give their pseudocode and discuss their implementation. Furthermore, I provide all implementations within *SeqAn*, the generic C++ template library for sequence analysis, which is freely available under <http://www.seqan.de/>. Subsequently, I experimentally evaluate all implemented methods, with the aim of guiding the engineering of new sequence alignment software. To the best of my knowledge, this is the first work providing a comprehensive exposition, implementation and evaluation of such methods.

In the second part of this thesis, I turn to the engineering and evaluation of read mapping tools. First, I present a novel method to find all mapping locations per read within a user-defined error rate; this method is published in the peer-reviewed journal *Nucleic Acids Research* and packaged in an open source tool nicknamed *Masai*. Afterwards, I generalize this method to quickly report all co-optimal or suboptimal mapping locations per read within a user-defined error rate; this method, packaged in a tool called *Yara*, provides a more practical, yet sound solution to the read mapping problem. Extensive evaluations, both on simulated and real datasets, show that *Yara* has better speed and accuracy than *de-facto* standard read mapping tools.

## Zusammenfassung

In den letzten Jahren wurde die Hochdurchsatz-Sequenzierung (HTS) zu einem unverzichtbaren Bestandteil der molekularmedizinischen Forschung. HTS Technologien ermöglichen es DNS-Proben von Individuen schnell und günstig in Form von Milliarden kurzer DNS-Reads zu sequenzieren. Die Fähigkeit, die Basenpaarabfolge einer DNS-Probe zu bestimmen, eröffnet viele neue Anwendungsgebiete, wie zum Beispiel die Genotypisierung, die Beurteilung von strukturellen Variationen, die Messung der Genexpressionslevel oder die Etablierung von epigenetischen Faktoren. Jedoch setzen sowohl Quantität als auch die Qualität der von HTS-Technologien produzierten Daten recheneffiziente und akkurate Analysemethoden voraus um als Standardverfahren im biomedizinischen Bereich eingesetzt werden zu können.

In dieser Arbeit stelle ich neue Methoden für das sogenannte Read Mapping von HTS Daten vor. Read Mapping ist ein essentielles Verfahren, bei dem aus den Produkten einer HTS Applikation mit Hilfe eines bereits bekannten Referenzgenoms die ursprüngliche DNS-Reads assembliert wird. Die Entwicklung von speziellen und neuartigen Algorithmen für die approximative Stringsuche spielt dabei eine entscheidende Rolle, um akkurate und effiziente Read Mapping Programme zu entwickeln.

Im ersten Teil dieser Arbeit beschreibe ich praktische Index-Datenstrukturen sowie Filtermethoden, die bei der approximativen Stringsuche angewendet werden. Ich stelle die Algorithmen im Pseudocode dar und bespreche deren Funktionsweise und Implementierung im Detail. Sämtliche Algorithmen und Datenstrukturen, die ich in dieser Arbeit vorstelle, wurden in *SeqAn* - einer generischen C++ Template-Bibliothek für Sequenzanalyse - implementiert und sind darüber verfügbar (siehe <http://www.seqan.de/>). Anschließend analysiere und vergleiche ich die vorgestellten Methoden in verschiedenen Experimenten ausführlich miteinander, mit dem Ziel neue verbesserte Alignentalgorithmen entwerfen zu können. Nach meinem besten Wissen und Gewissen, ist dies die erste Arbeit, die die genannten Methoden umfassend in ihrer Implementierung und Funktionsweise diskutiert und bewertet.

Im zweiten Teil dieser Arbeit beschreibe und diskutiere ich zwei neue Read Mapping Programme und vergleiche diese mit bisherigen Anwendungen. Dabei stelle ich zuerst eine neue Methode vor, welche alle potentiellen genomischen Ursprungspositionen für jeden DNS-Read mit einer spezifizierten Fehleranzahl lokalisiert. Die beschriebene Methode wurde im Open Source Programm *Masai* implementiert und wurde in der Zeitschrift *Nucleic Acids Research* publiziert. Anschließend generalisiere ich die Methode und zeige wie alle co-optimalen oder suboptimalen Positionen, gegeben einer vom Nutzer definierten Fehlerrate, effizient gefunden werden können. Das Open Source Programm *Yara* implementiert dieses Verfahren und bietet somit eine wesentlich praktischere und noch dazu solide Lösung für das Read Mapping Problem. Eine umfassende Analyse auf simulierten und realen Daten ergab, dass Yara schneller und genauer als De-facto-Standard Read Mapping Programme ist.

## Acknowledgments

The warmest thanks for this work go to my advisor, Knut Reinert. He constantly showed his interest, motivated and guided me gently during these years. Now I consider myself lucky and honored to have been working with him.

Special thanks go to Raffaele Giancarlo, first for having introduced me to the field of Bioinformatics, then for having recommended me to the International Max Planck Research School (IMPRS), and finally for having accepted to review this thesis. I am also grateful to my IMPRS advisors Martin Vingron and Peter N. Robinson, as well as to the IMPRS coordinator Kirsten Kelleher — this work has been possible thanks to the financial support of the IMPRS.

Huge thanks go to David Weese for constantly collaborating and teaching me his art of computer programming — without his help I would not have gone far. My gratitude goes also to Jochen Singer for jointly working on the FM-index and to Manuel Holtgrewe for working together on the evaluations of RazerS 3 and Masai. And thanks to Jacopo Pantaleoni, Jonathan Cohen and Mark Berger for the great experience at NVIDIA!

Thanks — for too many reasons — to the IMPRSers: Akdes, Alessandro, Anne-Katrin, Arthur, Birte, Chris, Leon, Matt, Mike, Sabrina, Sandro, Stephan, Temesgen, Victor, Xiao, and to the people at the Free University: Anja, Christian, Edzard, Hannes, Julianus, Jialu, Katrin, René.

I finally send a warm hug for their support to my parents Anna and Giovanni, my bros Alessandro and Daniele, and Vija, the sweetest thing in my life.





# CONTENTS

<b>1. Introduction</b>	1
1.1 High-throughput sequencing	1
1.1.1 Applications	2
1.2 Outline	2
1.2.1 Approximate string matching	3
1.2.2 Read mapping	3
 <b>Part I Approximate string matching</b>	 5
<b>2. Preliminaries</b>	7
2.1 Definitions	7
2.2 Transcripts, alignments and distances	8
2.3 Edit distance computation	9
2.4 String matching	11
2.4.1 Online methods	12
2.4.2 Indexed methods	13
2.4.3 Filtering methods	15
 <b>3. Indexed methods</b>	 19
3.1 Classic full-text indices	20
3.1.1 Suffix array	20
3.1.2 Suffix tree realizations	23
3.1.3 $q$ -Gram index	24
3.1.4 Trie and radix tree realizations	26
3.2 Succinct full-text indices	26
3.2.1 Burrows-Wheeler transform	27
3.2.2 Rank dictionaries	29
3.2.3 FM-index	33
3.3 Algorithms	35
3.3.1 Construction	35
3.3.2 Depth-first traversal bounded by depth	36
3.3.3 Exact string matching	37
3.3.4 Backtracking $k$ -mismatches	38
3.3.5 Multiple exact string matching	40
3.3.6 Multiple $k$ -mismatches	41

<b>4. Filtering methods</b>	45
4.1 Exact seeds	46
4.1.1 Principle	46
4.1.2 Efficiency	46
4.2 Approximate seeds	47
4.2.1 Principle	47
4.2.2 Filtration schemes	48
4.3 Contiguous $q$ -grams	49
4.3.1 Principle	49
4.3.2 Filtration schemes	49
4.3.3 Bucketing	50
4.4 Gapped $q$ -grams	50
4.4.1 Principle	51
4.4.2 Filtration schemes	52
4.4.3 Full sensitivity	53
4.4.4 Optimal threshold	54
4.4.5 Specificity	55
4.4.6 Families	59
4.5 Evaluation	60
4.5.1 Runtime	60
4.5.2 Verification versus filtration time	62
4.5.3 Positive predictive value	62
 <b>Part II Read mapping</b>	 65
<b>5. Background</b>	67
5.1 High-throughput sequencing data	67
5.1.1 Read sequences	67
5.1.2 Phred base quality scores	68
5.2 High-throughput sequencing data analysis	69
5.2.1 Data analysis pipelines	69
5.2.2 Secondary analysis paradigms	70
5.2.3 Best-mapping	72
5.2.4 All-mapping	73
5.3 Limits of high-throughput sequencing	74
5.3.1 Genome mappability	74
5.3.2 Genome mappability score	75
5.4 Popular read mappers	76
5.4.1 Bowtie and Bowtie 2	77
5.4.2 BWA	78
5.4.3 Soap	78
5.4.4 SHRiMP 2	78
5.4.5 RazerS and RazerS 3	78
5.4.6 mrFast and mrsFast	79
5.4.7 Hobbess 2	79
5.4.8 GEM	79

<b>6. Masai</b>	81
6.1 Engineering	81
6.1.1 Filtration	82
6.1.2 Indexing	83
6.1.3 Verification	85
6.2 Evaluation	85
6.2.1 Rabema benchmark on simulated data	85
6.2.2 Variant detection on simulated data	87
6.2.3 Performance on real data	88
6.2.4 Performance with different indices	89
6.2.5 Filtration efficiency	89
6.3 Discussion	91
<b>7. Yara</b>	93
7.1 Engineering	93
7.1.1 Stratified mapping	93
7.1.2 Adaptive filtration	95
7.1.3 Indexing	95
7.1.4 Paired-end and mate-pair protocols	96
7.1.5 Mapping qualities	96
7.2 Evaluation	98
7.2.1 Experimental setup	98
7.2.2 Rabema benchmark on real data	99
7.2.3 Accuracy on simulated data	100
7.2.4 Variant calling on real data	102
7.3 Discussion	104
<b>A. Read mappers parameterization</b>	107
A.1 Masai evaluation	107
A.2 Yara evaluation	108
<b>B. Curriculum Vitæ</b>	109
<b>C. Declaration</b>	113
<b>Bibliography</b>	115



The sequencing of the whole human genome has been one of the major scientific achievements of the last decades. In February 2001, the three billion dollars publicly funded *Human Genome Project* (HGP) published a first draft *covering more than 96 % of the euchromatic part of the human genome* [Consortium, 2001]. Concurrently, the privately funded company *Celera Genomics* published a *2.91 billion base pair consensus sequence of the euchromatic portion of the human genome* [Venter *et al.*, 2001]. The application of efficient computational methods has been crucial for the accomplishment of these projects.

These sequencing projects used the *Sanger method* [Sanger *et al.*, 1977] with *capillary electrophoresis*, a technology producing high fidelity DNA reads long 700 bp in average, at a throughput of about 150 *kilo base pairs per hour* (Kbp/h). Because of such technological limitation, the sequencing of a whole genome had to be coupled with the *shotgun* approach, which consists of chopping long DNA fragments up into smaller segments and then generating *reads* of these short nucleotide sequences. The whole human genome finally had to be reassembled from its sequencing reads using computational methods.

## 1.1 High-throughput sequencing

The success of these projects did not mark the end of the sequencing era, but rather its beginning. Since then, sequencing technology steadily improved and evolved into what is now called *high-throughput sequencing* (HTS) or *next-generation sequencing* (NGS). In 2004, *454 Life Science* commercialized the *Genome Sequencer FLX*, an instrument based on large-scale parallel *pyrosequencing*, capable of sequencing DNA in form of 400 bp reads at a throughput of about 20 Mbp/h. High-throughput sequencing was born.

In 2006, Solexa released its *1G Genetic Analyzer*, based on a massively parallel technique of *reversible terminator-based sequencing*. The instrument produced reads as short as 30 bp with lower accuracy than Sanger sequencing but at very high-throughput: it allowed *resequencing* a whole human genome in three months for about \$100,000. Following this success, Solexa was acquired by *Illumina*, which is nowadays the market leader. At the beginning of 2014, Illumina announced the *HiSeq X Ten*, allowing in less than three days the sequencing of many whole human genomes at \$1,000 each.

### 1.1.1 Applications

In the last few years, HTS has become an invaluable method of investigation for computational molecular biologists. Abundant and cost-effective production of sequencing data permits viewing not only genomic DNA but also transcripts and epigenetic features at single-base resolution.

*Whole genome sequencing* (WGS) allows discovery of genetic variations across the whole genome. These variations may be in the form of single nucleotide variants (SNVs), small insertions or deletions (INDELs), or large structural variants such as transversions, trans-locations, and copy number variants (CNVs).

*Whole exome sequencing* (WES) is a cost-effective, yet powerful alternative to WGS. This protocol consists in the *targeted sequencing* of the *exome*, i.e. the protein coding subset of a genome. WES has recently begun to be used for clinical diagnostics, e.g. for the interpretation of tumor samples or for the characterization of mendelian disorders.

*RNA sequencing* (RNA-seq) is a protocol to sequence the *transcriptome*, i.e. the set of RNA molecules of an organism, including mRNAs, rRNAs, tRNAs and other non-coding RNAs. Actual HTS technologies perform deep-sequencing of RNA molecules after they have been reverse-transcribed into cDNA. RNA-seq provides a myriad of information on gene expression, alternative splicing, intron/exon boundaries, untranslated regions (UTRs), and genetic variation with single-base accuracy.

*Chromatin immunoprecipitation* with next-generation sequencing (ChIP-seq) is a protocol that allows the selective sequencing of DNA bound by a specific protein. The process isolates chromatin-protein complexes and sheares them by sonication; all DNA fragments bound by the specific protein are subsequently pulled down using immunoprecipitation with a protein-specific antibody and finally sequenced using HTS. With ChIP-seq it is possible to determine global methylation patterns, identify transcription factor binding sites, histone modifications, and chromatin remodeling proteins.

## 1.2 Outline

This thesis presents novel methods for the *efficient* and *accurate* mapping of high-throughput sequencing DNA reads, based on state of the art *approximate string matching* algorithms and data structures. Read mapping is a non-trivial, ubiquitous task in all resequencing applications. Efficiency is mandatory to keep the pace of sequencing technologies, exponentially increasing in throughput. Accuracy is required to enable downstream data analysis at single-base resolution. The ingenuity of state of the art approximate string matching methods is crucial for the design and implementation of efficient and accurate read mapping programs.

The contributions of this work are of purely practical interest, nonetheless I follow a rigorous approach. I subdivide this work in two parts. Part I covers practical approximate string matching methods, whose interest is beyond HTS applications. Part II describes the application of such methods to engineer two HTS read mapping programs developed by myself. In the two following sections, I describe the contents of part I and II.

### 1.2.1 Approximate string matching

In chapter 2, I introduce basic stringology concepts. I give an overview of basic online, indexed and filtering methods for exact and approximate string matching. In particular, in section 2.4.2, I introduce the concept of full-text index and define a set of *generic* top-down traversal operations.

In chapter 3, I cover *indexed methods* for exact and approximate string matching. First, I describe some classic full-text indices (suffix arrays and  $q$ -gram indices) and succinct full-text indices (uncompressed variants of the FM-index). Afterwards, I introduce generic string matching algorithms, working on any of these data structures, and provide their experimental evaluation. My implementation of all these algorithms and data structures is publicly available in source form within the C++ library SeqAn [Döring *et al.*, 2008]. This work is being published as a book chapter [Weese and Siragusa, 2015].

In chapter 4, I cover *filtering methods* for approximate string matching. I consider two classes of full-sensitive filtering methods: those based on *seeds* and those based on  $q$ -grams. From the former class, I cover *exact seeds* [Baeza-Yates and Perleberg, 1992] and *approximate seeds* [Myers, 1994; Navarro and Baeza-Yates, 2000]. From the latter one *contiguous  $q$ -grams* [Jokinen and Ukkonen, 1991], *gapped  $q$ -grams* [Burkhardt and Kärkkäinen, 2001] and *multiple gapped  $q$ -grams* (also called  *$q$ -gram families*) [Kucherov *et al.*, 2005]. To the best of my knowledge, this is the first work providing a comprehensive exposition of these methods together with their experimental evaluation. In addition, I introduce a formal framework for (multiple) gapped  $q$ -grams, which leads to the formulation of efficient algorithms answering some combinatorially hard filter design questions.

### 1.2.2 Read mapping

In chapter 5, I give background information on HTS read mapping. I start by giving a quick overview of market-leading HTS technologies. Subsequently, I introduce two *de facto* standard paradigms for HTS data analysis: *best-mapping* and *all-mapping*. By reviewing some recent works [Derrien *et al.*, 2012; Lee and Schatz, 2012], I try to delineate which are the limits of HTS. Finally, I give an overview of the most popular read mappers. In the last two chapters, I consider these popular programs in the evaluation of my own tools, Masai [Siragusa *et al.*, 2013a] and Yara.

In chapter 6, I present the engineering and evaluation of a read mapping tool for the all-mapping paradigm. My method is packaged in a C++ tool nicknamed *Masai*, which stands for *multiple backtracking of approximate seeds on a suffix array index*. Masai is part of the SeqAn library [Döring *et al.*, 2008], it is distributed under the BSD license and can be downloaded from <http://www.seqan.de/projects/masai>. The result of this work has been published in the peer-reviewed journal *Nucleic Acids Research* [Siragusa *et al.*, 2013a].

In chapter 7, I present *Yara*, a non-heuristic best-mapper capable of quickly reporting all co-optimal or suboptimal mapping locations within a given error rate. The tool works with Illumina or Ion Torrent reads, supports paired-end and mate-pair protocols,

computes accurate mapping qualities, offers parallelization via multi-threading, has a low memory footprint thanks to the FM-index, and does not require ad-hoc parameterization. Yara is part of the SeqAn library [Döring *et al.*, 2008], it is distributed under the BSD license and can be downloaded from <http://www.seqan.de/projects/yara>. A manuscript for this work is in preparation and will be submitted to a peer-reviewed journal.



## **Part I**

### **APPROXIMATE STRING MATCHING**



In this chapter, I introduce fundamental definitions and problems of stringology. The reader familiar with basic stringology can skip this chapter and proceed to chapter 3.

## 2.1 Definitions

An *alphabet*  $\Sigma = \{a_0, \dots, a_{\sigma-1}\}$  is a finite ordered set of characters (or symbols)  $a_i$ . A *string* (or word)  $S = s_0 \dots s_{n-1}$  over  $\Sigma$  is a finite sequence of symbols  $s_i \in \Sigma$ . I denote with  $S_{i\dots j}$  the sequence of symbols  $s_i \dots s_{j-1}$  (excluding  $j$ ) for any  $0 \leq i < j \leq n$ , where  $n = |S|$  is the length of  $S$ . The set  $\Sigma^0 = \{\epsilon\}$  contains only the *empty string* s.t.  $|\epsilon| = 0$ ,  $\Sigma^n$  contains all strings of length  $n$  over  $\Sigma$ , and  $\Sigma^* = \cup_{n=0}^{\infty} \Sigma^n$  all strings over  $\Sigma$ .

In this manuscript, I use numerical intervals, arrays and matrices. Numerical intervals are sequences over  $\mathbb{N}_0$ , i.e.  $[i, j] = \{i, i+1, \dots, j\}$  and  $[i, j) = [i, j-1]$ . An *array* (or table, or vector)  $A = a_0 \dots a_{n-1}$  is a finite sequence of numbers  $a_i \in \mathbb{N}_0$ . I denote by  $A[i]$  the subscript element  $a_i$ . A *matrix*  $X$  is a rectangular array of numbers  $x_{i,j} \in \mathbb{N}_0$  arranged s.t.

$$X = \begin{pmatrix} x_{0,0} & \cdots & x_{0,n-1} \\ \vdots & \ddots & \vdots \\ x_{m-1,0} & \cdots & x_{m-1,n-1} \end{pmatrix}$$

I denote by  $X[i, j]$  the subscript element  $x_{i,j}$ .

**Definition 2.1 (String collection).** A *string collection* is an ordered multiset  $\mathbb{S} = \{S^0, \dots, S^{c-1}\}$  of non necessarily distinct strings over a common alphabet  $\Sigma$ . I denote by  $\|\mathbb{S}\| = \sum_{i=0}^{c-1} |S^i|$  the total length of the string collection. The above definitions of prefix, suffix and substring generalize to multisets, e.g.  $\mathbb{S}_{(d,i)\dots(d,j)}$  denotes the substring  $S^d_{i\dots j}$ .

**Definition 2.2 (String iterator).** Given a string  $S$ , I denote with  $s$  its iterator starting at position  $i = 0$  and providing the following operations

- $\text{GONEXT}(s)$  moves to position  $i + 1$ ;
- $\text{GOPREVIOUS}(s)$  moves to position  $i - 1$ ;
- $\text{ATEND}(s)$  returns *true* iff  $i = |S|$ ;
- $\text{VALUE}(s)$  returns  $s_i$ .

**Definition 2.3 (Substring).** A string  $X$  is a *prefix* of another string  $Y$  iff  $X = Y_{0\dots i}$  for some  $0 < i \leq |Y|$ , a *suffix* iff  $X = Y_{0\dots i}$  for some  $0 < i \leq |Y|$ , and a *substring* iff  $X = Y_{i\dots j}$  for

some  $0 < i < j \leq |Y|$ . For convenience, I denote the suffix of  $Y$  beginning at position  $i$  simply by  $Y_i$ . Moreover, I denote the concatenation of two strings  $X$  and  $Y$  by  $X \cdot Y$  or simply  $XY$ .

**Definition 2.4 (Occurrence).** A string  $X$  of length  $n$  occurs in a longer string  $Y$  iff  $X$  is a substring of  $Y$ . Given  $X = Y_{i \dots i+n}$ , the occurrence of  $X$  starts at position  $i$  and ends at position  $i + n - 1$  in  $Y$ .

**Definition 2.5.** Given an alphabet  $\Sigma$ , I define the function  $\delta : \Sigma \times \Sigma \rightarrow \{0, 1\}$  s.t.  $\delta(a, b) = 1$  for any distinct  $a, b \in \Sigma$  and 0 otherwise.

**Definition 2.6 (Lexicographic rank).** Given an alphabet  $\Sigma$  of size  $\sigma$ , I denote the *lexicographic rank* of any alphabet symbol by the function  $\text{rank} : \Sigma \rightarrow [0, \sigma)$ , s.t.  $\text{rank}(a) < \text{rank}(b) \Leftrightarrow a < b$  for any distinct  $a, b \in \Sigma$ . In addition, I denote by  $\min_{\text{lex}} \Sigma$  the lexicographically smallest symbol  $a \in \Sigma$  s.t.  $\text{rank}(a) = 0$ , and by  $\text{next}_{\text{lex}} a$  the lexicographical successor  $b$  of  $a$  in  $\Sigma$  s.t.  $\text{rank}(b) = \text{rank}(a) + 1$ .

**Definition 2.7 (Lexicographical order).** The *lexicographical order*  $<_{\text{lex}}$  between two non-empty strings  $X, Y$  is defined as  $X <_{\text{lex}} Y \Leftrightarrow x_0 < y_0$ , or  $x_0 = y_0$  and  $X_1 <_{\text{lex}} Y_1$ .

To well-define full-text indices, I use the definitions of padded string (2.8) and padded string collection (2.9). To this end, I introduce special *terminator* symbols  $\$ \notin \Sigma$ .

**Definition 2.8 (Padded string).** Given a string  $S$  over  $\Sigma$ , I call *padded string* the string  $S\$$  over  $\Sigma_\$$ , where  $\Sigma_\$ = \Sigma \cup \{\$\}$  and  $\text{rank}(\$) < \text{rank}(a)$  for any  $a \in \Sigma$ .

**Definition 2.9 (Padded string collection).** Given a string collection  $\mathcal{S} = \{S^0, \dots, S^{c-1}\}$  over  $\Sigma$ , I call *padded string collection* the collection  $\{S^0\$, \dots, S^{c-1}\$^{c-1}\}$  over  $\Sigma_\$$ , where  $\Sigma_\$ = \Sigma \cup \{\$, \dots, \$^{c-1}\}$  and  $\text{rank}(\$^i) < \text{rank}(\$^j) \Leftrightarrow i < j$ .

**Definition 2.10 (Reverse string).** Given a string  $S$ ,  $\bar{S}$  denotes its reversal, i.e.  $\bar{s}_i = s_{n-1-i}$  for any  $0 \leq i < n$ .

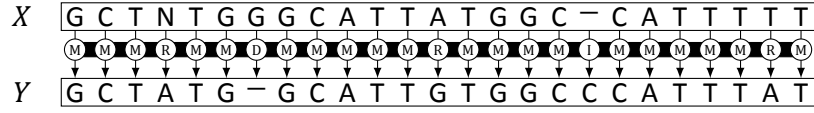
**Definition 2.11 (Reverse string collection).** Given a string collection  $\mathcal{S} = \{S^0, \dots, S^{c-1}\}$ ,  $\bar{\mathcal{S}}$  denotes the reversed collection  $\{\bar{S}^0, \dots, \bar{S}^{c-1}\}$ .

## 2.2 Transcripts, alignments and distances

I now define the basic edit operations to transform one string into another. Given two strings  $X, Y$  of equal length  $n$ , the string  $X$  is easily transformed into the string  $Y$  by substituting (or replacing) all symbols  $x_i$  s.t.  $x_i \neq y_i$  into  $y_i$ , for  $0 \leq i < n$ . However, if the two strings have different lengths, some symbols must be necessarily inserted or deleted from  $X$  in order to obtain  $Y$ . This fact motivates the following definition of *edit transcript*.

**Definition 2.12 (Edit transcript).** An edit transcript for any two given strings  $X, Y$  is a finite sequence of substitutions, insertions and deletions transforming  $X$  into  $Y$  [Gusfield, 1997].

**Figure 2.1:** Example of edit transcript and alignment. The string  $X$  is transformed into  $Y$ . The transcript character M indicates a match, R a replacement, I an insertion, and D a deletion.



An *alignment* is an alternative way of visualizing a transformation between strings. While an edit transcript is an explicit sequence of edit operations that transform one string into another, an alignment is an explicit relationship between pairs of symbols from the two strings. Nonetheless, when some symbols are inserted or removed, some symbols in one string are not related to any symbol in the other string. For this reason, it is necessary to introduce an additional gap symbol —, not being part of the string alphabet  $\Sigma$ . The definition of alignment follows. Figure 2.1 shows an example of edit transcript with its associated alignment.

**Definition 2.13 (Alignment).** An alignment of two strings of length  $m, n$  over  $\Sigma$  is a string of length between  $\min\{m, n\}$  and  $m + n$  over the pair alphabet  $(\Sigma \cup \{-\}) \times (\Sigma \cup \{-\}) \setminus \{(-, -)\}$ .

Hamming and edit distance are two fundamental distance functions on strings.

**Definition 2.14 (Hamming distance).** The *Hamming distance*  $d_H : \Sigma^n \times \Sigma^n \rightarrow \mathbb{N}_0$  between two strings  $X, Y \in \Sigma^n$  counts the number of substitutions necessary to transform  $X$  into  $Y$  [Hamming, 1950].

**Definition 2.15 (Edit distance).** The *Levenshtein* or *edit distance*  $d_E : \Sigma^* \times \Sigma^* \rightarrow \mathbb{N}_0$  between two strings  $X, Y \in \Sigma^*$  counts the *minimum* number of edit operations necessary to transform  $X$  into  $Y$  [Levenshtein, 1966].

The problem of finding an optimal alignment between two strings is equivalent to the problem of finding their minimum distance [Gusfield, 1997]. While the Hamming distance between any two strings of length  $n$  is easily computed in time  $\mathcal{O}(n)$ , computing the edit distance involves solving a non-trivial optimization problem.

## 2.3 Edit distance computation

The edit distance between two strings is efficiently computed via *dynamic programming* (DP) [Needleman and Wunsch, 1970]. Let  $X, Y$  be two strings of length  $n \geq m$ . The edit distance  $d_E(X_{0\dots i}, Y_{0\dots j})$  between any their prefixes  $X_{0\dots i}$  and  $Y_{0\dots j}$  is defined recursively. The base conditions of the recurrence relation are:

$$d_E(\epsilon, \epsilon) = 0 \quad (2.1)$$

$$d_E(X_{0\dots i}, \epsilon) = i \text{ for all } 1 < i \leq n \quad (2.2)$$

$$d_E(\epsilon, Y_{0\dots j}) = j \text{ for all } 1 < j \leq m \quad (2.3)$$

and the recursive case for all  $1 < i \leq n$  and  $1 < j \leq m$  is as follows:

$$d_E(X_{0\dots i}, Y_{0\dots j}) = \min \begin{cases} d_E(X_{0\dots i-1}, Y_{0\dots j}) & + 1 \\ d_E(X_{0\dots i}, Y_{0\dots j-1}) & + 1 \\ d_E(X_{0\dots i-1}, Y_{0\dots j-1}) & + \delta(x_i, y_j) \end{cases} \quad (2.4)$$

where the function  $\delta$  indicates whether the characters  $x_i, y_j$  match or mismatch (see definition 2.5).

Algorithm 2.1 computes the above recurrence relation in time  $\mathcal{O}(nm)$  using a dynamic programming table  $D$  of  $(n+1) \times (m+1)$  cells, where cell  $D[i, j]$  stores the value of  $d_E(X_{0\dots i}, Y_{0\dots j})$ . The sole edit distance without any alignment can be computed in space  $\mathcal{O}(m)$ ; indeed, column  $j-1$  alone is sufficient to compute column  $j$ . An optimal alignment can be computed in time  $\mathcal{O}(n+m)$  via *traceback* on the table  $D$ ; the traceback starts in the cell  $D[m, n]$  and goes backwards (either left, up-left, or up) to the previous cell by deciding which condition of equation 2.4 yielded the value of  $D[m, n]$ .

Further, if no alignment at all is needed, there is a more efficient solution to the edit distance problem. Myers' algorithm [Myers, 1999], instead of computing DP cells one after another, encodes one whole DP column in two bit-vectors and computes the adjacent column in a constant number of 12 logical and 3 arithmetical operations. The time complexity of Myers' algorithm is  $\mathcal{O}(nm/b)$ , where  $b$  is the processor's bit-width, e.g. 64 on modern processors.

Finally, if the edit distance problem is relaxed to simply checking if  $d_E(X, Y) \leq k$ , with  $k \in \mathbb{N}_0$ , an answer is given in time  $\mathcal{O}(nk)$ . Indeed, the DP matrix  $D$  needs not to be filled entirely, as the cells outside a *band* of width  $k+1$  around the 0-th diagonal are not *coaccessible* from the *sink* cell  $D[m, n]$ . The reader is referred to [Siragusa *et al.*, 2013b]

---

**Algorithm 2.1** EDITDISTANCE( $X, Y$ )

---

**Input**     $X$  : string of length  $n$

$Y$  : string of length  $m$

**Output**   the edit distance  $d_E(X, Y)$  between  $X$  and  $Y$

1:  $D[0, 0] \leftarrow 0$

2: **for**  $j \leftarrow 1$  **to**  $m$  **do**

3:      $D[j, 0] \leftarrow D[j-1, 0] + 1$

4: **for**  $i \leftarrow 1$  **to**  $n$  **do**

5:      $D[0, i] \leftarrow D[0, i-1] + 1$

6:     **for**  $j \leftarrow 1$  **to**  $m$  **do**

7:          $D[i, j] \leftarrow \min \begin{cases} D[i-1, j] & + 1, \\ D[i, j-1] & + 1, \\ D[i-1, j-1] & + \delta(x_i, y_j) \end{cases}$

8: **return**  $D[m, n]$

---

for a proof and to [Weese *et al.*, 2012] for a banded version of Myers' algorithm running in time  $\mathcal{O}(nk/b)$ .

## 2.4 String matching

**Definition 2.16 (Exact string matching).** Given a string  $P$  of length  $m$ , called the *pattern*, and a string  $T$  of length  $n$ , called the *text*, the *exact string matching* problem is to find all occurrences of  $P$  into  $T$  [Gusfield, 1997].

Exact string matching is one of the most fundamental problems in stringology. This problem has been extensively studied from the theoretical standpoint and is well solved in practice [Faro and Lecroq, 2013]. Nonetheless, the definition of distance functions between strings lends to a more challenging problem: *approximate string matching*.

**Definition 2.17 (Approximate string matching).** Given a text  $T$ , a pattern  $P$ , and a *distance threshold*  $k \in \mathbb{N}$ , the approximate string matching problem is to find all occurrences of  $P$  into  $T$  within distance  $k$  [Galil and Giancarlo, 1988].

The approximate string matching problem under the Hamming distance is commonly called the *k-mismatches* problem, while under the edit distance is called the *k-differences* problem. Figure 2.2 shows an example of occurrence for *k-differences*. For *k-mismatches* and *k-differences*, it must hold  $k > 0$  as the case  $k = 0$  corresponds to exact string matching, and  $k < m$  as a pattern trivially occurs at any position in the text if all its  $m$  characters are substituted. Frequently, the problem's input respects the condition  $k \ll m \ll n$ .

**Definition 2.18 (Error rate).** Under the edit or Hamming distance, the *error rate* is defined as  $\epsilon = k/m$ , with  $0 < \epsilon < 1$  given the above conditions.

String matching problems are subdivided in two categories, *online* and *offline*, depending on which string, the pattern or the text, is given first. Algorithms for online string matching work by preprocessing the pattern and scanning the text from left to right (or right to left). Algorithms for offline string matching are instead allowed to preprocess the text, hence they build an index of the text beforehand to speed up subsequent searches. In practice, if the text is long, static and searched frequently, offline methods outperform online methods in terms of runtime, provided the necessary amount of memory for text indexing.

It goes without saying that offline string matching algorithms are tightly bound to text indexing data structures. Almost all of these algorithms require a *full-text index*, i.e. a data structure representing all substrings of the text. Very often, such full-text index is realized as the *suffix tree* [Weiner, 1973], a fundamental data structure in stringology. Among its virtues [Apostolico, 1985], the suffix tree natively provides exact string matching in optimal time and approximate string matching via backtracking [Ukkonen, 1993]. Often, the suffix tree finds its use within hybrid *filtering methods* rather than on its own.

Filtering methods first discard uninteresting portions of the text and subsequently verify only narrower areas. These methods work either online or offline. Online filtering





condition 2.2 of the edit distance recurrence relation becomes

$$d(X_{0\dots i}, \epsilon) = 0 \text{ for all } 1 < i \leq n, \quad (2.5)$$

thus algorithm 2.1 (line 5) must be adapted to initialize the top row of the DP matrix  $D$  accordingly. Then, as an occurrence of the pattern can end anywhere in the text, algorithm 2.1 (line 8) must be changed to check any cell in the bottom row for the condition  $D[i, m] \leq k$ .

## 2.4.2 Indexed methods

I now introduce *suffix tries*, idealized data structures to index full texts. I define a set of generic operations to traverse them in a top-down fashion. In chapter 3, I introduce various data structures replacing suffix tries in practical implementations. Thanks to these generic traversal operations, I later formulate generic string matching algorithms that are independent of the practical suffix trie implementations.

Consider a padded string collection  $\mathbb{S}$  (definition 2.9) consisting of  $c$  strings. Note that padding is necessary to ensure that no string  $S^i \in \mathbb{S}$  is a prefix of another string  $S^j \in \mathbb{S}$ .

**Definition 2.19 (Trie).** The *trie*  $RT(\mathbb{S})$  of  $\mathbb{S}$  is a lexicographically ordered tree data structure having one node designated as the root and  $c$  leaves, denoted as  $0, \dots, c - 1$ , where leaf  $i$  points to string  $S^i$ . In addition,

1. the incoming edge of any inner node of  $RT(\mathbb{S})$  is labeled with a symbol in  $\Sigma$ ;
2. the incoming edge of any leaf of  $RT(\mathbb{S})$  is labeled with a terminator symbol;
3. the path from the root to leaf  $i$  spells the string  $S^i\$^i$ .

**Definition 2.20 (Suffix trie).** Given a padded string  $S$  (definition 2.8) of length  $n$ , the suffix trie  $ST(S)$  is the trie of all suffixes of  $S$ .  $ST(S)$  has  $n$  leaves, where leaf  $i$  points to suffix  $S_i$ .

**Definition 2.21 (Generalized suffix trie).** Given a padded string collection  $\mathbb{S}$  (definition 2.9), its generalized suffix tree is the suffix tree of all strings in  $\mathbb{S}$ . Each of the  $|\mathbb{S}|$  leaves is labeled by a pair  $(i, j)$  and points to suffix  $\mathbb{S}_{(i,j)}$ .

Note that the (generalized) suffix trie here defined contains  $\mathcal{O}(n^2)$  nodes, yet I only consider the suffix trie as an idealized data structure to formulate generic algorithms. The optimal data structure to represent in linear space all suffixes of a given string is the *suffix tree* [Morrison, 1968]. However, the suffix tree comes with the restriction that internal nodes are branching, i.e. they must have more than one child, and with the property that edges can be labeled by strings of arbitrary length (see figure 2.3). This latter property slightly complicates the exposition of string matching algorithms but it does not affect their runtime complexity nor their result. For this reason, in the remainder of this manuscript, I consider w.l.o.g. always suffix tries instead of suffix trees.

### Top-down traversal

I define a set of generic operations to traverse tries in a top-down fashion. Note that the following traversal operations do not walk over edges labeled terminator symbols, as these symbols are not necessary in string matching applications. I define the following operations on the node  $v$  pointed during the traversal of  $ST(S)$ :

- $ISROOT(v)$  returns true iff the pointed node is the root;
- $ISLEAF(v)$  returns true iff all outgoing edges are labeled by terminator symbols;
- $PARENT(v)$  returns the parent node of  $v$ ;
- $LABEL(v)$  returns the symbol labeling the edge entering  $v$ ;
- $OCCURRENCES(v)$  returns the list of positions pointed by leaves below  $v$ ;
- $CONCAT(v)$  returns the concatenation of edge labels on the path from the root to  $v$ .

Moreover, I define the following operations moving from pointed node  $v$ , and returning *true* on success and *false* otherwise:

- $GODOWN(v)$  moves to the lexicographically smallest child of  $v$ ;
- $GODOWN(v, c)$  moves to the child of  $v$  whose entering edge is labeled by  $c$ ;
- $GORIGHT(v)$  moves to the lexicographically next child of  $v$ ;
- $GOUP(v)$  moves to the parent node of  $v$ .

Time complexities of the above operations depend on the data structure implementing the trie. Usually  $LABEL$  is  $\mathcal{O}(1)$ , both variants of  $GODOWN$  and  $GORIGHT$  range from  $\mathcal{O}(1)$  to  $\mathcal{O}(\log n)$ ,  $OCCURRENCES$  is linear in the number of occurrences,  $GOUP$  is  $\mathcal{O}(1)$  but with an additional  $\mathcal{O}(n)$  space complexity to stack all parent nodes. Algorithms 2.2 and 2.3 show how to implement respectively  $GODOWN$  and  $GORIGHT$  using  $GOUP$  a symbol and  $GOUP$ , although with a worst-case time complexity of  $\mathcal{O}(\sigma)$ . In chapter 3, I consider various data structures to implement suffix tries, I show how to implement these operations and give their complexities.

**Figure 2.3:** Suffix trie and suffix tree of the string  $S = ANANAS\$$ .



**Figure 2.4:** Generalized suffix trie of the string collection  $\mathcal{S} = \{ \text{ANANAS}\$, \text{CACAO}\$ \}$ .



#### Algorithm 2.2 goDOWN( $v$ )

**Input**  $v$  : iterator of  $ST(\mathcal{S})$   
**Output** boolean indicating success  
1: **if not** ISLEAF( $v$ )  
2:    $c \leftarrow \min_{lex} \Sigma$   
3:   **do**  
4:     **if** goDOWN( $v, c$ )  
5:       **return true**  
6:   **while**  $c \leftarrow \text{next}_{lex} c$   
7: **return false**

#### Algorithm 2.3 goRIGHT( $v$ )

**Input**  $v$  : iterator of  $ST(\mathcal{S})$   
**Output** boolean indicating success  
1: **if not** ISROOT( $v$ )  
2:    $c \leftarrow \text{LABEL}(v)$   
3:    $w \leftarrow \text{PARENT}(v)$   
4:   **while**  $c \leftarrow \text{next}_{lex} c$  **do**  
5:     **if** goDOWN( $w, c$ )  
6:       **return true**  
7: **return false**

### 2.4.3 Filtering methods

Filtering methods work in two stages: the *filtration stage* discards portions of the text unlikely or unable to contain an occurrence of the pattern; subsequently the *verification stage* checks the remaining portions. The filtration stage proceeds online by scanning the text or alternatively offline using an index of the text. The verification stage uses a conventional method, e.g. the online dynamic programming method or some variation of it. The crux of filtering methods is thus to accurately and efficiently classify text portions as containing or not some occurrence of the pattern.

#### Specificity and sensitivity

Any filtering method is thus a binary classification method. Any text location is *true* if it coincides with the beginning of an occurrence of the pattern and *false* if it does not. The

outcome of the classification method is *positive* for text locations filtered in and *negative* for locations filtered out. Therefore, as shown in table 2.1, any text location belongs either to the set of *true positives* (TP), *false positives* (FP), *true negatives* (TN), or *false negatives* (FN). Thus, standard *specificity* and *sensitivity* measure the accuracy of filtering methods. The specificity of a filter  $f$  on a text  $T$  is defined as:

$$\frac{|TN_f(T)|}{|TN_f(T)| + |FP_f(T)|} \quad (2.6)$$

and the sensitivity as:

$$\frac{|TP_f(T)|}{|TP_f(T)| + |FN_f(T)|} \quad (2.7)$$

**Definition 2.22.** A filter is *lossless* or *full-sensitive* if its sensitivity is 1, i.e. it produces no false negatives, otherwise it is *lossy*.

Many practical applications, e.g. read mapping, do not require strictly lossless filtration. Nonetheless, a predictable or controlled lossy filter helps to interpret the results and insure their quality. For this reason, I focus on criteria yielding filters which are lossless or lossy in a predictable fashion.

### Efficiency

The total runtime of a filtering method is given by the sum of the runtimes of its filtration and verification stages. Filtration specificity determines how much time is spent in the verification stage. Clearly, to keep verification time small, the number of false positives must be low. Nonetheless, the filtration stage must also run in a reasonable amount of time. Hence, the runtime of an efficient filtering method is usually balanced between filtration and verification time. In any way, no filtering method can be efficient when the number of true positive locations approaches the text length.

Filtering methods for string matching work under the assumption that patterns occur in the text with a *low average probability*. For  $k$ -differences, the occurrence probability is a function of the error rate  $\epsilon$  and the alphabet size  $\sigma$ , and can be computed or estimated under the assumption of the text being generated by a specific random source. Under the uniform Bernoulli model, where each symbol of  $\Sigma$  occurs with probability  $1/\sigma$ , Navarro [2001] experimentally finds that  $\epsilon < 1 - 1/\sqrt{\sigma}$  is a tight upper bound on the error rate which ensures few occurrences and for which filtering algorithms are effective. For higher error rates, non-filtering methods, either online or offline, work better.

**Table 2.1:** Classification of text locations by filtering methods.

Occurrence	Filtered in	Filtered out
Yes	True positive (TP)	False negative (FN)
No	False positive (FP)	True negative (TN)

---

According to the above considerations, filtration efficiency is bound to 50 % error rate on random DNA. However, practical HTS read mapping applications require much lower error rates, usually ranging from 3 to 10 %. For this reason, filtering methods are expected to be very efficient in HTS applications.

### **Seeds versus $q$ -grams**

Filtering methods apply combinatorial criteria to determine which portions of the text might contain some occurrence of the pattern. These criteria are in general valid for both online and offline variants of the problem. In practice, one specific criterion might be more convenient for one variant of the problem rather than the other. The combinatorial criterion underlying a filter is of paramount importance as it provides guarantees on filtration sensitivity.

In chapter 4, I consider two classes of combinatorial filtering methods: those based on *seeds* and those based on  $q$ -grams. Filters in the former class partition the pattern into *non-overlapping* factors called seeds; application of the pigeonhole principle yields full-sensitive partitioning strategies. Instead, filters in the latter class consider all *overlapping* substrings of the pattern having length  $q$ , the so-called  $q$ -grams; simple counting lemmata give lower bounds on the number of  $q$ -grams that must be present in a narrow window of the text, as necessary conditions for an approximate occurrence of the pattern.



Suffix trees are elegant data structures but they are rarely used in practice. Although suffix trees provides theoretically optimal construction and query time, their high space consumption prohibits practical indexing of large string collections. A practical study on suffix trees by Kurtz [1999] reports that efficient implementations achieve sizes between  $12n$  and  $20n$  bytes per character. For instance, two years before completing the sequencing of the human genome, Kurtz conjectured the resources required for computing the suffix tree for the complete human genome (consisting of about  $3 \cdot 10^9$  bp) in 45.31 GB of memory and nine hours of CPU time, and concluded that *“it seems feasible to compute the suffix tree for the entire human genome on some computers”*.

One might be tempted to think that such memory requirements are not anymore a limiting factor as, at the time of writing, even standard workstations come with 32 GB of main memory. Indeed, over the last decades, the semiconductors industry followed the exponential trends dictated by Moores’ law and yielded not only exponentially faster microprocessors but also larger memories. Unfortunately, memory latency improvements have been more modest, leading to the so called memory wall effect [Wilkes, 1995]: data access times are taking an increasingly fraction of total computation times. Thus, if Knuth [1973] wrote that *“space optimization is closely related to time optimization in a disk memory”*, forty years later one can simply say that space optimization is always related to time optimization.

Over the last years, significant effort has been devoted to the engineering of more space-efficient data structures to replace the suffix tree in practical applications. In particular, much research has been done into designing succinct (or even compressed) data structures providing efficient query times using space proportional to that of the uncompressed (or compressed) input. Thanks to these advances, a succinct index of the human genome consumes as little as 3.5 GB of memory and often even improves query time over classic indices.

In this chapter, I introduce some classic full-text indices (suffix arrays and  $q$ -gram indices) and subsequently succinct full-text indices (uncompressed variants of the FM-index). Afterwards, I introduce generic string matching algorithms that work on any of these data structures, and at the same time provide their experimental evaluation. My implementation of all these algorithms and data structures is publicly available in source form within the C++ library SeqAn [Döring *et al.*, 2008].

## 3.1 Classic full-text indices

### 3.1.1 Suffix array

The suffix array (SA) has been initially proposed by Manber and Myers [1990] as a practical replacement of the suffix tree in exact string matching. The SA reduces the space complexity of the suffix tree by a constant factor at the cost of an increase of the search time complexity by a logarithmic factor. It achieves that by explicitly representing only the leaves of the suffix trie while omitting its internal nodes and outgoing edges. Figure 3.1 shows an example. The leaves referring to sorted suffixes are sufficient to emulate suffix trie traversals. Any path from the root to some internal node is indeed computed on the fly via binary search over the sorted suffixes. Below, I define the (generalized) suffix array and then show how to emulate all top-down traversal operations.

**Definition 3.1 (Suffix array).** The *suffix array* of a padded string  $S$  of length  $n$  is an array  $A$  containing a permutation of the interval  $[0, n)$ , s.t.  $S_{A[i-1]} <_{lex} S_{A[i]}$  for all  $1 \leq i < n$ .

**Definition 3.2 (Generalized suffix array).** The *generalized suffix array* (GSA) of a padded string collection  $\mathbb{S}$  (definition 2.9), consisting of  $c$  strings of total length  $n$ , is an array  $A$  of length  $n$  containing a permutation of pairs  $(i, j)$  for all  $0 \leq i < c$  and  $0 \leq j < |S^i|$ . Pairs are ordered s.t.  $\mathbb{S}_{A[i-1]} <_{lex} \mathbb{S}_{A[i]}$  for all  $1 \leq i < n$ .

**Figure 3.1:** Suffix array and generalized suffix array. (a) Suffix array of the string  $S = \text{ANANAS\$}$ . (b) Generalized suffix array of the string collection  $\mathbb{S} = \{\text{ANANAS\$}^0, \text{CACAO\$}^1\}$ .

(a) Suffix array.			(b) Generalized suffix array.		
$i$	$A[i]$	$S_{A[i]}$	$i$	$A[i]$	$\mathbb{S}_{A[i]}$
0	6	\$	0	(0, 6)	$\$^0$
1	0	ANANAS\$	1	(1, 5)	$\$^1$
2	2	ANAS\$	2	(1, 1)	$\text{ACAO\$}^1$
3	4	AS\$	3	(0, 0)	$\text{ANANAS\$}^0$
4	1	NANAS\$	4	(0, 2)	$\text{ANAS\$}^0$
5	3	NAS\$	5	(1, 3)	$\text{AO\$}^1$
6	5	S\$	6	(0, 4)	$\text{AS\$}^0$
			7	(1, 0)	$\text{CACAO\$}^1$
			8	(1, 2)	$\text{CAO\$}^1$
			9	(0, 1)	$\text{NANAS\$}^0$
			10	(0, 3)	$\text{NAS\$}^0$
			11	(1, 4)	$\text{O\$}^1$
			12	(0, 5)	$\text{S\$}^0$



The SA is constructed in  $\mathcal{O}(n)$  time, for instance using the [Kärkkäinen and Sanders, 2003] algorithm, or using non-optimal but practically faster algorithms, e.g. [Schürmann and Stoye, 2007]. The space consumption of the SA is  $\Theta(n \log n)$  bits. When  $n < 4 \text{ GB}$ , a 32 bit integer is sufficient to encode any value in the range  $[0, n)$ . Consequently, the space consumption of suffix arrays for strings shorter than 4 GB is  $4n$  bytes.

Weese [2013] gives a generalization of Kärkkäinen and Sanders algorithm to construct the GSA in  $\mathcal{O}(n)$  time. The space consumption of the GSA is  $\Theta(n \log n^*)$  bits, where  $n^* = \max |S^i|$ . For instance, a pair of 1 + 4 bytes suffices to encode any suffix position in a collection consisting of up to 256 strings, each one shorter than 4 GB.

### Top-down traversal

I now concentrate on describing suffix trie functionalities, as I implemented them within the SeqAn library. Any suffix trie node is univocally identified by an interval of the suffix array  $A$ . For instance, the root node is represented by the interval  $[0, n)$  containing the whole SA. Thus, while traversing the trie, I remember the interval  $[l, r)$  associated to the current node. In addition, I also remember the right interval  $r_p$  of the parent node and the depth  $d$  of the current node. The incoming edge label of an internal node at depth  $d$  is identified by the  $(d - 1)$ -th symbol within the corresponding suffix  $S_{A[i]}$  with  $i \in [l, r)$ . Any leaf node is defined to have all outgoing edges labeled by terminator symbols (see section 2.4.2). The occurrences below any node correspond by definition to the interval  $[l, r)$  of  $A$ . Summing up, I represent the current node  $v$  by the integers  $\{l, r, r_p, d\}$  and define the following operations on it:

- **GOROOT**( $v$ ) initializes  $v$  to  $\{0, n, n, 0\}$ ;
- **ISLEAF**( $v$ ) returns true iff  $|S_{A[v.r-1]}| = v.d + 1$ ;
- **LABEL**( $v$ ) returns  $S_{A[v.l]+v.d}$ ;
- **OCCURRENCES**( $v$ ) returns  $A[v.l \dots v.r]$ .

Algorithm 3.1 $L(v, c)$	Algorithm 3.2 $R(v, c)$
<b>Input</b> $v$ : SA node iterator $c$ : character to query <b>Output</b> integer denoting the left interval 1: $(l_1, l_2) \leftarrow (v.l, v.r)$ 2: <b>while</b> $l_1 < l_2$ <b>do</b> 3: $i \leftarrow \lfloor \frac{l_1 + l_2}{2} \rfloor$ 4: <b>if</b> $S_{A[i]+v.d} <_{lex} c$ 5: $l_1 \leftarrow i + 1$ 6: <b>else</b> 7: $l_2 \leftarrow i$ 8: <b>return</b> $l_1$	<b>Input</b> $v$ : SA node iterator $c$ : character to query <b>Output</b> integer denoting the right interval 1: $(r_1, r_2) \leftarrow (v.l, v.r)$ 2: <b>while</b> $r_1 < r_2$ <b>do</b> 3: $i \leftarrow \lfloor \frac{r_1 + r_2}{2} \rfloor$ 4: <b>if</b> $S_{A[i]+v.d} \leq_{lex} c$ 5: $r_1 \leftarrow i + 1$ 6: <b>else</b> 7: $r_2 \leftarrow i$ 8: <b>return</b> $r_1$

**Algorithm 3.3**  $\text{goDOWN}(v)$ 


---

**Input**  $v$  : SA node iterator  
**Output** boolean indicating success  
1: **if**  $\text{ISLEAF}(v)$   
2:     **return false**  
3:  $v.l \leftarrow R(v, \$)$   
4:  $c_l \leftarrow S_{A[v.l]+v.d}$   
5:  $c_r \leftarrow S_{A[v.r-1]+v.d}$   
6:  $v.r_p \leftarrow v.r$   
7: **if**  $c_l \neq c_r$   
8:      $v.r \leftarrow R(v, c_l)$   
9:  $v.d \leftarrow v.d + 1$   
10: **return true**

---

**Algorithm 3.4**  $\text{goRIGHT}(v)$ 


---

**Input**  $v$  : SA node iterator  
**Output** boolean indicating success  
1: **if**  $v.r = v.r_p$   
2:     **return false**  
3:  $v.l \leftarrow v.r$   
4:  $c_l \leftarrow S_{A[v.l]+v.d}$   
5:  $c_r \leftarrow S_{A[v.r-1]+v.d}$   
6:  $v.r \leftarrow v.r_p$   
7: **if**  $c_r \neq c_p$   
8:      $v.r \leftarrow R(v, c_r)$   
9: **return true**

---

Binary search is the key to implement function  $\text{goDOWN}$  a symbol. Functions  $L$  (algorithm 3.1) and  $R$  (algorithm 3.2) compute in  $\mathcal{O}(\log n)$  binary search steps the position in  $A$  of the left and right intervals corresponding to the child node that is reached by going down the edge labeled by a given symbol  $c$ . Algorithms 3.3 and 3.4 show how to implement respectively  $\text{goDOWN}$  and  $\text{goRIGHT}$  with a time complexity independent of the alphabet size  $\sigma$ . As they rely on a single call of  $R$ , their time complexity is  $\mathcal{O}(\log n)$ . In this way, the SA supports exact string matching in  $\mathcal{O}(m \log n)$  time (see algorithm 3.11).

Note that the suffix array can be binary searched by spelling a full pattern within a single call of  $L$  and  $R$ : in line 6 of algorithms 3.1–3.2, instead of comparing a single character, it suffices to compare the full pattern to the current suffix. Nonetheless, the worst case runtime of algorithm 3.11 stays  $\mathcal{O}(m \log n)$ , as each step of the binary search requires

**Algorithm 3.5**  $\text{goDOWN}(v, c)$ 


---

**Input**  $v$  : SA node iterator  
 $c$  : character to query  
**Output** boolean indicating success  
1: **if**  $\text{ISLEAF}(v)$   
2:     **return false**  
3:  $l \leftarrow L(v, c)$   
4:  $r \leftarrow R(v, c)$   
5: **if**  $l < r$   
6:      $v \leftarrow \{l, r, v.r, v.d + 1\}$   
7:     **return true**  
8: **else**  
9:     **return false**

---

a full lexicographical comparison between the pattern and any suffix of the text, which takes  $\mathcal{O}(m)$  time in the worst case. As shown in [Manber and Myers, 1990], the worst case runtime can be decreased to  $\mathcal{O}(m + \log n)$  at the expense of additional  $\Theta(n \log n)$  bits, by storing the precomputed longest common prefixes (LCP) between any two consecutive suffixes  $S_{A[i]}, S_{A[i+1]}$  for all  $1 \leq i < n$ . Alternatively, the *average case* runtime is reduced to  $\mathcal{O}(m + \log n)$ , without storing any additional information, by using the MLR heuristic [Manber and Myers, 1990]. In practice, the MLR heuristic outperforms the SA + LCP algorithm, due to the higher cost of fetching additional data from the LCP table [Weese, 2013].

### 3.1.2 Suffix tree realizations

I briefly introduce two suffix tree realizations from the literature: the *enhanced suffix array* (ESA) [Abouelhoda *et al.*, 2004] and the *lazy suffix tree* (LST) [Giegerich *et al.*, 2003]. These realizations explicitly implement or implicitly emulate a suffix tree rather than a suffix trie. The string matching algorithms of section 3.3 work on tries, yet any tree can be easily traversed as a trie. The implementation of these data structures within SeqAn, equally generalized to multiple sequence, is due to Weese. Hence, for an extensive illustration, the reader is invited to consult [Weese, 2013].

#### Enhanced suffix array

The enhanced suffix array (ESA) [Abouelhoda *et al.*, 2004] supplements the SA and LCP tables (see section 3.1.1) with another table called *child* table. Each SA value represents one leaf of the suffix tree, while each LCP value represent the length of one edge of the suffix tree. What is still missing, in order to represent a full suffix tree, are the SA intervals of the children of each inner node. These intervals would have to be computed in logarithmic time by `goRight` during a top-down SA traversal. As proposed by Abouelhoda *et al.* [2004], these intervals are computed in linear time, within one single bottom-up traversal, and stored in the child table, which consumes additional  $\Theta(n \log n)$  bits, thus  $4n$  bytes for collections smaller than 4 GB.

#### Lazy suffix tree

The lazy suffix tree (LST) [Giegerich *et al.*, 1999] variant proposed by [Weese, 2013] is composed by a partially sorted SA plus a node directory. The SA initially reflects the ordering of the suffixes up to depth 1, and the node directory table contains only the root node. During a top-down traversal, the current node at depth  $i$  is expanded by means of the *wotd* algorithm [Giegerich *et al.*, 1999], which calls one round of radix sort to refine the ordering of the suffixes up to depth  $i + 1$  and inserts the newly computed children nodes in the directory. The construction of the full LST takes  $\mathcal{O}(n^2 + \sigma n)$  time in the worst case.

### 3.1.3 $q$ -Gram index

If the traversal of the idealized suffix trie is bounded to a maximum depth  $q$ , the logarithmic factor paid by using the SA vanishes. The idea is to supplement the SA with a so-called  $q$ -gram *directory*: an additional array  $D$  of  $\sigma^q + 1$  integers that represent the SA ranges computed by algorithm 3.1 for all possible words of length  $q$  over  $\Sigma$ .

With the aim of addressing  $q$ -grams in the directory  $D$ , I impose a canonical code on  $q$ -grams through a bijective function  $h : \Sigma^q \rightarrow [0, \sigma^q)$  defined as in [Knuth, 1973]:

$$h(P) = \sum_{i=0}^{q-1} \text{rank}(p_i) \cdot \sigma^{q-i} \quad (3.1)$$

where  $P \in \Sigma^q$  is some  $q$ -gram and  $\text{rank}$  is the lexicographic rank defined on  $\Sigma$  (see definition 2.6). The canonical code assigned by  $h$  preserves the lexicographical ordering for all words not longer than  $q$ , i.e.  $X <_{lex} Y$  iff  $h(X) < h(Y)$  for all  $X, Y \in \Sigma^{\leq q}$ . The hash function  $h$  allows to store in and retrieve from  $D$  the left SA interval returned by algorithm 3.1 for each  $q$ -gram, i.e.  $D[h(P)] = L(0, n, P)$  for all  $P \in \Sigma^q$ . Note that the right interval returned by algorithm 3.2 is equivalent to the left interval of the following  $q$ -gram in lexicographical order and therefore available in  $D[h(P) + 1]$ .

In practice, the  $q$ -gram index is applicable only to relatively small alphabets and tree depths. For instance, parameters  $\sigma = 4$  and  $q = 14$  imply a  $q$ -gram directory consisting of 268 M entries. Using a 32 bits integer encoding, such directory alone consumes 1 GB of memory.

**Figure 3.2:** 2-Gram index of the string  $S = \text{ANANAS\$}$  over the alphabet  $\Sigma = \{A, N, S\}$ . The example shows the lookup of the 2-gram NA. The hash value  $h(\text{NA}) = 3$  addresses two lookups, respectively in  $D[3]$  and  $D[4]$ , that in turn provide the range  $[4, 6)$  in  $A$ .

$P$	$h(P)$	$D[h(P)]$	$i$	$A[i]$	$S_{A[i]}$
AA	0	1	0	6	\$
AN	1	1	1	0	ANANAS\$
AS	2	3	2	2	ANAS\$
NA	3	4	3	4	AS\$
NN	4	6	4	1	NANAS\$
NS	5	6	5	3	NAS\$
SA	6	6	6	5	S\$
SN	7	6			
SS	8	6			
	9	6			

### Top-down traversal

I now describe how I extended the SA traversal operations of section 3.1.1 to use the  $q$ -gram directory  $D$  within the generic text indexing framework of the SeqAn library. Again, while traversing the trie, I maintain the current range  $[l, r)$  and the current depth  $d$ . In addition, I maintain the interval  $[l_h, r_h)$  in  $D$  and the incoming edge label  $e$  node in order to answer LABEL( $v$ ). Summing up, I represent the current node  $v$  by the elements  $\{l, r, d, l_h, r_h, e\}$ . I define the basic node operations as follows:

- GOROOT( $v$ ) initializes  $v$  to  $\{0, n, 0, 0, \sigma^q, \epsilon\}$ ;
- ISLEAF( $v$ ) returns true iff  $v.d = q$ ;
- LABEL( $v$ ) returns  $v.e$ ;
- OCCURRENCES( $v$ ) returns  $A[v.l \dots v.r]$ .

Algorithms 3.6 and 3.7 show respectively functions L and R using the directory  $D$  instead of  $A$ . In this way, both variants of GOWDOWN (algorithms 3.5 and 3.3) and GORIGHT (algorithm 3.4) take  $\mathcal{O}(1)$  time. Conversely, GOWDOWN and GORIGHT are implemented by the generic algorithms 2.2 and 2.3 running in time  $\mathcal{O}(\sigma)$ .

The directory  $D$  alone is sufficient for top-down traversals bounded to a maximum depth  $q$ ; the suffix array  $A$  is accessed only to locate text locations pointed by the leaves. In this case, the total ordering of the text suffixes in the SA can be relaxed to prefixes of length  $q$ . This gives a twofold advantage, as one can (i) construct the SA more efficiently using bucket sorting and (ii) maintain leaves in each bucket sorted by their relative text positions. The latter property allows to compress the SA bucket-wise e.g. using Elias  $\delta$  encoding [Elias, 1975] or to devise cache-oblivious strategies to process the occurrences [Hach *et al.*, 2010].

The directory  $D$  is still usable, even if the traversal needs to go below depth  $q$ . An hybrid traversal can use the directory  $D$  up to depth  $q$  and later continue with binary searches on the suffix array  $A$ . This hybrid traversal cuts the most expensive binary searches and increases memory locality. Furthermore, this traversal becomes useful whenever the SA is too big to fit in main memory and has to reside in external memory. In section 3.3, I show the experimental evaluation of this hybrid  $q$ -gram index.

Algorithm 3.6 L( $v, c$ )	Algorithm 3.7 R( $v, c$ )
<b>Input</b> $v$ : $q$ -gram index iterator $c$ : character to query	<b>Input</b> $v$ : $q$ -gram index iterator $c$ : character to query
<b>Output</b> integer denoting the left interval 1: $v.l_h \leftarrow v.l_h + \text{rank}(c) \cdot \sigma^{v.d}$ 2: <b>return</b> $D[v.l_h]$	<b>Output</b> integer denoting the right interval 1: $v.r_h \leftarrow v.r_h - \text{rank}(c) \cdot \sigma^{v.d}$ 2: <b>return</b> $D[v.r_h]$

---

**Algorithm 3.8**  $\text{GoDown}(v, c)$ 


---

**Input**      $v$  :  $q$ -gram index iterator  
                $c$  : character to query

**Output**    boolean indicating success

```

1: if ISLEAF( $v$ )
2:   return false
3:  $v.d \leftarrow v.d + 1$ 
4:  $v.e \leftarrow c$ 
5:  $v.l \leftarrow L(v, v.e)$ 
6:  $v.r \leftarrow R(v, v.e)$ 
7: return  $v.l < v.r$ 

```

---

### 3.1.4 Trie and radix tree realizations

Before turning to succinct full-text indices, I briefly describe how I reused the text-indices just exposed to implement tries and radix trees in the SeqAn library. A trie is easily emulated by means of a partial SA. I index only the first suffix of each string in the collection and subsequently construct the SA-based trie via quicksort in time  $\mathcal{O}(n \log n)$ , where  $n$  is the total length of the string collection. The top-down traversal based on binary search still works as described in section 3.1.1. This trie is also extendable by a  $q$ -gram directory as in section 3.1.3. A radix tree is constructed in an analogous way, starting from the LST of section 3.1.2. I fill the LST's partial SA as described above and subsequently apply the *wotd* algorithm [Giegerich *et al.*, 1999] to construct the radix tree in time  $\mathcal{O}(\sigma n)$ .

## 3.2 Succinct full-text indices

The *Burrows-Wheeler transform* (BWT) [Burrows and Wheeler, 1994] is a transformation defining a permutation of an input string. Such transformation exposes two important properties: it is *reversible* and it tends to produce *runs of equal characters*. The former property allows to recover the original string from its BWT alone, while the latter property makes the transformed string more amenable to compression. Because of these two properties, the BWT is a fundamental method for text compression.

Some years after its introduction, Ferragina and Manzini proposed the BWT for full-text indexing. They show in [Ferragina and Manzini, 2000] that the BWT allows to perform exact string matching and engineer in [Ferragina and Manzini, 2001] a compressed full-text index called FM-index. Over the last years, the FM-index has been widely employed under different re-implementations by many popular bioinformatics tools e.g. Bowtie [Langmead *et al.*, 2009] and BWA [Li and Durbin, 2009], and is now considered a fundamental method for the indexing of genomic sequences.

In the following, I give the fundamental ideas behind the BWT. Subsequently, I discuss my generalized FM-index implementation covering strings and string collections.

### 3.2.1 Burrows-Wheeler transform

Let  $S$  be a padded string (definition 2.8) of length  $n$  over  $\Sigma_\$$ . In the following, consider the string  $S$  to be *cyclic* and its subscript  $s_i$  to be *modular*, i.e.  $s_i = s_{i \bmod n}$  for any  $i \in \mathbb{Z}$ .

**Definition 3.3 (BWT).** Given a string  $S$  and its SA  $A$ , the BWT  $L$  of  $S$  is a string over  $\Sigma_\$$  with  $l_i = s_{A[i]-1}$ .

To visualize the BWT, consider the square matrix consisting of all cyclic shifts of the string  $S$  sorted in lexicographical order. Note how the cyclic shifts matrix is related to the suffix array  $A$  of  $S$ : the  $i$ -th cyclic shift is  $S_{A[i] \dots n} S_{0 \dots A[i]}$ . The BWT is the string obtained by concatenating the symbols in the last column of the sorted cyclic shifts matrix of  $S$ . Figure 3.3 shows an example.

Note that the BWT easily generalizes to string collections. Indeed, definition 3.3 still holds for a padded string collection  $\mathbb{S}$  (definition 2.9) and its cyclic shifts matrix sorted in lexicographical order.

The cyclic shifts matrix is conceptual and does not have to be constructed explicitly to derive the BWT of  $S$ . The BWT can be obtained in linear time by scanning the SA  $A$  and assigning the symbol  $s_{A[i]-1}$  to the  $i$ -th BWT symbol. However, constructing the BWT from the SA is not desirable in practice, especially for strings over small alphabets like DNA, as the SA consumes  $\Theta(n \log n)$  bits in addition to the  $\Theta(n \log \sigma)$  bits of the BWT. Therefore, some direct BWT construction algorithms working within  $o(n \log \sigma)$  bits plus constant space have been recently proposed in [Bauer *et al.*, 2013; Crochemore *et al.*, 2013].

#### Inversion

I now describe how to invert the BWT to reconstruct the original string. For convenience, I denote the first column  $s_{A[i]}$  by  $F$  and the last column  $s_{A[i]-1}$  by  $L$ . Inverting the BWT means being able to know where any BWT character occurs in the original text. To this intent, I define two permutations  $LF : [0, n) \rightarrow [0, n)$  and  $\Psi : [0, n) \rightarrow [0, n)$ , with  $LF = \Psi^{-1}$ , where the value of  $LF(i)$  gives the position  $j$  in  $F$  where character  $l_i$  occurs and the value  $\Psi(j)$  gives back the position  $i$  in  $L$  where  $f_j$  occurs. Figure 3.4 illustrates.

**Figure 3.3:** Cyclic shifts matrix of the string  $S = \text{ANANAS\$}$ . Column  $s_{A[i]-1}$  is the BWT of  $S$ .

$i$	$A[i]$	$s_{A[i]}$	...	$s_{A[i]-1}$
0	6	\$	ANANA	S
1	0	A	NANAS	\$
2	2	A	NAS\$A	N
3	4	A	S\$ANA	N
4	1	N	ANAS\$	A
5	3	N	AS\$AN	A
6	5	S	\$ANAN	A

**Figure 3.4:** Permutations  $LF$  and  $\Psi$  of the string  $S = \text{ANANAS\$}$ . The example shows that  $LF = \Psi^{-1}$ , e.g.  $LF(\Psi(5)) = 5$  and  $\Psi(LF(3)) = 3$ . Moreover, the example shows that the relative order of characters between  $L$  and  $F$  is preserved, e.g. the first occurrence of  $N$  in  $L$  corresponds to the first occurrence in  $F$ .

$i$	$\Psi(i)$	$LF(i)$	$S_{A[i]}$	...	$S_{A[i]-1}$
0	1	6	\$	ANANA	S
1	4	0	A	NANAS	\$
2	5	(4)	A	NAS\$A	(N)
3	6	5	A	S\$ANA	N
4	(2)	1	(N)	ANAS\$	A
5	3	2	N	AS\$AN	A
6	0	3	S	\$ANAN	A

I define the iterated  $\Psi$  as

$$\begin{aligned}\Psi^0(j) &= j \\ \Psi^{i+1}(j) &= \Psi(\Psi^i(j))\end{aligned}\tag{3.2}$$

and the iterated  $LF$  as

$$\begin{aligned}LF^0(j) &= j \\ LF^{i+1}(j) &= LF(LF^i(j)).\end{aligned}\tag{3.3}$$

The character  $s_i$  corresponds to  $f_{\Psi^{i+1}(j)}$ , where  $j$  is the position of  $\$$  in  $F$ ; analogously,  $\bar{s}_i$  corresponds to  $l_{LF^{i+1}(j)}$ , where  $j$  is the position of  $\$$  in  $L$ . Therefore, the full string  $S$  is recovered by starting in  $F$  at the position of  $\$$  and following the cycle of length  $n$  defined by the permutation  $\Psi$ ; analogously, the reverse string  $\bar{S}$  is recovered by starting in  $L$  at the position of  $\$$  and following the cycle of length  $n$  defined by the permutation  $LF$ . Figure 3.5 exemplifies.

Inverting the generalized BWT works in the same way. Indeed, permutations  $\Psi$  and  $LF$  are composed of  $c$  cycles, where each cycle corresponds to a distinct string  $S^i$  in the collection. The string  $S^i$  is recovered by starting at the position of  $\$^i$  and following the cycle of  $\Psi$  (or  $LF$ ) associated to  $S^i$ . See [Mantaci *et al.*, 2007] for further details on the generalized BWT.

### Permutation $LF$

Permutation  $LF$  needs not to be stored explicitly. Differently from  $\Psi$ , it is possible to deduce  $LF$  from the BWT string  $L$  with the help of some additional character counts. The counting argument relies on two simple observations of Burrows and Wheeler [1994] on the cyclic shifts matrix.

**Observation 3.1.** For all  $i \in [0, n)$ , the character  $l_i$  precedes the character  $f_i$  in  $S$ .



**Figure 3.5:** Recovering the string  $S = \text{ANANAS\$}$  from permutation  $\Psi$ . The example shows only the first two steps of the inversion recovering AN.

$s_{A[i]}$	$i$	$\Psi(i)$
\$	0	1
(A)	1	4
A	2	5
A	3	6
(N)	4	2
N	5	3
S	6	0

**Observation 3.2.** For all characters  $c \in \Sigma$ , the  $i$ -th occurrence of  $c$  in  $F$  corresponds to the  $i$ -th occurrence of  $c$  in  $L$ .

These observation are evident, indeed  $F = s_{A[i]}$  and  $L = s_{A[i]-1}$  (see figure 3.4). Given the two above observations, Ferragina and Manzini [2000] define the permutation  $LF$  as:

$$LF(i) = C(l_i) + \text{Occ}(l_i, i) \quad (3.4)$$

where  $C : \Sigma \rightarrow [0, n)$  denotes the total number of occurrences in  $S$  of all characters alphabetically smaller than  $c$ , and  $\text{Occ} : \Sigma \times [0, n) \rightarrow [0, n)$  the number of occurrences of character  $c$  in the prefix  $L_{0 \dots i}$ . The key problem of encoding the permutation  $LF$  lies in representing function  $\text{Occ}$ , as function  $C$  is easily tabulated by a small array of size  $\Theta(\sigma \log n)$  bits. In the next subsection, I address the problem of representing function  $\text{Occ}$  efficiently. Subsequently, I explain how to implement generic full-text index traversal using the permutation  $LF$ .

### 3.2.2 Rank dictionaries

Function  $\text{Occ}$ , answering the question “how many times a given character  $c$  occurs in the prefix  $L_{0 \dots i}$ ?”, has to be answered efficiently, ideally in constant time and succinct space. The general problem on arbitrary strings has been tackled by several studies on the succinct representation of data structures [Jacobson, 1989]. This specific question takes the name of *rank query* and a data structure answering rank queries is called *rank dictionary* (RD).

**Definition 3.4 (Rank query).** Given a string  $S$  over an alphabet  $\Sigma$  and a character  $c \in \Sigma$ ,  $\text{RANK}_c(S, i)$  returns the number of occurrences of  $c$  in the prefix  $S_{0 \dots i}$ <sup>1</sup>.

The key idea of RDs is to maintain a succinct (or even compressed) representation of the input string and attach a dictionary to it. By doing so, Jacobson [1989] shows

<sup>1</sup> Note that symbol  $s_i$  is excluded from  $\text{RANK}_c(S, i)$ .

how to answer rank queries in constant time (on the RAM model) using  $o(n)$  additional bits for an input binary string of  $n$  bits. Here, I cover only the most practical succinct RDs and discuss some implementation aspects, crucial to obtain practical efficiency. I first consider the binary alphabet  $\mathbb{B} = \{0, 1\}$  and subsequently the case of an arbitrary alphabet.

### Binary alphabet

Here, I follow the explanation of [Navarro and Mäkinen, 2007]. Hence, I start by describing a simple *one-level* rank dictionary answering rank queries in constant time but consuming  $\mathcal{O}(n)$  additional bits. Subsequently, I describe an extended *two-levels* RD consuming only  $o(n)$  additional bits. Finally, I briefly discuss my implementation of practical *multi-level* RDs.

The one-level binary RD partitions the binary input string  $S \in \mathbb{B}^*$  in blocks of  $b$  symbols and complements it with an array  $R$  of length  $\lfloor n/b \rfloor + 1$ . The  $j$ -th entry of  $R$  provides a summary of the number of occurrences of the bit 1 in  $S$  before position  $jb$ , i.e.  $R[0] = 0$  and  $R[j] = \text{RANK}_1(S, jb)$  for any  $j > 0$ . Note that  $R$  summarizes only  $\text{RANK}_1$ , as  $\text{RANK}_0(S, i) = i - \text{RANK}_1(S, i)$ . Therefore, the rank query is rewritten as:

$$\text{RANK}_1(S, i) = R[\lfloor i/b \rfloor] + \text{POPCOUNT}_1(S_{\lfloor i/b \rfloor \cdot b}, i \bmod b). \quad (3.5)$$

The query is answered in time  $\mathcal{O}(b)$  by (i) fetching the rank summary from  $R$  in constant time and (ii) counting the number of occurrences of the bit 1 within a block of  $\mathcal{O}(b)$  bits. Figure 3.6a illustrates. Jacobson poses  $b = \lfloor \log n \rfloor$  and answers  $\text{POPCOUNT}$  in time  $\mathcal{O}(1)$  with the four-Russians tabulation technique [Arlazarov *et al.*, 1970]. As the array  $R$  stores  $\mathcal{O}(\lfloor n/\log n \rfloor)$  positions and each position in  $S$  requires  $\lfloor \log n \rfloor$  bits,  $R$  consumes  $\mathcal{O}(n)$  bits.

A binary *two-levels* RD squeezes additional space consumption down to  $o(n)$  bits. The idea is to divide  $S$  in *superblocks* of length  $b^2$  and add another array  $R^2$  to store the number of occurrences of the bit 1 in  $S$  before each superblock. The initial array  $R$  now counts bit 1 only within the overlapping superblock defined by  $R^2$ . Figure 3.6b exemplifies. The rank query becomes:

$$\text{RANK}_1(S, i) = R^2[\lfloor i/b^2 \rfloor] + R[\lfloor i/b \rfloor] + \text{POPCOUNT}_1(S_{\lfloor i/b \rfloor \cdot b}, i \bmod b). \quad (3.6)$$

Each entry of  $R$  now represents only  $b^2$  possible values and thus consumes only  $\mathcal{O}(\log b^2) = \mathcal{O}(\log b)$  bits. Summing up,  $R^2$  adds  $\mathcal{O}(n \log n/b^2)$  bits and  $R$  adds  $\mathcal{O}(n \log b/b)$  bits. By posing  $b = \log n$  as above, it follows  $\mathcal{O}(n/\log n)$  bits for  $R^2$  and  $\mathcal{O}(n \log \log n/\log n)$  bits for  $R$ . Hence, the two-levels binary RD consumes  $o(n)$  additional bits.

I implemented generic *multi-levels* RDs, where the block size  $b$  is a template parameter adjustable at compile time. Typically, the input string is smaller than 4 GB, thus I employ the one-level RD with  $b = 32$  bits or the two-level RD with  $b = 16$  bits and  $b^2 = 32$  bits; otherwise, for longer strings, I employ the two-levels RD with  $b = 32$  bits and  $b^2 = 64$  bits, or the three-levels RD with  $b = 16$  bits,  $b^2 = 32$  bits and  $b^3 = 64$  bits. In order to reduce the number of cache misses, the succinct representation of the input string is *interleaved* with the lowest level summaries array  $R$ . Moreover, I use the SSE 4.2 `popcnt` instruction [Intel, 2011] to count symbols within a block in time  $\mathcal{O}(b/w)$ , where  $w$  is the total SSE register width (on modern processors  $w = 256$  bits).

**Figure 3.6:** Binary rank dictionaries (RDs) of the string  $S = 010101100100$ . (a) One-level RD with  $b = 4$ . In the example,  $\text{RANK}_1(S, 6) = R[1] + \text{POPCOUNT}_1(S_4, 2) = 3$ . (b) Two-levels RD with  $b = 2$ . In the example,  $\text{RANK}_1(S, 6) = R^2[1] + R[2] + \text{POPCOUNT}_1(S_4, 2) = 3$ . Note that  $R$  in (a) is different from  $R$  in (b).

(a) One-level rank dictionary.

$i$	$s_i$	$R[\lfloor i/4 \rfloor]$
0	0	0
1	1	
2	0	
3	1	
4	$\begin{pmatrix} 0 \\ 1 \end{pmatrix}$	$\begin{pmatrix} 2 \end{pmatrix}$
5	$\begin{pmatrix} 1 \end{pmatrix}$	
6	1	
7	0	
8	0	4
9	1	
10	0	
11	0	

(b) Two-levels rank dictionary.

$i$	$s_i$	$R[\lfloor i/2 \rfloor]$	$R^2[\lfloor i/4 \rfloor]$
0	0	0	0
1	1		
2	0	1	
3	1		
4	$\begin{pmatrix} 0 \\ 1 \end{pmatrix}$	$\begin{pmatrix} 0 \end{pmatrix}$	$\begin{pmatrix} 2 \end{pmatrix}$
5	$\begin{pmatrix} 1 \end{pmatrix}$		
6	1	1	
7	0		
8	0	0	4
9	1		
10	0	1	
11	0		

### Small alphabets

The extension of binary RDs to arbitrary alphabets is easy. However, the space consumption of such RD has a linear dependency in the alphabet size. This fact renders such extension appealing only for small alphabets like DNA. Here, I show how to extend the one-level RD.

Consider an input string  $S$  of length  $n$  over  $\Sigma$ , thus consisting of  $\Theta(n \log \sigma)$  bits. As in the binary case, this one-level RD partitions  $S$  in blocks of  $b$  bits. It complements the string  $S$  with a matrix  $R_\sigma$  of size  $\lfloor n/b \rfloor + 1 \times \sigma$ , summarizing the number of occurrences for each symbol in  $\Sigma$ . The rank query is rewritten accordingly:

$$\text{RANK}_c(S, i) = R_\sigma[\lfloor i/b \rfloor, \text{rank}(c)] + \text{POPCOUNT}_c(S_{\lfloor i/b \rfloor \cdot b}, i \bmod b). \quad (3.7)$$

Figure 3.7 shows an example of one-level DNA RD. Answering this query requires counting the number of occurrences of the character  $c$  inside a block of  $b$  bits. In order to answer this query in constant time, I consider blocks of  $\lfloor \log n \rfloor / \lfloor \log \sigma \rfloor$  symbols, i.e. each block consumes not more than  $\lfloor \log n \rfloor$  bits as in the binary RD case. The matrix  $R_\sigma$  has thus  $\sigma$  rows and  $\mathcal{O}(n \log \sigma / \log n)$  columns, each one consuming  $\lfloor \log n \rfloor$  bits. Thus,  $R_\sigma$  adds  $\mathcal{O}(\sigma n \log \sigma)$  bits to  $S$ .

**Figure 3.7:** One-level DNA rank dictionary of the string  $S = \text{CTCGCA}$  with  $b = 2$ . In the example,  $\text{RANK}_C(S, 4) = R_\sigma[1, 1] + \text{POPCOUNT}_C(S_2, 2) = 2$ .

$i$	$s_i$	$R_\sigma[\lfloor i/2 \rfloor, \{A, C, G, T\}]$
0	C	[0, 0, 0, 0]
1	T	
2	C	[0, 1, 0, 1]
3	G	
4	C	[0, 2, 1, 1]
5	A	

### Wavelet tree

To mitigate the factor  $\sigma$  affecting the RD just exposed, Grossi *et al.* [2003] propose a *hierarchical* RD, called the *wavelet tree* (WT). This tree data structure recursively partitions the alphabet  $\Sigma$  in balanced subsets and therefore decomposes the input string in *subsequences* containing symbols from one subset. Any tree node represents one alphabet partition and its associated subsequence. I first give the formal definition of WT and then discuss how to answer rank queries.

**Definition 3.5 (Wavelet tree).** The wavelet tree of a string  $S \in \Sigma^*$  is a balanced binary tree of height  $\lceil \log \sigma \rceil$ . The root represents all symbols in  $\Sigma$  and each leaf exactly one symbol  $c \in \Sigma$ . Any non-leaf node  $v$  represents some subset of symbols  $\Sigma_v$  whose lexicographic rank is in range  $[i, j]$  i.e.  $\Sigma_v = \{c \in \Sigma : \text{rank}(c) \in [i, j]\}$ , its left child  $l$  represents the subset  $\Sigma_l$  of symbols in range  $[i, \frac{i+j}{2})$  while its right child  $r$  represents in  $\Sigma_r$  those in range  $[\frac{i+j}{2}, j)$ . Node  $v$  implicitly represents the subsequence  $S^v$  of all symbols of  $S$  in  $\Sigma_v$  and explicitly encodes its decomposition as a binary string  $B^v$  s.t.  $b_i^v = 0$  if  $s_i^v \in \Sigma_l$  and 1 otherwise [Grossi *et al.*, 2003; Navarro and Mäkinen, 2007].

Any query  $\text{RANK}_C(S, i)$  is decomposed as a sequence of  $\mathcal{O}(\log \sigma)$  binary rank queries. The sequence of queries starts in the root node and follows the path to the leaf corresponding to symbol  $c$ . On any non-leaf node  $v$ , the traversal goes left if  $c$  belongs to  $\Sigma_l$ , otherwise it goes right. Suppose w.l.o.g. that  $c$  belongs to  $\Sigma_l$ . The rank of symbol  $c$  in  $S^v$  is established as  $\text{RANK}_0(B^v, j)$ , where  $j$  is the rank of  $c$  in the parent node or  $i$  in the root node. Figure 3.8 illustrates.

The WT encodes any binary string  $B^v$  associated to some non-leaf node  $v$  using a separate binary RD. The WT contains  $\lceil \log \sigma - 1 \rceil$  non-leaf levels and any such level encodes  $n$  bits overall. Using two-levels binary RDs (section 3.2.2), the WT consumes  $(n + o(n)) \log \sigma$  bits, i.e.  $n \log \sigma(1 + o(1))$  bits, and answers any rank query in time  $\mathcal{O}(\log \sigma)$ . At the same time, the WT does not need to store the original input string  $S$  of  $\Theta(n \log \sigma)$  bits.

### 3.2.3 FM-index

I now come back to the problem of implementing a full-text index based on the permutation  $LF$ . First, I show how to emulate a top-down traversal of the suffix trie, which is sufficient to count the number of occurrences of any substring in the original text. Later, I focus on how to represent the leaves of the suffix trie, which are necessary to locate occurrences in the original text.

#### Top-down traversal

The top-down traversal of  $ST(\$)$  spells all suffixes of  $\$$ . Therefore, to emulate the traversal, it is sufficient to enumerate all suffixes of  $\$$  by inverting the BWT using the counting argument of permutation  $LF$ . This enumeration provides not only all suffixes of  $\$$ , but also their associated SA intervals, as permutation  $LF$  and the SA are closely related. However, we have seen that permutation  $LF$  enumerates suffixes of  $\bar{\$}$  in *backward* direction, while we would rather need permutation  $\Psi$  to enumerate suffixes of  $\bar{\$}$  in *forward* direction. Since  $LF = \Psi^{-1}$ , we consider the BWT of  $\bar{\$}$ , so that its permutation  $LF$  enumerates suffixes of  $\bar{\$}$ .  $LF$  of  $\bar{\$}$  computes intervals on the backward SA  $\bar{A}$  of  $\bar{\$}$ , so OCCURRENCES has to reverse the occurrences of  $\bar{A}$  back to those of  $A$ .

I represent the current node  $v$  by the elements  $\{l, r, e\}$ , where  $[l, r)$  represents the current SA interval and  $e$  is the entering edge label. Therefore, I define the following node operations:

- $GO\_ROOT(v)$  initializes  $v$  to  $\{0, n, \epsilon\}$ ;
- $IS\_LEAF(v)$  returns true iff  $Occ(\$, v.r) - Occ(\$, v.l) = v.r - v.l$ ;
- $LABEL(v)$  returns  $v.e$ ;

**Figure 3.8:** Wavelet tree of the DNA string  $S = \text{CTCGCA}$ . The alphabet  $\{A, C, G, T\}$  is recursively partitioned as  $\{\{A\}, \{C\}\}, \{\{G\}, \{T\}\}$ . In the example,  $RANK_C(S, 4) = 2$  is decomposed as  $RANK_0(B, 4) = 2$  on the root node and then  $RANK_1(B^{AC}, 2) = 2$  on the left inner node.



---

**Algorithm 3.9**  $\text{goDOWN}(v, c)$ 


---

**Input**      $v$  : FM-index iterator  
               $c$  : character to query

**Output**    boolean indicating success

```

1: if ISLEAF( $v$ )
2:   return false
3:  $v.l \leftarrow C(c) + Occ(c, v.l)$ 
4:  $v.r \leftarrow C(c) + Occ(c, v.r)$ 
5:  $v.e \leftarrow c$ 
6: return  $v.l < v.r$ 

```

---

- $\text{OCCURRENCES}(v)$  returns each value  $(i, j)$  of  $\bar{A}[l, r]$  as  $(i, |S^i| - j - v.d)$ .

The traversal easily goes from the root node to a child node following the edge labeled by character  $c$ , as it suffices to determine the interval  $[C(c), C(\text{next}_{lex}c))$ . Now, suppose the traversal is on an arbitrary suffix trie node  $v$  of known interval  $[v.l, v.r)$  s.t. the path from the root to  $v$  spells the substring  $\text{CONCAT}(v)$ . The traversal goes down to a child node  $w$  of unknown interval  $[w.l, w.r)$  s.t. the path from the root to  $w$  spells  $\text{CONCAT}(v) \cdot c$  for some  $c \in \Sigma$ . Thus, the known interval  $[v.l, v.r)$  contains all suffixes of  $\mathbb{S}$  starting with  $\text{CONCAT}(v)$ , i.e. all prefixes of  $\mathbb{S}$  ending with  $\text{CONCAT}(v)$ , while the unknown interval  $[w.l, w.r)$  contains all suffixes of  $\mathbb{S}$  starting with  $\text{CONCAT}(v) \cdot c$ , i.e. all prefixes of  $\mathbb{S}$  ending with  $c \cdot \text{CONCAT}(v)$ . All these characters  $c$  are in  $L_{w.l \dots w.r}$ , since  $l_i$  is the character  $\mathbb{S}_{A[i]-1}$  preceding the suffix pointed by  $A[i]$ . Moreover, these characters  $c$  are contiguous and in relative order in  $F$  (see observations 3.1–3.2). If  $i, \dots, j$  are all positions in  $L$  s.t.  $l_i = c$  and  $v.l \leq i < \dots < j < v.r$ , then  $w.l = LF(i)$  and  $w.r = LF(j) + 1$ . Therefore,  $w.l$  becomes:

$$\begin{aligned}
 LF(i) &= C(l_i) + Occ(l_i, i) \\
 &= C(c) + Occ(c, i) \\
 &= C(c) + Occ(c, v.l)
 \end{aligned} \tag{3.8}$$

and analogously  $w.r$  becomes:

$$\begin{aligned}
 LF(j) + 1 &= C(l_j) + Occ(l_j, j) + 1 \\
 &= C(c) + Occ(c, j) + 1 \\
 &= C(c) + Occ(c, v.r)
 \end{aligned} \tag{3.9}$$

Algorithm 3.9 uses the above equations for the operation  $\text{goDOWN}$  a symbol. Conversely,  $\text{goDOWN}$  and  $\text{goRIGHT}$  are provided by the generic algorithms 2.2 and 2.3 which enumerate all symbols in  $\Sigma$ . Thus,  $\text{goDOWN}$  a symbol runs in time  $\mathcal{O}(1)$ , while  $\text{goDOWN}$  and  $\text{goRIGHT}$  run in time  $\mathcal{O}(\sigma)$ .

### Sparse suffix array

The suffix array  $A$  is required to implement function OCCURRENCES, however the whole SA would take  $n \log n$  bits of space. As proposed by Ferragina and Manzini [2000], I maintain a *sparse* suffix array  $A^\epsilon$  containing positions sampled at regular intervals from the input string. In order to determine if and where I sampled any  $A[i]$  in  $A^\epsilon$ , I employ a binary rank dictionary  $I$  of length  $n$  s.t.  $I[i] = 1$  iff  $A[i]$  is in  $A^\epsilon$  [ $\text{RANK}_1(I, i)$ ]. I obtain any  $A[i]$  by finding the smallest  $j \geq 0$  such that  $LF^j(i)$  is in  $A^\epsilon$ , and then  $A[i] = A[LF^j(i)] + j$ .

By sampling one text position out of  $\log^{1+\epsilon} n$ , for some  $\epsilon > 0$ , then  $A^\epsilon$  consumes  $\mathcal{O}(n/\log^\epsilon n)$  space and OCCURRENCES( $v$ ) returns all occurrences in  $\mathcal{O}(o \cdot \log^{1+\epsilon} n)$  time [Ferragina and Manzini, 2000]. In practice, I sample text positions at rates between  $2^{-3}$  and  $2^{-5}$ . The rank dictionary  $I$  consumes  $n + o(n)$  extra space, independently of the sampling rate.

## 3.3 Algorithms

In this section, I give string matching algorithms that use the generic suffix trie traversal operations defined in section 2.4.2. Thus the following algorithms can be applied to all of the suffix trie implementations presented so far. I first consider a simple algorithm performing a top-down traversal bounded by depth. Then, I present algorithms for exact string matching and  $k$ -mismatches. I finally give, for the first time to the best of my knowledge, algorithms solving multiple variants of indexed exact string matching and  $k$ -mismatches. At the same time, I show the results of an experimental evaluation of all of these algorithms on various suffix trie implementations.

As data structures, I consider the fully-constructed lazy suffix tree (LST), the enhanced suffix array (ESA), the suffix array (SA), the  $q$ -gram index for  $q = 12$  (q-Gram), the FM-index with a two-levels DNA rank dictionary (FM-TL), the FM-index with a wavelet tree composed of two-levels binary rank dictionaries (FM-WT). As text, I take the *C. elegans* reference genome (WormBase WS195), i.e. a collection of 6 DNA strings of about 100 Mbp total length. As patterns, I use sequences extrapolated from an Illumina sequencing run (SRA/ENA id: SRR065390).

All experiments run on a desktop computer running Linux 3.10.11, equipped with one Intel® Core i7-4770K CPU @ 3.50 GHz, 32 GB RAM and a 2 TB HDD @ 7200 RPM. The plots show always *average* runtimes per pattern, both in single and multiple string matching variants. Moreover, they consider only traversal times, while they exclude the time to locate the patterns in the text by following leaf pointers, e.g. uncompressing SA values.

### 3.3.1 Construction

Table 3.1 shows construction times and memory consumption for all indices. The construction of the LST uses the quadratic-time in-memory *wotd* algorithm [Giegerich *et al.*, 1999], while the construction of all other indices relies on the linear-time external-memory

**Table 3.1:** Index construction times and memory footprints.

	LST	ESA	SA	$q$ -Gram	FM-WT	FM-TL
Time [s]	53.94	68.66	50.58	52.17	59.48	58.26
Memory [MB]	1404.02	1272.94	489.59	555.13	122.43	122.42

$DC7$  suffix array construction algorithm [Dementiev *et al.*, 2008]. In addition, the construction of the ESA uses the linear-time algorithms in [Kasai *et al.*, 2001; Abouelhoda *et al.*, 2004], while the FM-WT construction follows the linear-time algorithm proposed in [Grossi *et al.*, 2003]. For additional information on the SA and ESA construction algorithms and their runtimes, refer to [Weese, 2013].

On the chosen text, the practical runtime of the *wotd* algorithm (LST) is in line with that one of the  $DC7$  algorithm (SA). The construction of the  $q$ -gram index directory adds only 3% runtime over the plain SA construction, the FM-indices adds a 15–18% additional runtime, while the more involved ESA's LCP and child tables add 36% overhead. As expected, the FM-TL and FM-WT are the most compact data structures, while the LST followed by the ESA are the most space inefficient ones.

### 3.3.2 Depth-first traversal bounded by depth

Before turning to proper string matching algorithms, I present a simple algorithm that helps to comprehend subsequent backtracking algorithms. Algorithm 3.10 performs the top-down traversal of a suffix trie in depth-first order. The traversal is bounded, i.e. after reaching the nodes at depth  $d$  it stops going down and goes right instead.

The experimental evaluation shown in figure 3.9 provides a first glimpse on what are the practical performances of various suffix trie implementations. As expected, the WT FM-index is always slower than the TL FM-index. The  $q$ -gram index is never slower than the SA alone, however the contribution of the  $q$ -gram directory becomes insignificant for deep traversals.

---

#### Algorithm 3.10 DFS( $v, d$ )

---

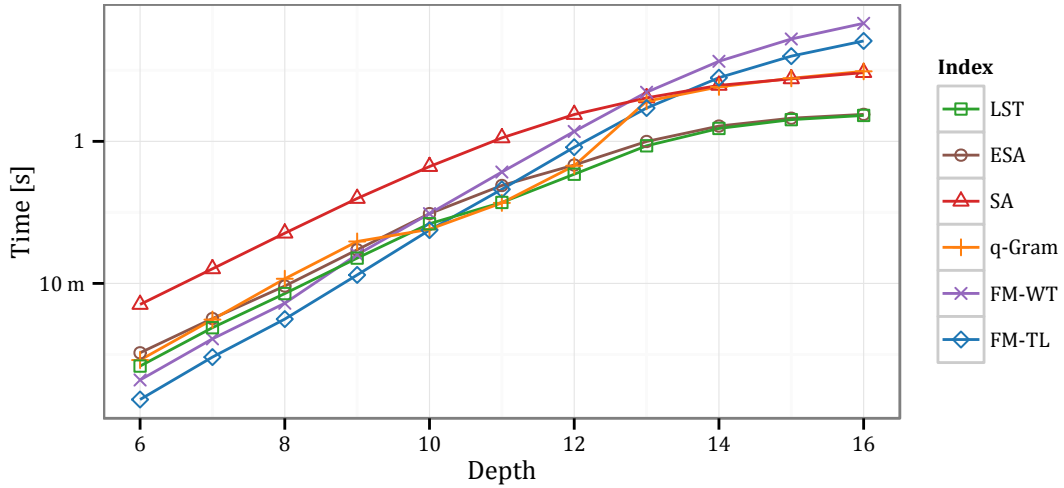
**Input**     $v$  : iterator of  $ST(S)$   
              $d$  : integer bounding the traversal depth

- 1: **if**  $d > 0$
- 2:     **if** GoDOWN( $v$ )
- 3:         **repeat**
- 4:             DFS( $v, d - 1$ )
- 5:         **until** GoRIGHT( $v$ )

---



**Figure 3.9:** Runtime of the bounded depth-first traversal of various suffix trie implementations.



Depth 12 marks the turning point, as the indices become sparse. The TL FM-index is the fastest index up to depth 10, while the  $q$ -gram index and the LST are the fastest at depths 10–12. Below depth 12, the tree indices (ESA and LST) become significantly faster than the trie indices (SA,  $q$ -gram and FM-indices). In particular, the FM-indices, which are based on backward search, become more than one order of magnitude slower than tree indices.

### 3.3.3 Exact string matching

I now give a simple algorithm performing exact string matching on a generic suffix trie. In the following, I assume the text  $T$  to be indexed by its suffix trie  $ST(T)$ . Algorithm 3.11 searches the pattern  $P$  by starting on the root node of  $ST(T)$  and following the path spelling the pattern. If the search ends up in a node  $v$ , then each leaf  $i$  below  $v$  points to a distinct suffix  $T_i$  such that  $T_{i...i+m}$  equals  $P$ . If `goDown` is implemented in constant time and `OCCURRENCES` in linear time, all occurrences of  $P$  into  $T$  are found in optimal time  $\mathcal{O}(m + o)$ , where  $m$  is the length of  $P$  and  $o$  its number of occurrences in  $T$ .

Figure 3.10 shows the results of the experimental evaluation of algorithm 3.11. On forward indices (LST, ESA, SA and  $q$ -gram index) the search time becomes practically constant for patterns of length above 15, i.e. when the tree becomes sparse. Conversely, on backward (FM) indices the practical search time stays linear in the pattern length.

The ESA and LST are never faster than the  $q$ -gram index despite their higher memory consumption. The SA alone is at least 20 % slower than the  $q$ -gram index, hence never competitive. In particular, the SA shows a runtime peak for patterns of length 10, due to the fact that binary search algorithms 3.1–3.2 converge more slowly for shorter patterns.

Concerning FM-indices, the WT variant is almost twice as slow as the TL variant, as

**Algorithm 3.11** EXACTSEARCH( $t, p$ )

---

**Input**      $t$  : iterator of  $ST(T)$   
               $p$  : iterator of  $P$

**Output**   list of all occurrences of  $P$  in  $T$

```

1: if ATEND( $p$ )
2:   report OCCURRENCES( $t$ )
3: else if GOWDOWN( $t$ , VALUE( $p$ ))
4:   GONEXT( $p$ )
5:   EXACTSEARCH( $t, p$ )

```

---

the WT-based rank dictionary performs twice the number of random memory accesses than the levels rank dictionary. Summing up, the TL FM-index is the fastest index to match exact patterns within length 30, while the  $q$ -gram index is the fastest for patterns above length 30.

### 3.3.4 Backtracking $k$ -mismatches

I now give an algorithm that solves  $k$ -mismatches by backtracking a generic suffix trie. The idea of backtracking a suffix tree has been first proposed in [Ukkonen, 1993]. Recently, various popular bioinformatics tools, e.g. Bowtie [Langmead *et al.*, 2009] and BWA [Li and Durbin, 2009], adopted variations of this method in conjunction with an FM-index. Yet, the idea dates back to more than twenty years ago.

Algorithm 3.12 performs a top-down traversal on the suffix trie  $ST(T)$ , spelling in-

**Figure 3.10:** Runtime of exact string matching on various suffix trie implementations.



**Algorithm 3.12** KMISMATCHES( $t, p, k$ )

---

**Input**      $t$  : iterator of  $ST(T)$   
               $p$  : iterator of  $P$   
               $k$  : integer bounding the number of mismatches

**Output**   list of all occurrences of  $P$  in  $T$

```

1: if  $k = 0$ 
2:   EXACTSEARCH( $t, p$ )
3: else
4:   if ATEND( $p$ )
5:     report OCCURRENCES( $t$ )
6:   else if GODOWN( $t$ )
7:     repeat
8:        $d \leftarrow \delta(\text{LABEL}(t), \text{VALUE}(p))$ 
9:       GONEXT( $p$ )
10:      KMISMATCHES( $t, p, k - d$ )
11:      GOPREVIOUS( $p$ )
12:    until GORIGHT( $t$ )

```

---

crementally all distinct substrings of  $T$ . While traversing each branch of the trie, this algorithm incrementally computes the distance between the query and the spelled string. If the computed distance exceeds  $k$ , the traversal backtracks and proceeds on the next branch. Conversely, if the pattern  $P$  is completely spelled and the traversal ends up in a node  $v$ , each leaf  $i$  below  $v$  points to a distinct suffix  $T_i$  such that  $d_H(T_{i..i+m}, P) \leq k$ .

**Figure 3.11:** Runtime of 1-mismatch search on various suffix trie implementations.



Figure 3.11 shows the results of the experimental evaluation of algorithm 3.12 for  $k = 1$ . The TL FM-index is always faster than any other index: for instance, on patterns of length 30, the SA is 3 times slower; even the  $q$ -gram index is 50 % slower than the TL FM-index. Despite their tree structure, ESA and LST are always slower than  $q$ -gram and FM-indices.

On the TL FM-index, 1-approximate matching of patterns of length 30 is 16 times slower than exact matching: on average, 1-approximate matching spends 21 microseconds ( $\mu s$ ), while exact matching takes  $1.3 \mu s$ . On the ESA, 1-approximate matching shows a slow-down of 24 times: 1-approximate matching spends  $40.8 \mu s$ , compared to the  $1.7 \mu s$  for exact matching.

### 3.3.5 Multiple exact string matching

Before turning to multiple  $k$ -mismatches, I describe a simpler algorithm for multiple exact string matching. In addition to the text  $T$ , multiple exact string matching provides a collection of patterns  $\mathbb{P}$ . Hence, in addition to the suffix trie  $ST(T)$ , algorithm 3.13 considers the trie  $RT(\mathbb{P})$ . Algorithm 3.13 matches simultaneously in  $ST(T)$  all patterns indexed in  $RT(\mathbb{P})$ . The traversal performed by algorithm 3.13 visits pairs of nodes in  $ST(T) \times RT(\mathbb{P})$  whose entering edges have the same label. Such traversal implicitly *intersects* the two tries. However, algorithm 3.13 is not symmetric:  $ST(T)$  and  $RT(\mathbb{P})$  cannot be interchanged. The traversal stops whenever it reaches a leaf node in  $RT(\mathbb{P})$  and reports the occurrences pointed by all the leaves beneath the current node in  $ST(T)$ .

The experimental evaluation compares algorithm 3.13 (Multiple) with algorithm 3.11 processing patterns in random order (Single) and in lexicographic order (Sorted). Figure 3.12 shows the results. These three methods ran on 10 M patterns of length 30: runtimes shown in figure 3.12 (histogram Single) correspond to runtimes shown in figure 3.10 (plots at pattern length 15).

---

#### Algorithm 3.13 MULTIPLEEXACTSEARCH( $t, p$ )

---

**Input**      $t$  : iterator of  $ST(T)$   
               $p$  : iterator of  $RT(\mathbb{P})$

**Output**   list of all occurrences of  $\mathbb{P}$  in  $T$

```

1: if ISLEAF( $p$ )
2:   report OCCURRENCES( $t$ )  $\times$  OCCURRENCES( $p$ )
3: else
4:   GOWDOWN( $p$ )
5:   repeat
6:     if GOWDOWN( $t$ , LABEL( $p$ ))
7:       MULTIPLEEXACTSEARCH( $t, p$ )
8:       GOWUP( $t$ )
9:   until GOWRIGHT( $p$ )

```

---

Figure 3.12 shows that a simple lexicographical sort of the patterns (histogram Sorted) speeds up algorithm 3.11 on the SA and ESA by a factor of 2. The same trick does not yield a significant speed-up on FM-indices nor on the  $q$ -gram index, as the  $q$ -gram directory already provides a cache local access pattern.

Algorithm 3.13 (histogram Multiple) further reduces the traversal time. Nonetheless, its runtime is dominated by the additional preprocessing time paid to construct the trie of the patterns. This algorithm becomes more useful as a primitive within the multiple  $k$ -mismatches algorithm, presented in the following section.

### 3.3.6 Multiple $k$ -mismatches

Algorithm 3.14 is the straightforward generalization of algorithm 3.13 to  $k$ -mismatches. The algorithm receives a collection of patterns  $\mathbb{P}$  and performs backtracking on  $\mathbb{T}$  as in algorithm 3.12, this time using the associated trie  $\mathbb{P}$ .

The experimental evaluation compares algorithm 3.14 (Multiple) with algorithm 3.12 processing patterns in random order (Single) and in lexicographic order (Sorted). All three methods ran on 10 M patterns of length 30, with  $k$  fixed to 1. Thus, runtimes shown in figure 3.13 (histogram Single) correspond to runtimes shown in figure 3.11 (plots at pattern length 30).

Algorithm 3.12 on lexicographically sorted patterns (histogram Sorted) is faster by a factor of 2 or more, on all indices. The time to sort the patterns becomes insignificant

---

#### Algorithm 3.14 MULTIPLEKMISMATCHES( $t, p, k$ )

---

**Input**      $t$  : iterator of  $ST(T)$   
               $p$  : iterator of  $RT(\mathbb{P})$   
               $k$  : integer bounding the number of mismatches

**Output**   list of all occurrences of  $\mathbb{P}$  in  $T$

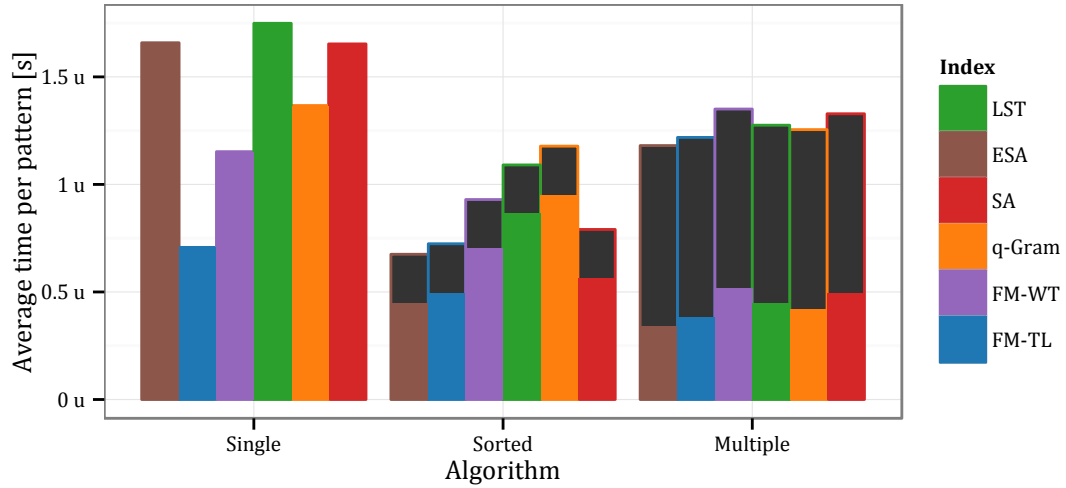
```

1: if  $k = 0$ 
2:   MULTIPLEEXACTSEARCH( $t, p$ )
3: else
4:   if ISLEAF( $p$ )
5:     report OCCURRENCES( $t$ )  $\times$  OCCURRENCES( $p$ )
6:   else if GODOWN( $t$ )
7:     repeat
8:       GODOWN( $p$ )
9:       repeat
10:         $d \leftarrow \delta(\text{LABEL}(t), \text{LABEL}(p))$ 
11:        MULTIPLEKMISMATCHES( $t, p, k - d$ )
12:      until GORIGHT( $p$ )
13:    GOUP( $p$ )
14:  until GORIGHT( $t$ )

```

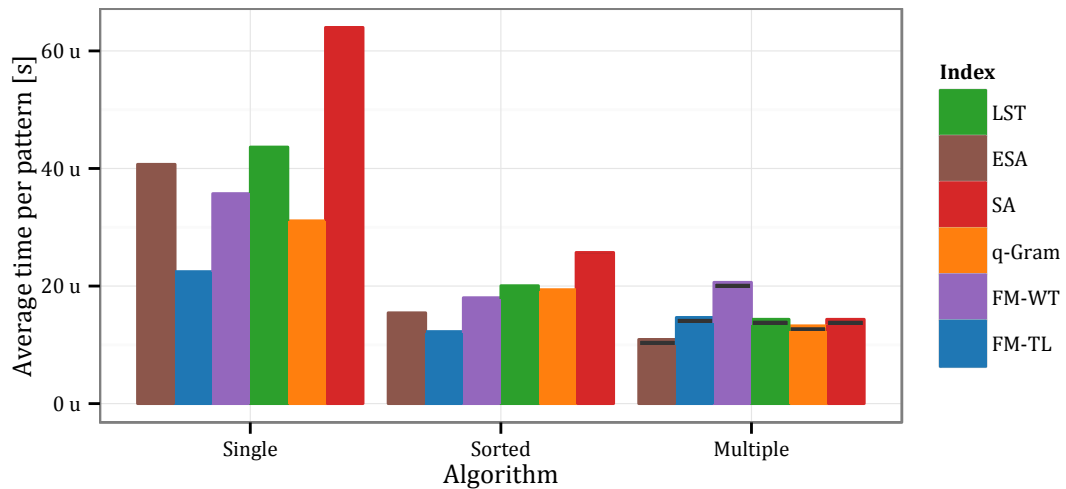
---

**Figure 3.12:** Runtime of multiple exact string matching on various suffix trie implementations. Pattern length is fixed to 15. Preprocessing times are shown in black.



compared to the traversal time. Algorithm 3.13 (histogram Multiple) reduces traversal time by a factor of 4–5 on LST, ESA, and SA. Thus, the time to construct the trie of the patterns is easily justified. In practice, algorithm 3.13 fills the gap between the runtime of the SA and the  $q$ -gram index. Surprisingly, algorithm 3.13 increases traversal time on FM-indices.

**Figure 3.13:** Runtime of multiple 1-mismatch on various suffix trie implementations. Pattern length is fixed to 30. Preprocessing times are shown in black.



**Figure 3.14:** Speed-up of multiple 1-mismatch by number of patterns on the SA. Pattern length is fixed to 30. Traversal times without preprocessing are shown by dashed lines.



This algorithm works according to a *cache-friendly* memory access pattern, which holds for forward search but not for backward search. Using forward search, the traversal of a suffix trie becomes less expensive as it proceeds towards bottom nodes. Indeed,

**Figure 3.15:** Speed-up of multiple 1-mismatch by number of patterns on the TL FM-Index. Pattern length is fixed to 30. Traversal times without preprocessing are shown by dashed lines.



traversal towards a child node involves the computation of a subinterval of the current suffix array interval; such computation accesses memory locations within the current interval, having good chances to be in the cache. Conversely, using backward search, the traversal becomes more expensive as it proceeds deeper in the trie; traversal downwards involves the computation of intervals outside of the current one, unlikely to be in the cache as they are accessed less often than top intervals. Multiple backtracking factorizes the traversal of top nodes, thus it pays off with forward search rather than with backward search.

Figure 3.14 shows the average runtime of the three approaches on the SA by varying the number of patterns. While the average runtime of the single method is constant, both multiple methods clearly benefit from receiving a higher number of patterns. In particular, method Multiple is constantly faster than Sorted, and the runtime gap increases with the number of patterns. The speed-up of multiple methods slowly decreases, though there is still some space of improvement with more than 10 M patterns.

Figure 3.15 presents the same evaluation of figure 3.14, but for the TL FM-index. Multiple methods exhibit again decreasing average runtimes by number of patterns. However, here method Sorted is constantly faster than Multiple, but the runtime gap decreases with the number of patterns. Moreover, the speed-up of both multiple methods slowly increases with the number of patterns instead of decreasing.



In this chapter, I present various filtering methods for approximate string matching. I consider two classes of filtering methods: those based on *seeds* and those based on *q-grams*. Filters of the former class partition the pattern into *non-overlapping* factors called seeds, while filters of the latter class consider all *overlapping* substrings of the pattern having length  $q$ , the so-called *q-grams*. Both classes include various combinatorial filtering methods of increasing specificity and complexity; all methods provide filtration schemes with guarantees on filtration sensitivity.

I present two seed filtering methods: *exact seeds* [Baeza-Yates and Perleberg, 1992] and *approximate seeds* [Myers, 1994; Navarro and Baeza-Yates, 2000]. Exact seeds partition the pattern in  $k + 1$  non-overlapping seeds to be searched exactly in the text. Approximate seeds increase filtration specificity by factorizing the pattern in less than  $k + 1$  non-overlapping seeds to be searched within a distance threshold smaller than  $k$ .

I present the following *q-gram* filtering methods: *contiguous q-grams* [Jokinen and Ukkonen, 1991], *gapped q-grams* [Burkhardt and Kärkkäinen, 2001] and *multiple gapped q-grams* (also called *q-gram families*) [Kucherov *et al.*, 2005]. Contiguous *q-grams* rely on counting arguments to filter out text regions containing less than a given threshold of *q-gram* occurrences. Gapped *q-grams* introduce *don't care positions* to lower the correlation between occurrences of consecutive *q-grams*. Multiple gapped *q-grams* adopt multiple patterns of don't care positions to further increase filtration specificity.

It will become clear through this chapter that seed filters are more practical, flexible, straightforward to design and implement than *q-gram* filters. All seed filters provide full-sensitive filtration schemes for the  $k$ -differences problem, while (multiple) gapped *q-grams* only for  $k$ -mismatches. The design of highly specific yet full-sensitive filtration schemes for *q-gram* filters is combinatorially hard, while it is quite straightforward for seed filters. Also implementation-wise, *q-gram* filters are more involved than seeds filter. In fact, seed filters lend themselves well to both online and offline variants of the problem, while *q-gram* filters are better suited for the online variant. Finally, the experimental evaluation shows that seed filters outperform *q-gram* filters for most practical inputs. For these reasons, I design the applications of chapters 6 and 7 around seed filtering methods.

## 4.1 Exact seeds

Filtration with exact seeds is one of the naïvest filtering methods for approximate string matching. I first explain the underlying combinatorial principle, then I discuss implementation details and lastly give some insights on the efficiency of this method.

### 4.1.1 Principle

I consider the case of two arbitrary strings  $X, Y$  within edit distance  $k$ . The generalization to  $k$ -differences is straightforward.

**Lemma 4.1 (Exact seeds).** *Let  $X, Y$  be two strings s.t.  $d_E(X, Y) = k$ . If  $Y$  is partitioned w.l.o.g. into  $k + 1$  non-overlapping seeds, then at least one seed occurs as a factor of  $X$  [Baeza-Yates and Perleberg, 1992].*

It is immediate to see that any edit distance error can cover at most one seed. Therefore, at least one seed of  $Y$  will not be covered by any seed and hence occur as a factor of  $X$ . Figure 4.1 shows an example.

This filtering method reduces the approximate search into multiple smaller exact searches. It solves  $k$ -differences by partitioning the pattern into  $k + 1$  seeds, searching all seeds in the text, and verifying a text window around each occurrence of any seed in the text. As lemma 4.1 is valid for *any substring* of the text within distance  $k$  from the pattern, this method finds all approximate occurrences of the pattern in the text.

### 4.1.2 Efficiency

The efficiency of this method strongly depends on the number of verifications. It is straightforward to derive the expected number of verifications under the assumption of the text being generated according to the uniform Bernoulli model. I introduce the random variable  $C$ , counting the number of occurrences of a word in a text. The emission probability of any symbol in  $\Sigma$  is  $p = 1/\sigma$  and under i.i.d. assumptions the emission (and occurrence) probability of any word of length  $q$  is simply

$$\Pr[C > 0] = \frac{1}{\sigma^q} \quad (4.1)$$

**Figure 4.1:** Filtration with 6 exact seeds solves 5-differences. In the illustration, pattern  $P$  occurs in text  $T$  at edit distance 5. The seed in grey is not covered by any error and thus preserved.



thus the expected number of occurrences of a seed of length  $q$  in a text of length  $n$  is

$$E[C] = \sum_{i=0}^{n-q} \Pr[C > 0] = \frac{n - q + 1}{\sigma^q} \leq \frac{n}{\sigma^q}. \quad (4.2)$$

Lemma 4.1 requires to partition the pattern into  $k + 1$  seeds but leaves the freedom to choose their length. This leads to the problem of finding an optimal pattern partitioning that minimizes the expected number of verifications. I fix<sup>1</sup> the length of all seeds to be

$$q = \left\lfloor \frac{m}{k + 1} \right\rfloor \quad (4.3)$$

to minimize the expected number of occurrences of any seed. Under these conditions, the expected number of verifications produced by filtration with exact seeds is

$$E[V] = E[C] \cdot (k + 1) < \frac{n(k + 1)}{\sigma^q}. \quad (4.4)$$

Nonetheless, inputs of practical interest like genomes and natural texts do not fit well the uniform Bernoulli model. On those texts, uniform seed length often leads to suboptimal filtration.

## 4.2 Approximate seeds

The simple analysis of section 4.1.2 shows that filtration specificity is strongly correlated to the seed length. Therefore, the crux of designing a stronger filter lies into increasing the seed length while respecting full-sensitivity constraints. Myers [1994], subsequently followed by Navarro and Baeza-Yates [2000], proposed *approximate seeds* as a practical and effective generalization of exact seeds that yield stronger filters for  $k$ -differences. The key idea of filtration with approximate seeds is to reduce the approximate search into smaller approximate searches, as opposed to filtration with exact seeds that reduces the approximate search into smaller exact searches.

### 4.2.1 Principle

Again, I start by considering two arbitrary strings  $X, Y$  within edit distance  $k$ . The result then holds for any substring of the text within distance  $k$  from the pattern.

**Lemma 4.2 (Approximate seeds).** *Let  $X, Y$  be two strings s.t.  $d_E(X, Y) = k$ . If  $Y$  is partitioned w.l.o.g. into  $s$  non-overlapping seeds s.t.  $1 \leq s \leq k + 1$ , then at least one seed occurs as a factor of  $X$  within distance  $\lfloor k/s \rfloor$  [Myers, 1994; Navarro and Baeza-Yates, 2000].*

To prove full-sensitivity it suffices to see that, if none of the seeds occurs within its assigned distance, the total distance must be greater than  $s \cdot \lfloor k/s \rfloor = k$ . Figure 4.2 illustrates.

<sup>1</sup> For simplicity I ignore that some seed could have length  $\lfloor \frac{m}{k+1} \rfloor$ .

**Figure 4.2:** Filtration with approximate seeds. A filtration scheme with thresholds  $T = (1, 1, 1)$  solves 5-difference. In the illustration, pattern  $P$  occurs in text  $T$  at edit distance 5. The seed in grey is covered only by one error and thus preserved.



### 4.2.2 Filtration schemes

Approximate seeds provide filtration schemes of variable specificity. The fastest but weakest filtration scheme is given by  $s = k + 1$ , while the most specific filtration is obtained for  $s = 1$  i.e. perfect filtration scheme without any verification step. Alternatively, filtration specificity is controlled by acting on the minimum seed length  $q$ . Fixing  $q$  yields  $s = \lfloor m/q \rfloor$ , or vice versa, fixing the number of seeds  $s$  gives  $q = \lfloor m/s \rfloor$ . Filtration specificity is expected to increase with seed length.

Lemma 4.2 assigns the same distance threshold to all seeds, yet this is not obligatory. Hence, I give a more general definition of *filtration scheme* for approximate seeds.

**Definition 4.1 (Filtration scheme).** A seeds filtration scheme is an integer vector  $T = (t_0, \dots, t_{s-1})$ , where integer  $t_i \in \mathbb{N}_0$  represents the threshold assigned to the  $i$ -th seed.

**Lemma 4.3.** Any filtration scheme  $T = (t_0, \dots, t_{s-1})$  s.t.

$$s + \sum_{i=0}^{s-1} t_i > k \quad (4.5)$$

is full-sensitive for  $k$ -differences (and  $k$ -mismatches).

**Example 4.1.** The filtration schemes  $(0, 0, 0, 0, 0)$ ,  $(1, 1, 0)$ ,  $(2, 1)$ ,  $(4)$  are full-sensitive for 4-differences. For instance, given a pattern of length  $m = 100$ , according to equation 4.3,  $q$  is respectively 20, 33, 50, 100.

How to choose a *good* filtration scheme in practice? Both Myers [1994] and Navarro and Baeza-Yates [2000] carried out involved analysis to estimate the optimal parameterization. Navarro and Baeza-Yates find out that a number of seeds  $s = \Theta(m/\log_\sigma n)$  yields an overall time complexity sublinear for an error rate  $\epsilon < 1 - e/\sqrt{\sigma}$ . Myers reports an analogous sublinear time when  $q = \Theta(\log_\sigma n)$  is the seed length. Yet, these results do not necessarily translate into optimal filtration schemes in practice. The parameterization depends on the full-text index, the verification algorithm, the statistical properties of the text. Missing the optimal number of seeds by one often results in a runtime penalty of an order of magnitude.

Having established the number of seeds, or their length, thresholds have to be assigned. Lemma 4.3 allows to assign arbitrary distance thresholds. In practice, it is convenient to distribute distance thresholds evenly, as seeds with the highest threshold dominate the overall filtration time. The most strict threshold assignment is to give distance  $\lfloor k/s \rfloor$  to  $(k \bmod s) + 1$  seeds and distance  $\lfloor k/s \rfloor - 1$  to the remaining seeds [Siragusa *et al.*, 2013a].

## 4.3 Contiguous $q$ -grams

$q$ -Gram filters rely on counting arguments to filter out text regions containing less than a given threshold of  $q$ -gram occurrences. The first  $q$ -gram counting filter for approximate string matching has been proposed in [Jokinen and Ukkonen, 1991]. Filters for more general alignment problems have been proposed and implemented in *QUASAR* [Burkhardt *et al.*, 1999], *SWIFT* [Rasmussen *et al.*, 2006] and *STELLAR* [Kehr *et al.*, 2011].

### 4.3.1 Principle

The counting argument of contiguous  $q$ -gram filters is based on  *$q$ -gram similarity*: the number of substrings of length  $q$  common to two given strings. The following lemma relates  $q$ -gram similarity to edit distance, by giving a *lower bound* on the  $q$ -gram similarity of any two strings  $X, Y$  within edit distance  $k$ . As for seed filters, this result then easily translates to  $k$ -differences.

**Lemma 4.4 (The  $q$ -gram lemma).** *Let  $X, Y$  be two strings s.t.  $d_E(X, Y) = k$  and assume w.l.o.g.  $|X| \leq |Y|$  and  $|X| = m$ . Then  $X$  and  $Y$  have  $q$ -gram similarity  $\tau_q(m, k) \geq m - q + 1 - kq$  [Jokinen and Ukkonen, 1991].*

The first part of the threshold function  $\tau_q$  counts the number of  $q$ -grams of  $X$  (i.e.  $m - q + 1$ ), while the second part counts how many  $q$ -grams can be covered by  $k$  errors (i.e. at most  $q$  per error, hence  $kq$  in total). The position of errors in the transcript solely determines which  $q$ -gram occurrences are affected or preserved. The  $q$ -gram lemma considers one *worst case* positioning of the errors that *minimizes* the threshold. Figure 4.3 exemplifies.

### 4.3.2 Filtration schemes

I denote by a pair  $(q, t)$  the filtration scheme counting  $q$ -grams with threshold  $t$ . According to lemma 4.4, if  $t = \tau_q(m, k) \geq 1$ , then  $(q, t)$  is full-sensitive for any  $k$ -differences instance where  $|P| = m$ . In this case, I say that  $(q, t)$  *solves* instance  $(m, k)$ .

The following question arises: which is the longest  $q$ -gram solving instance  $(m, k)$ ? In order to satisfy lemma 4.4, the  $q$ -gram threshold must be greater than zero, i.e. it must hold  $\tau_q(m, k) \geq 1$ . Thus, by substituting  $\tau_q(m, k)$ , it follows that the  $q$ -gram length must be  $q \leq \lfloor \frac{m}{k+1} \rfloor$ , analogously to seed filters (see equation 4.3).

However, the longest  $q$ -gram does not yield always the most specific filtration scheme. For instance, a threshold of 1 completely discards the counting argument of lemma 4.4

**Figure 4.3:** Filtration with contiguous  $q$ -grams. A filtration scheme  $(q, t) = (4, 2)$  solves the instance  $(25, 5)$ . In the illustration, pattern  $P$  of length 25 occurs in text  $T$  at edit distance 5. In this worst-case positioning of the errors, the two grey  $q$ -grams are preserved.



and makes filtration very unspecific in practice. Hence, on certain  $(m, k)$  instances, filtration schemes with non-optimal  $q$ -gram length yield more specific filtration. Example 4.2 shows alternative filtration schemes solving a given  $(m, k)$  instance.

**Example 4.2.** The following  $(q, t)$  filtration schemes solve  $(100, 4)$ -differences:  $(20, 1)$ ,  $(19, 6)$ ,  $(18, 11)$ .

### 4.3.3 Bucketing

Filtration with  $q$ -grams requires *bucketing* the text in windows, in order to apply the counting argument of lemma 4.4. Buckets are obtained by subdividing the implicit DP matrix in parallelograms and projecting them on the text. Figure 4.4 illustrates this concept: any approximate occurrence of the pattern in the text spans at most  $k$  diagonals and is thus enclosed inside a parallelogram of width  $k + 1$  [Rasmussen *et al.*, 2006]. Hence the projection of any text bucket has length  $2k + 1$  and any occurrence has length between  $m - k$  and  $m + k$ . The implementations described in [Rasmussen *et al.*, 2006; Kehr *et al.*, 2011; Weese *et al.*, 2009] use more efficient bucketing strategies with larger, overlapping parallelograms.

This method lends itself to work in a multiple online fashion rather than offline. The filtration stage scans the text and counts how many  $q$ -grams of each pattern fall into each parallelogram bucket. As long as the filter scans the text, it remembers only the buckets that span the patterns' lengths. The verification stage then verifies only those parallelograms exceeding threshold  $t$ . Conversely, the program QUASAR [Burkhardt *et al.*, 1999] uses a  $q$ -gram index of the text to speed up the filtration phase. Such implementation requires more memory, as it must bucket the whole text and keep the text index in memory.

## 4.4 Gapped $q$ -grams

Califano and Rigoutsos [1993] first introduced *gapped  $q$ -grams* in sequence analysis. Since then, a surprisingly high number of research papers have been published on this

**Figure 4.4:** Parallelogram buckets. Picture from [Weese et al., 2009].



topic (see [Brown, 2008] for a survey). Almost all works focus on lossy filtration for homology search, rather than full-sensitive filtration for approximate string matching. Here, I consider gapped  $q$ -grams only in the context of full-sensitive filtration for  $k$ -mismatches. This case has been first considered by Burkhardt and Kärkkäinen [2001].

#### 4.4.1 Principle

Gapped  $q$ -grams introduce fixed *don't care positions* where text and pattern characters are ignored. A comparison between figures 4.3 and 4.5 illustrates the advantage of such don't care positions. While in figure 4.3 an error in a transcript affects (at most) a cluster of  $q$  consecutive  $q$ -gram occurrences, in figure 4.5 a mismatch does not affect those gapped  $q$ -gram occurrences that ignore its position. This fact relaxes the full-sensitivity threshold and opens the door to more specific filtration schemes for  $k$ -mismatches.

The counting argument of gapped  $q$ -grams generalizes  $q$ -gram similarity (section 4.3) from substrings to *subsequences*, i.e. from contiguous to *non-contiguous* sequences of symbols. Filtration with gapped  $q$ -grams indeed counts the number of subsequences of length  $q$  common to two strings. An additional set  $Q$  determines which symbols are taken in the subsequences. The formal definition of gapped  $q$ -gram follows.

**Definition 4.2 ( $Q$ -gram).** A gapped  $q$ -gram (abbreviated as  $Q$ -gram) is a finite sequence  $Q$  of natural numbers starting with the unit element, i.e.  $Q \subset \mathbb{N}_0$  and  $0 \in Q$ . The length  $|Q|$  is called the *weight* of  $Q$  and denoted as  $w(Q)$ . The maximum element of  $Q$  is named *span* and indicated by  $s(Q)$ .

**Figure 4.5:** Filtration with gapped  $q$ -grams. A filtration scheme  $(Q, t) = (\{0, 2, 3, 4\}, 3)$  solves instance  $(25, 5)$ . In the illustration, pattern  $P$  of length 25 occurs in text  $T$  at edit distance 5. The three  $q$ -grams in grey are not covered by any mismatch.



Figure 4.5 shows an example of  $Q$ -gram. As in lemma 4.4, the threshold for  $Q$ -grams still depends on the worst-case positioning of the errors in the transcript, which in turn depends on parameters  $(m, k)$ . Thus, I still consider filtration schemes  $(Q, t)$  solving a  $(m, k)$  instance. However, contrary to contiguous  $q$ -grams, function  $\tau_Q$  now gives only a lower bound to the full-sensitivity threshold.

#### 4.4.2 Filtration schemes

Which is the most specific filtration scheme  $(Q, t)$  solving a given instance  $(m, k)$ ? This question turns out to be surprisingly hard to answer. As discussed in section 4.3.2, the most specific filtration schemes for contiguous  $q$ -grams are easily found. The choice falls on a few values of  $q$  that are close to the maximum and do not yield a threshold too close to 1. Conversely, such choice is non-trivial for  $Q$ -grams, as the search space of  $Q$  is exponentially large in the span  $s(Q)$  and a full-sensitivity threshold for arbitrary  $Q$ -grams is hard to compute. In addition, it is not easy to determine which filtration scheme is the most specific one in a set of full-sensitive candidates.

Given a  $Q$ -gram, I consider the following problems:

- FULL SENSITIVITY Does filtration scheme  $(Q, 1)$  solve instance  $(m, k)$ ?
- OPTIMAL THRESHOLD Which is the optimal threshold  $t^*$  s.t.  $(Q, t^*)$  solves  $(m, k)$ ?
- SPECIFICITY Which is the expected specificity of scheme  $(Q, t)$ ?

Nicolas and Rivals [2005] considered the decision problem FULL SENSITIVITY associated to OPTIMAL THRESHOLD. FULL SENSITIVITY is easy for contiguous  $q$ -grams, i.e. the answer is no iff  $\tau_q(m, k) = 0$ . Nicolas and Rivals show, by performing an indirect reduction from EXACT COVER BY 3-SETS, that FULL SENSITIVITY is *strongly* NP-complete for arbitrary  $Q$ -grams. Strong NP-completeness implies that no *fully polynomial-time approximation scheme* (FPTAS) nor any *pseudo-polynomial* algorithm for FULL SENSITIVITY exist, under the assumption that  $P \neq NP$ .

Burkhardt and Kärkkäinen [2001] first considered the optimization problem OPTIMAL THRESHOLD. They give a DP algorithm solving OPTIMAL THRESHOLD in time  $\mathcal{O}(m \cdot k \cdot$



$2^{s(Q)}$  for any  $Q$ -gram. Subsequently, they use their DP algorithm to explore the search space of full-sensitive  $Q$ -grams for some specific instances of  $(m, k)$ . They reduce the search space using a *branch-and-bound* algorithm based on the observation that:

$$Q_1 \subseteq Q_2 \Leftrightarrow \tau_{Q_1}(m, k) \geq \tau_{Q_2}(m, k). \quad (4.6)$$

Finally, Burkhardt and Kärkkäinen guess the most specific filtration scheme  $(Q, t)$  by computing, again via branch-and-bound, the minimum number of characters covered by  $t$  distinct  $Q$ -gram occurrences (MINIMUM COVERAGE criterion).

I propose alternative algorithms for the design of  $Q$ -gram filtration schemes. I give an *integer linear program* (ILP) that solves exactly OPTIMAL THRESHOLD and is orders of magnitude quicker than the DP algorithm from [Burkhardt and Kärkkäinen, 2001] (see table 4.2). Afterwards, I give a simple polynomial time *approximation algorithm* for OPTIMAL THRESHOLD usable as an additional search space filter to quickly discard lossy  $Q$ -grams and eventually providing good ILP starting points. Finally, I propose a randomized algorithm to estimate in polynomial time the specificity of a filtration scheme, instead of using the MINIMUM COVERAGE criterion.

Having said that, I have no practical interest into designing  $Q$ -gram filtration schemes. On the one hand, almost all HTS applications require methods that solve  $k$ -differences in order to detect indels. On the other hand, the experimental evaluation of section 4.5 shows that  $Q$ -gram filters are not better than approximate seeds filters. Therefore, I do not need these algorithms in the applications presented in chapters 6 and 7.

### 4.4.3 Full sensitivity

I start by modeling the decision problem FULL SENSITIVITY before turning to the optimization problem OPTIMAL THRESHOLD. I consider any Hamming distance transcript over  $\{R, M\}$  as an  $m$ -dimensional vector  $X = (x_0, \dots, x_{m-1})$  over  $\mathbb{B} = \{0, 1\}$ . I denote by  $\mathbb{B}_k^m \subset \mathbb{B}^m$  the set containing all transcripts  $X$  of length  $m$  whose Hamming distance is  $k$ .

**Definition 4.3.** A  $Q$ -gram *occurs* at position  $i$  in transcript  $X$  iff  $\forall j \in Q \ x_{i+j} = 1$ .

**Definition 4.4.** A filtration scheme  $(Q, t)$  *detects* transcript  $X$  iff the  $Q$ -gram occurs at least  $t$  times in  $X$ .

I introduce a *boolean function* to characterize the set of transcripts detected by a filtration scheme of the form  $(Q, 1)$ . Let  $T_Q^m : \mathbb{B}^m \rightarrow \mathbb{B}$  denote the boolean function such that  $T_Q^m(X)$  is true iff the  $Q$ -gram occurs in a transcript  $X$  of length  $m$ . I define such boolean function as the disjunction

$$T_Q^m(X) = \bigvee_{i=0}^{m-s(Q)} \bigwedge_{j \in Q} x_{i+j} \quad (4.7)$$

where each *clause* of  $T_Q^m$  represents a single possible occurrence of  $Q$  in  $X$ . According to definition 4.4, filtration scheme  $(Q, t)$  detects  $X$  iff  $X$  satisfies at least  $t$  clauses of  $T_Q^m$ .

**Problem 4.1** (FULL SENSITIVITY)

**Instance** A  $Q$ -gram, an  $(m, k)$  instance.

**Question**  $\exists X \in \mathbb{B}_k^m$  s.t.  $T_Q^m(X) = 0$ ?

#### 4.4.4 Optimal threshold

I now consider the *pseudo-boolean function*, counterpart of function 4.7, that associates a filtration threshold to any transcript. Let the function  $t_Q^m : \mathbb{B}^m \rightarrow \mathbb{N}_0$  be the boolean function  $T_Q^m$  acting on  $\mathbb{N}_0$ . Here,  $t_Q^m(X)$  counts how many times a  $Q$ -gram occurs in a transcript  $X$  of length  $m$ . I define such pseudo-boolean function as

$$t_Q^m(X) = \sum_{i=0}^{m-s(Q)} \prod_{j \in Q} x_{i+j} \quad (4.8)$$

The formal definition of OPTIMAL THRESHOLD follows.

**Problem 4.2 (OPTIMAL THRESHOLD)**

**Instance** A  $Q$ -gram, an  $(m, k)$  instance.

**Solution**  $\min t_Q^m(X)$  subject to  $X \in \mathbb{B}_k^m$

**Exact solution**

I propose the following ILP to solve OPTIMAL THRESHOLD.

$$\begin{aligned} \min \quad & \sum T \\ \text{subject to} \quad & \end{aligned} \quad (4.9)$$

$$T \in \mathbb{B}^{m-s(Q)+1} \quad (4.10)$$

$$X \in \mathbb{B}_k^m \quad (4.11)$$

$$t_i \geq x_{i+j} \quad \forall 0 \leq i \leq m - s(Q), j \in Q \quad (4.12)$$

Vector  $X$  represents any Hamming transcript: each  $x_j$  is subject to an integer linear constrain s.t. the Hamming distance of  $X$  is within  $k$  (constraint 4.11). Vector  $T$  represents function  $t_Q^m$ : each  $t_i$  indicates the truthfulness of the  $i$ -th clause of  $T_Q^m$  (constraint 4.10) and its sum (objective function 4.9) the value of  $t_Q^m(X)$ . The inequalities of the form  $t_i \geq x_{i+j}$  (constraint 4.12) bind the satisfiability of each clause  $t_i$  to its associated transcript values  $x_j$ . The solution  $T^*$  to the above ILP provides the optimal threshold  $t = \sum T^*$  for a  $Q$ -gram on instance  $(m, k)$ . Table 4.2 shows the performance of ILP 4.9 on some problem instances.

**Approximate solution**

Pseudo-boolean function  $t_Q^m$  has two important properties: it is *monotone non-decreasing* and *supermodular*. In particular, supermodularity is the discrete analog of concavity. These properties allow to compute approximate solutions to constrained optimization problems in polynomial time. A pseudo-boolean function  $f : \mathbb{B}^m \rightarrow \mathbb{Z}$  is monotone non-decreasing iff

$$f(Y) \leq f(Y \vee Z) \quad \forall Y, Z \in \mathbb{B}^m \quad (4.13)$$

---

**Algorithm 4.1** APPROXIMATETHRESHOLD( $Q, m, k$ )

---

**Input**      $Q$  :  $Q$ -gram  
               $m$  : integer denoting the pattern length  
               $k$  : integer bounding the number of mismatches

**Output**   integer indicating the optimal threshold

```

1:  $X \leftarrow \mathbf{1}^m$ 
2: while  $k > 0$  do
3:    $i \leftarrow \arg \max_{j=1}^m \frac{\partial t_Q^m(X)}{\partial x_j}$ 
4:    $x_i \leftarrow 0$ 
5:    $k \leftarrow k - 1$ 
6: return  $t_Q^m(X)$ 

```

---

and it is *supermodular* iff

$$f(Y) + f(Z) \leq f(Y \vee Z) + f(Y \wedge Z) \quad \forall Y, Z \in \mathbb{B}^m. \quad (4.14)$$

It is immediate to see that function  $t_Q^m$  satisfies all the above inequalities, as it does not contain any term in negative form by definition (see equation 4.8). OPTIMAL THRESHOLD thus involves the minimization of a monotone supermodular function subject to linear constraints.

There is a well-known *gradient ascent* algorithm for the maximization of a monotone submodular function subject to linear constraints with an APX-ratio of  $1 - 1/e$  [Nemhauser *et al.*, 1978], i.e. the ratio between the value of the greedy approximation and the value of the optimal solution is within 0.632. Algorithm 4.1 computes an approximate solution to OPTIMAL THRESHOLD using the same principle. Nonetheless, its solution is not guaranteed to be within a constant ratio from the optimal one. Indeed, algorithm 4.1 minimizes function  $t_Q^m$  which is monotone non-decreasing supermodular and has its unconstrained minimum at 0.

In practice, the threshold found by algorithm 4.1 turns out to be a tight upper bound for OPTIMAL THRESHOLD (see table 4.2). Such threshold cannot be used in a filtration scheme guaranteed to be full-sensitive: if the threshold is non-optimal, the filtration scheme is lossy. Instead, algorithm 4.1 might find its use in the design of  $Q$ -grams, as a search space filter to quickly discard lossy  $Q$ -grams.

#### 4.4.5 Specificity

Assuming the text to be generated according to the uniform Bernoulli model, the expected specificity of any filtration scheme is proportional to the number of transcripts detected. Thus, among multiple full-sensitive filtration schemes, one clearly prefers the one that minimizes the number of transcripts detected.

**Table 4.1:** Best  $Q$ -grams found by Burkhardt and Kärkkäinen [2001] to solve instances with  $m = 50$  and  $k \in [4, 5]$ .  $Q$ -grams are visualized as strings over  $\{\#, -\}$ , e.g.  $Q = \{0, 2, 3, 4\} = \#-###$ . In the experimental evaluation of section 4.5,  $Q$ -gram (iii) is the best one for  $m = 100$  and  $k \in [8, 10]$ .

	Weight	Span	$Q$ -gram
(i)	13	14	#####-####
(ii)	11	30	###-#--#-----#-----#-#-#--##
(iii)	12	37	#-#-#---#-----#-#-#---#-----#-#-#---#

**Table 4.2:** Optimal threshold computation for the  $Q$ -grams in table 4.1 with  $m = 100$  and  $k \in [8, 10]$ . Threshold column shows the lower bound given the  $q$ -gram lemma (LEMMA), the optimal threshold (EXACT) and the upper bound given by algorithm 4.1 (APX). Runtime column shows the runtime in seconds of the DP algorithm by Burkhardt and Kärkkäinen [2001] and the ILP program 4.9 solved with IBM CPLEX 12.2; algorithm DP on  $Q$ -gram (iii) required 44 GB of RAM and could not complete within one hour; algorithm APX always completed within 6 milliseconds (data not shown).

Input		Threshold			Runtime [s]	
$Q$	$k$	LEMMA	EXACT	APX	ILP	DP
(i)	8	0	3	4	0.30	0.03
(i)	9	0	2	3	1.24	0.05
(i)	10	0	0	2	0.02	0.07
(ii)	8	0	7	9	0.20	31.60
(ii)	9	0	5	6	0.39	93.77
(ii)	10	0	1	3	0.32	244.83
(iii)	8	0	4	6	0.06	223.52
(iii)	9	0	3	4	0.08	885.37
(iii)	10	0	2	2	0.11	–

### Problem 4.3 (SPECIFICITY)

**Instance** A filtration scheme  $(Q, t)$ , an  $(m, k)$  instance.

**Solution**  $|\{X \in \mathbb{B}^m : t_Q^m(X) \geq t\}|$

### Approximate solution

I propose a randomized algorithm to estimate the number of transcripts detected by a filtration scheme of the type  $(Q, 1)$ . Such number coincides with the number of true assignments to the boolean function  $T_Q^m$ , which is in *disjunctive normal form* (DNF) by

definition (see equation 4.7). Counting the number of true assignments to an arbitrary boolean function in DNF is a classic *#P-complete* counting problem that however allows approximability to any degree [Karp *et al.*, 1989], i.e. it allows a *fully polynomial-time randomized approximation scheme* (FPRAS) [Vazirani, 2001]. Algorithm 4.2 shows the algorithm by Karp *et al.* for boolean function  $T_Q^m$ . Here, I follow the proof given by [Vazirani, 2001, chapter 28].

I denote with  $\mathbb{D}$  the set of transcripts detected by  $T_Q^m$

$$\mathbb{D} = \{X \in \mathbb{B}^m : T_Q^m(X)\} \quad (4.15)$$

and with  $\mathbb{D}_i$ , for  $0 \leq i \leq m - s(Q)$ , the set of transcripts detected by the  $i$ -th clause of  $T_Q^m$ , i.e. by the  $Q$ -gram starting at position  $i$ ,

$$\mathbb{D}_i = \left\{ X \in \mathbb{B}^m : \bigwedge_{j \in Q} x_{i+j} \right\} \quad (4.16)$$

whose cardinality is

$$|\mathbb{D}_i| = 2^{m-w(Q)} \quad (4.17)$$

as a 0 at any  $x_{i+j}$  would invalidate the  $Q$ -gram starting at position  $i$ . Problem SPECIFICITY for  $(Q, 1)$  involves computing

$$|\mathbb{D}| = \left| \bigcup_{i=0}^{m-1} \mathbb{D}_i \right| \quad (4.18)$$

which is less than the sum of all  $|\mathbb{D}_i|$  and hard to compute in general.

One might think of a randomized algorithm that estimates  $|\mathbb{D}|$  by repeatedly sampling  $X \in \mathbb{B}^m$  and counting whether  $X \in \mathbb{D}$  too. However, it is immediate to see that such an algorithm cannot estimate  $|\mathbb{D}|$  precisely and in polynomial time when  $\mathbb{D}$  is very small. The idea of algorithm 4.2 is thus to boost this estimate by sampling transcripts that are surely in  $\mathbb{D}$  and subsequently adjusting the counts. Algorithm 4.2 guarantees that

$$\Pr [|\text{COUNTTRANSCRIPTS}(Q, m, \epsilon) - |\mathbb{D}|| \geq \epsilon |\mathbb{D}|] \leq \frac{1}{4} \quad \forall \epsilon > 0. \quad (4.19)$$

Recall that  $t_Q^m(X)$  indicates the number of clauses of  $T_Q^m$  which are satisfied by  $X$ . I denote by  $\mathbb{M}$  the multi-set

$$\mathbb{M} = \{(X, t_Q^m(X)) : X \in \mathbb{D}_i \text{ for } 0 \leq i \leq m - s(Q)\} \quad (4.20)$$

whose cardinality is easily computed (line 2) as

$$|\mathbb{M}| = \sum_{i=0}^{m-s(Q)} |\mathbb{D}_i| = (m - s(Q) + 1) \cdot 2^{m-w(Q)}. \quad (4.21)$$

Algorithm 4.2 first (line 6) picks a  $\mathbb{D}_i$  with probability

$$\frac{|\mathbb{D}_i|}{|\mathbb{M}|} = \frac{1}{m - s(Q) + 1} \quad (4.22)$$

and then (lines 7–9) picks an  $X \in \mathbb{D}_i$  uniformly

$$\frac{1}{|\mathbb{D}_i|} = \frac{1}{2^{m-w(Q)}} \quad (4.23)$$

which results in picking an  $X \in \mathbb{M}$  with probability

$$\sum_{X \in \mathbb{B}^m : T_Q^m(X)} \frac{|\mathbb{D}_i|}{|\mathbb{M}|} \cdot \frac{1}{|\mathbb{D}_i|} = \sum_{X \in \mathbb{B}^m : T_Q^m(X)} \frac{1}{|\mathbb{M}|} = \frac{t_Q^m(X)}{|\mathbb{M}|}. \quad (4.24)$$

Algorithm 4.2 counts  $X$  (line 10) with weight

$$T = \frac{|\mathbb{M}|}{t_Q^m(X)} \quad (4.25)$$

so that the expected value of  $T$  is

$$E[T] = \sum_{X \in \mathbb{B}^m : T_Q^m(X)} \frac{t_Q^m(X)}{|\mathbb{M}|} \cdot \frac{|\mathbb{M}|}{t_Q^m(X)} = |\{X \in \mathbb{B}^m : T_Q^m(X)\}| = |\mathbb{D}|. \quad (4.26)$$

The standard deviation of  $T$  is bound w.r.t. the expected value (see [Vazirani, 2001, chapter 28]) as

$$\frac{\sigma(T)}{E[T]} \leq m - s(Q), \quad (4.27)$$

hence the number of times  $N$  (line 3) that algorithm 4.2 has to sample  $T$  (line 10) in order to guarantee the bound claimed in equation 4.19 is derived by applying Chebyshev's inequality

$$\Pr[|T - E[T]| \geq \epsilon E[T]] \leq \left( \frac{1}{\epsilon \sqrt{N}} \cdot \frac{\sigma(T)}{E[T]} \right)^2 \leq \frac{1}{4} \quad (4.28)$$

and solving it for  $N$ , which gives

$$N \geq 4 \cdot \frac{m - s(Q)}{\epsilon^2}. \quad (4.29)$$

A straightforward generalization of algorithm 4.2 for any filtration scheme  $(Q, t)$  consists into expanding boolean function  $T_Q^m$  to have  $\binom{m-s(Q)+1}{t}$  clauses, each one accounting for the joint occurrence of  $t$   $Q$ -grams. Whether the exponential dependency on  $t$  can be removed is an open question. In practice, algorithm 4.2 would be used in  $Q$ -grams design to choose the most specific  $Q$ -gram among a set of full-sensitive candidates. However, as my choice for the applications of chapters 6–7 felt upon the more practical seeds filters, I did not have to perform  $Q$ -grams design.

---

**Algorithm 4.2** COUNTTRANSCRIPTS( $Q, m, \epsilon$ )

---

**Input**      $Q$  :  $Q$ -gram  
               $m$  : integer denoting the pattern length  
               $\epsilon$  : real  $\in [0, 1]$  denoting the approximation factor

**Output**   integer indicating the number of transcripts detected

```

1:  $C \leftarrow m + s(Q) + 1$ 
2:  $M \leftarrow C \cdot 2^{m-w(Q)}$ 
3:  $N \leftarrow 4 \cdot (C - 1)^2 / \epsilon^2$ 
4:  $T \leftarrow 0$ 
5: repeat
6:    $i \leftarrow \text{RANDOM}([0, C])$ 
7:    $X \leftarrow \text{RANDOM}(\{0, 1\}^m)$ 
8:   for all  $j \in Q$  do
9:      $x_{i+j} \leftarrow 1$ 
10:   $T \leftarrow T + t_Q^m(X)$ 
11:  $N$  times
12: return  $N \cdot M / T$ 

```

---

#### 4.4.6 Families

To obtain even more specific filtration, Kucherov *et al.* [2005] propose *q-gram families* (also known as *multiple gapped q-grams*). Filtration with a  $q$ -gram family adopts disjunctively a set of multiple distinct gapped  $q$ -grams. The generalized counting argument now adds all occurrences of all gapped  $q$ -gram in the set. Figure 4.6 illustrates.

**Definition 4.5 ( $\mathbb{F}$ -gram).** A  $q$ -gram family (abbreviated as  $\mathbb{F}$ -gram) is a finite set  $\mathbb{F} = \{Q_0, \dots, Q_{f-1}\}$  of  $Q$ -grams. Its counting threshold  $\tau_{\mathbb{F}}$  is defined as:

$$\tau_{\mathbb{F}}(m, k) = \sum_{Q_i \in \mathbb{F}} \tau_{Q_i}(m, k) \quad (4.30)$$

All design problems introduced in section 4.4.2 apply also to  $\mathbb{F}$ -grams. The methods presented in section 4.4.4 and 4.4.5 naturally generalize to  $q$ -gram families. To this extent, it is sufficient to define the boolean function for an  $\mathbb{F}$ -gram as the disjunction

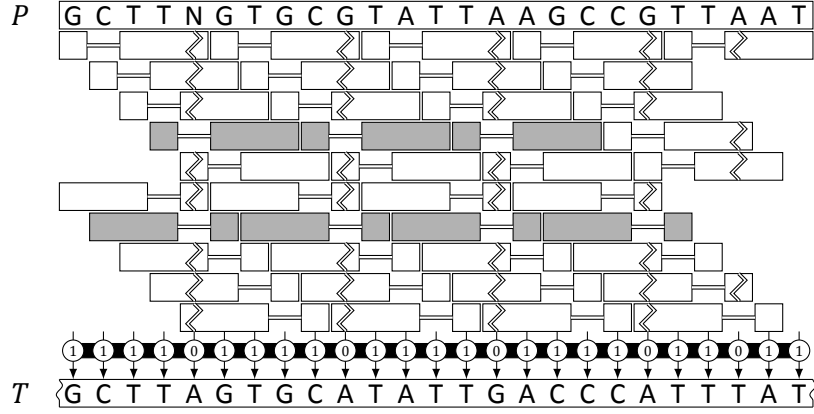
$$T_{\mathbb{F}}^m(X) = \bigvee_{Q_i \in \mathbb{F}} T_{Q_i}^m(X) \quad (4.31)$$

and its pseudo-boolean function as the sum

$$t_{\mathbb{F}}^m(X) = \sum_{Q_i \in \mathbb{F}} t_{Q_i}^m(X). \quad (4.32)$$

Function  $t_{\mathbb{F}}^m$  is still supermodular, because all supermodular functions are closed under non-negative linear combination. Said that, I do not use  $\mathbb{F}$ -grams in practice.

**Figure 4.6:** Filtration with multiple gapped  $q$ -grams. A filtration scheme  $(\mathbb{F}, t) = (\{\{0, 2, 3, 4\}, \{0, 1, 2, 4\}\}, 7)$  solves instance  $(25, 5)$ . In the illustration, pattern  $P$  of length 25 occurs in text  $T$  at edit distance 5. The seven  $q$ -grams in grey are not covered by any mismatch.



## 4.5 Evaluation

In this section, I show the results of an experimental evaluation of the filtration methods exposed so far. I consider seed filtration schemes with exact, 1 and 2-approximate seeds (Exact, 1-Apx and 2-Apx seeds), contiguous  $q$ -grams filtration schemes using the maximum lossless value of  $q$  for a threshold of 1 ( $q$ -Grams,  $t \geq 1$ ) and 4 ( $q$ -Grams,  $t \geq 4$ ), and the gapped  $q$ -gram (iii) from table 4.1 with the optimal thresholds indicated in table 4.2 ( $Q$ -Grams).

For indexed filters, I use the  $q$ -gram index with  $q = 10$  (see 3.1.3). To perform edit distance verifications, I use a banded version of the Myers' algorithm [Myers, 1999] by [Weese *et al.*, 2012]. As text, I take the *C. elegans* reference genome (WormBase WS195), i.e. a collection of 6 DNA strings of 100 Mbp total length. As patterns, I extrapolated 200k DNA sequences of length 100 bp from an Illumina sequencing run (Sequence Read Archive ID SRR065390).

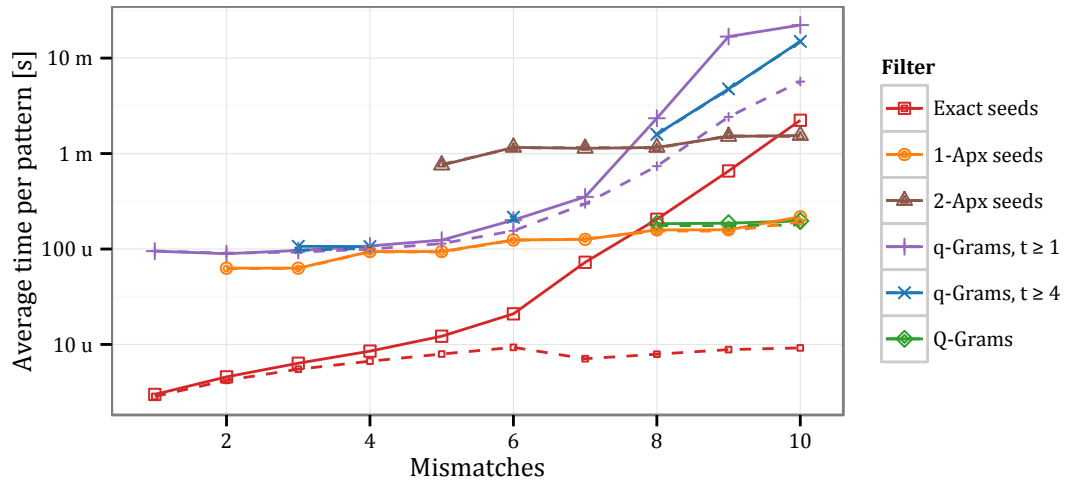
All experiments run on a desktop computer running Linux 3.10.11, equipped with one Intel® Core i7-4770K CPU @ 3.50 GHz, 32 GB RAM and a 2 TB HDD @ 7200 RPM. I repeated the experiments for  $k$ -mismatches and  $k$ -differences, varying  $k$  in the range  $[1, 10]$ . I measured the runtime of the filtration phase only, and then of the filtration plus verification phase. The plots show always *average* runtimes (or values) per pattern.

### 4.5.1 Runtime

Figure 4.7 shows the results for  $k$ -mismatches. Exact seeds are the best filtration method for  $k \leq 7$ , mainly due to their superior filtration speed, while 1-approximate seeds are better for  $k \geq 8$ . 2-Approximate seeds start to dominate exact seeds only for  $k \geq 10$ , i.e. they provide too strong filtration on this text. Both contiguous  $q$ -gram filtration schemes



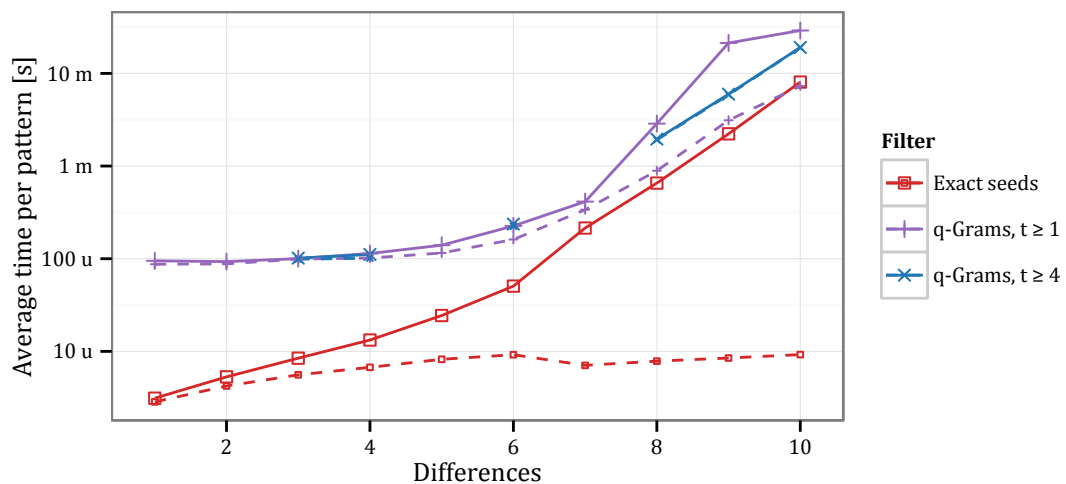
**Figure 4.7:** Filters runtime on  $k$ -mismatches. Pattern lengths are fixed to 100 bp. Solid lines represent total runtimes, while dashed lines represent filtration times only.



are always slower than 1-approximate seeds; it can be seen that enforcing  $t \geq 4$  improves the total runtime for those instances where  $t \geq 1$  renders the filter too weak. Gapped  $q$ -grams are on par with 1-approximate seeds, slightly better only for  $k = 10$ .

Figure 4.8 shows the results for  $k$ -differences. The more involved edit distance verification slightly fills the gap between exact seeds and contiguous  $q$ -grams, yet exact seeds continue to be always the fastest alternative.

**Figure 4.8:** Filters runtime on  $k$ -differences. Pattern lengths are fixed to 100 bp. Solid lines represent total runtimes, while dashed lines represent filtration times only.



### 4.5.2 Verification versus filtration time

Figure 4.9 shows the ratios between the runtimes of the verification and filtration phases of each filtration scheme. For  $k \leq 6$ , all schemes spend more time on filtration rather than verification. The weakest scheme, filtration with exact seeds, shows the closest ratio to 1. As shown in the runtime plots (figures 4.7-4.8), quick filtration pays off for low error rates. Here, a quicker full-text index would be beneficial. For  $k \geq 7$ , contiguous  $q$ -grams with  $t \geq 1$  show the closest ratio to 1, nonetheless the fastest alternative is provided by filtration with 1-approximate seeds, for which only 10 % of the runtime goes in verifications. Here, a judicious mix of exact and 1-approximate seeds could improve the total runtime. The ratios for  $k$ -difference show analogous patterns (data not shown).

### 4.5.3 Positive predictive value

Instead of measuring filtration specificity, as introduced in section 2.4.3, I measure the *positive predictive value* (PPV). As shown in table 4.3, I define *true positives* (TP), *false positives* (FP), *true negatives* (TN), or *false negatives* (FN) in terms of the number of verifications  $V_f(T)$  produced in the filtration phase and the number of approximate occurrences  $C(T)$  of the pattern in the text. Therefore, I define the PPV as:

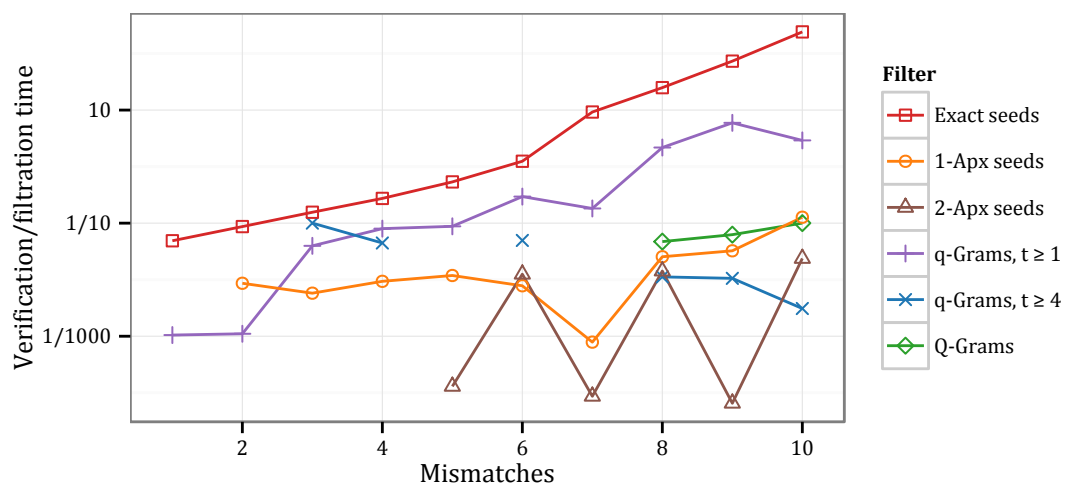
$$\frac{|TP_f(T)|}{|TP_f(T)| + |FP_f(T)|} = \frac{C(T)}{V_f(T)} \quad (4.33)$$

Figure 4.10 shows the results for  $k$ -mismatches. As expected, 2-approximate seeds have always the highest PPV, followed by 1-approximate seeds. The PPV of gapped  $q$ -grams is comparable to that of 2-approximate seeds, while that of contiguous  $q$ -grams with  $t \geq 1$  oscillates around that one of exact seeds. In particular, when  $t$  approaches 1, e.g. for  $k = 9$ , the PPV of contiguous  $q$ -grams stays below that one of exact seeds. Enforcing  $t \geq 4$  boosts the PPV of contiguous  $q$ -grams to that one of 1-approximate seeds. The PPVs for  $k$ -difference show analogous patterns (data not shown).

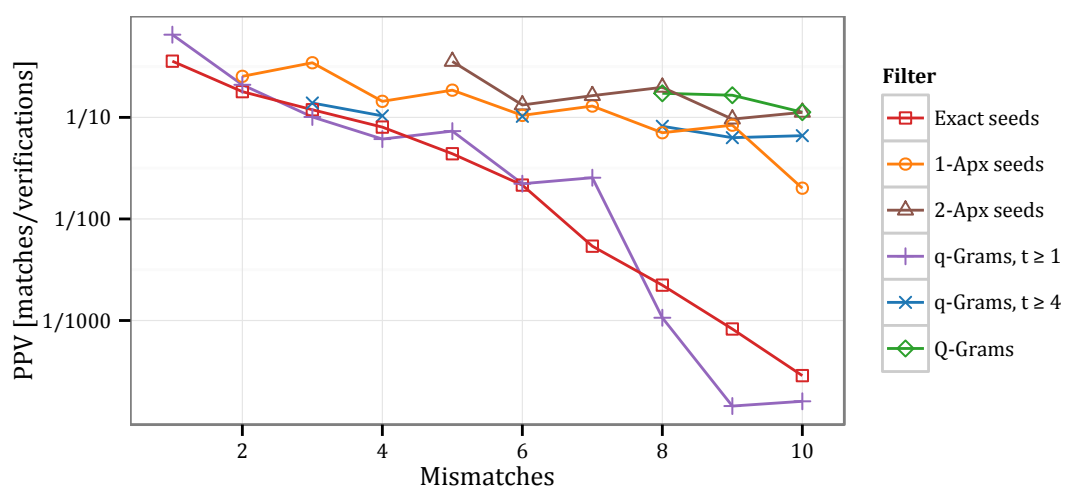
**Table 4.3:** Measurement of filtering methods efficiency.  $V_f(T)$  counts the number of verifications,  $C(T)$  the number of approximate occurrences, and  $|T|$  the text length. Since all considered filtering methods are full-sensitive, the number of false negatives is always 0.

	Positive	Negative
True	$C(T)$	$ T  - C(T)$
False	$V_f(T) - C(T)$	0

**Figure 4.9:** Filtration versus verification time on  $k$ -mismatches. Pattern lengths are fixed to 100 bp.



**Figure 4.10:** Filters specificity on  $k$ -mismatches. Pattern lengths are fixed to 100 bp.





## **Part II**

### **READ MAPPING**



In this chapter, I provide the reader with background knowledge in the fast-moving field of high-throughput sequencing (HTS). It goes without saying that the reader familiar with HTS can skip this chapter. In section 5.1, I briefly introduce the two most prominent HTS technologies and the kind of sequencing data they produce. Afterwards, in section 5.2, I explain how standard HTS data *analysis pipelines* are structured; in particular, in section 5.2.2, I present two popular paradigms for reference-guided assembly: best-mapping and all-mapping. In section 5.3, I review two studies on the limits of HTS data analysis. Finally, in section 5.4, I give an overview of the most popular read mapping tools.

## 5.1 High-throughput sequencing data

As anticipated in section 1.1, actual HTS technologies produce DNA reads which are shorter than Sanger sequencing and more likely to contain systematic sequencing artefacts. HTS data does not consist of read sequences only, but includes also base quality scores annotating the quality of the sequencing process.

### 5.1.1 Read sequences

I consider only sequencing data produced by the most prominent HTS instruments. *Illumina* is the actual market leader for HTS instruments, followed by *Life Technologies's Ion Torrent*. Some instruments, e.g. the *GS* by *454 Life Sciences* or the *SOLiD* by *Applied Biosystems*, became popular around 2006–2010, but are now discontinued. Other *third-generation* instruments, like the *single-molecule real-time* sequencing (SMRT) *RS II* by *Pacific Biosciences*, have great potential but still low impact on the HTS market.

#### **Illumina**

All Illumina instruments use the *sequencing by synthesis* (SBS) technology [Bentley *et al.*, 2008] developed by Solexa. In the library preparation phase, the DNA sample is sheared into smaller fragments. During the sequencing phase, single-stranded DNA fragments are attached on a slide called *flow cell* and amplified *in situ* using a variant of *polymerase chain reaction* (PCR) called *bridge amplification*. Clusters of amplified fragments are then

used as templates for multiple cycles of synthetic sequencing with four differentially labeled fluorescent reversible terminator *deoxyribonucleotide triphosphates* (dNTPs). In each sequencing cycle, dNTPs are incorporated into the fragments in each cluster, their corresponding fluorescent reversible terminators are imaged by a high-resolution camera, and then dNTPs are cleaved to allow incorporation of the next base.

After the sequencing phase, a base calling software converts all images corresponding to one cluster into one read. The software measures the intensities of the four colors imaged during the  $i$ -th cycle, calls the  $i$ -th read base, and assigns a base quality score. Eventual replication errors, made by DNA polymerase during bridge amplification, result in mixed signal intensities within a cluster, and hence to base calls of lower confidence.

Illumina's SBS technology further allows the sequencing of both ends of each DNA fragment. Two kind of libraries are available: *paired-end* libraries, consisting of reads from short-insert 300–600 bp fragments, and *mate pair* libraries, consisting of reads separated by several kilobases. The former libraries are adopted for high-resolution genome resequencing, while the latter ones provide accurate *de novo* sequence assembly or detection of large-scale structural variation.

## Ion Torrent

*Life Technologies's Ion Torrent* instruments use a semiconductor-based technology [Rusk, 2010]. During library preparation, the sheared single-stranded template DNA is embedded into *microwells* on a semiconductor chip with DNA polymerase. Microwells are sequentially flooded with unmodified A, C, G or T dNTP. The DNA polymerase incorporates the introduced dNTP into the growing strand only if this is complementary to the leading template nucleotide. Such polymerization reaction releases an hydrogen ion, which in turn changes the pH of the solution. A *ion-sensitive field-effect transistor* (ISFET) in the sequencing instrument measures this pH change. Any *homopolymer* in the template causes multiple dNTP to be incorporated in a single cycle and results in a higher pH change, which is not precisely measured by the instrument and thus causes systematic sequencing errors.

### 5.1.2 Phred base quality scores

Base quality scores have been introduced by the base calling tool *Phred* [Ewing *et al.*, 1998; Ewing and Green, 1998] to assess sequencing quality of single bases in capillary reads. Instead of discarding low-quality regions present in capillary reads, Phred output each base and annotates it with the probability that it has been wrongly called. The tool encodes the probability  $\epsilon_i$  of miscalling the  $i$ -th base in a read under the form of a base quality  $Q_i$  in logarithmic space:

$$Q_i = -10 \log_{10} \epsilon_i. \quad (5.1)$$

This method has been unanimously accepted. All sequencing technologies complement DNA reads with base quality scores, encoded in Phred-scale or similar.



## 5.2 High-throughput sequencing data analysis

### 5.2.1 Data analysis pipelines

The analysis of HTS data consists of numerous steps, ranging from the initial instrument-specific data processing to the final application-specific interpretation of the results. It is conventional wisdom to subdivide such data analysis pipelines in three stages of analysis. The *primary stage* of analysis consists of instrument specific steps for the generation, collection, and processing of raw sequencing signals. The *secondary stage* applies sequence analysis methods to the raw sequencing data in order to reconstruct the original sequence of the donor genome. The *tertiary stage* characterizes application-specific features of the donor genome and then provides interpretations, e.g. in whole exome sequencing (WES) it consists of calling genetic variations and predicting their pathogenicity. Below, I give more

#### Primary analysis

Primary analysis consists of instrument specific steps to call base pairs and compute quality metrics. The base calling software converts raw sequencing signals into bases, i.e. A, C, G, T, or N if the signal is unclear. The software assigns a quality value to each called base, estimating the probability of a base calling error. On early generation instruments, users could provide their own base calling tool. Now this process happens automatically on special hardware (e.g. FPGAs or GPUs) bundled within the instrument. The result of primary analysis is a standard *FASTQ* file containing DNA read sequences and their associated quality scores in *Phred* format (see section 5.1.2).

#### Secondary analysis

Secondary analysis aims at reconstructing the original sequence of the donor genome from its HTS reads. There are two main *plans* to reassemble the original genome: (A) *de novo* assembly, and (B) *reference-guided* assembly (commonly called *read mapping*). *De novo* assembly is very involved as it essentially requires finding a *shortest common superstring* (SCS) of the reads, which is a NP-complete problem [Maier and Storer, 1977; Gallant *et al.*, 1980; Turner, 1989]. Computational methods in plan A first scaffold the reads by performing *overlap alignments* [Myers, 2005] or equivalently by constructing *de Bruijn graphs* [Pevzner *et al.*, 2001]. The knowledge of a *reference genome*, highly similar to the donor, simplifies the problem and opens the door to plan B. The reads are simply aligned to the reference, tolerating a few base pair errors. It is worth mentioning that plans A and B can be combined, e.g. as proposed by Li [2012].

Plan B is always preferred in resequencing projects, as it is computationally more viable than plan A and directly provides a way to assess genetic variation w.r.t. a reference genome. In this work, I consider only plan B and present methods that work within this specific plan. Nonetheless, many of the algorithmic components that I introduced in the first part of this manuscript are ubiquitous in bioinformatics, thus applicable, if not

already applied, to plan A. According to plan B, the secondary analysis stage consists of three steps: *quality control*, *read mapping*, and *consensus alignment*.

In the *quality control* (Q/C) step, the quality of raw reads is checked. Reads produced by current HTS technologies contain sequencing artefacts in form of single bases or stretches of miscalled oligonucleotides. In order to circumvent this problem, various techniques have been developed, ranging from the simple *trimming* of low quality read stretches to sophisticated methods of *read error correction* [Weese *et al.*, 2013]. Sometimes, Q/C is simply omitted.

In the *read mapping* step, the reads are aligned to the reference genome. Read mapping tools adopt state of the art *approximate string matching* methods to efficiently analyze the deluge of data produced by HTS instruments. Approximate matching accounts for two kinds of errors: residual sequencing artefacts not removed by Q/C and small genetic variations in the donor genome. Consequently, a read mapper must take into account these errors when aligning the reads.

Mapped reads, often stored in de-facto standard BAM files [Li *et al.*, 2009a], are sorted by genomic coordinate and eventually multiply aligned in order to construct a *consensus sequence* of the donor genome. The height of the pileup denotes the sequencing depth or coverage at each locus in the donor genome; higher coverage implies more confidence in the consensus sequence of the donor genome and thus more accuracy in the tertiary analysis stage.

### Tertiary analysis

Tertiary analysis aims at interpreting the information provided by the secondary analysis stage. This pipeline stage groups a wide range of analyses specific to the sequencing application. In some pipelines, downstream analysis aggregates data coming from multiple samples.

Within DNA resequencing, *genotyping* consists of determining the variations between the donor and the reference genome. The result of *variant calling* is a set of variations characterizing the donor genome, usually stored in de-facto standard VCF files. Subsequently, *genome-wide association studies* (GWAS) associate genetic variants with phenotypic traits by examining *single-nucleotide polymorphisms* (SNPs), variants relatively common among individuals of the same population.

In RNA-seq, tertiary analysis consists of computing transcripts abundance by measuring *reads per kilobase per million mapped reads* (RPKM) [Mortazavi *et al.*, 2008]; subsequently, relative gene expression is determined by comparing multiple RNA samples. In ChIP-seq, this stage consists of calling peaks in correspondence of mapped reads, determining which peaks identify feasible transcription factor binding sites, then interesting sites affecting gene regulation.

## 5.2.2 Secondary analysis paradigms

The goal of secondary analysis is to reconstruct the original sequence of the donor genome. The *reference-guided* assembly plan assumes prior knowledge of a reference genome

which is close the donor genome. Reads are thus mapped (i.e. aligned) to the reference genome w.r.t. a given scoring scheme and threshold. The scoring scheme accounts for eventual genomic variation, as well as for any sequencing artefacts. Under these assumptions, an alignment of optimal score for a read implies its *original location* on the reference genome. Conversely, no alignment within the score threshold implies too many sequencing errors, too much genetic variation, or sample contamination. A problem arises in presence of co-optimal or close suboptimal alignments: the read cannot be mapped confidently to one single location.

The problem of confidently mapping high-throughput sequencing reads comes from the non-random nature of genomic sequences. Genomes evolved through multiple types of duplication events, including whole-genome duplications [Wolfe and Shields, 1997; Dehal and Boore, 2005] or large-scale segmental duplications in chromosomes [Bailey *et al.*, 2001; Samonte and Eichler, 2002], transposition of repetitive elements as short tandem repeats (microsatellites) [Wang *et al.*, 1994; Wooster *et al.*, 1994] and interspersed nuclear elements (LINE, SINE) [Smit, 1996], proliferation of repetitive structural elements such as telomeres and centromeres [Meyne *et al.*, 1990]. As a result of these events, for instance, about 50 % of the human genome is composed of repeats.

Repeats present in general technical challenges for all *de novo* assembly and sequence alignment programs [Treangen and Salzberg, 2011]. Due to repetitive elements, a non-ignorable fraction of high-throughput sequencing reads cannot be mapped confidently. In general, the shorter the reads, the higher the challenges due to repeats. I quantify this phenomenon more precisely in section 5.3. Here I focus on analysis strategies to deal with *multi-reads*, i.e. reads that cannot be mapped confidently as they align equally well to multiple locations.

It is not evident how to treat *multi-reads*. According to Treangen and Salzberg [2011], common strategies to deal with multi-reads are (i) to discard them all, (ii) to randomly pick one best mapping location, (iii) to consider all or up to  $k$  best mapping locations within a given distance threshold. A *de facto* standard strategy, combining strategies (i) and (ii), emerged over the last years. The read mapper randomly picks one best mapping location and complements it with its *mapping quality*, i.e. the probability of the mapping location being correct (see section 5.2.3). Subsequently, downstream analysis tools either (a) apply a mapping quality score cutoff to discard reads not mapping confidently to any location, or (b) annotate in turn their results with quality scores. The other popular strategy adopted by analysis tools is to consider all mapping locations within an edit distance threshold. In this case, it is not clear whether downstream analysis tools consider all mapping locations equal regardless of their distance.

In the light of these facts, I define two broad paradigms for the secondary and tertiary analysis of HTS data: *best-mapping* and *all-mapping*. The best-mapping paradigm considers a single mapping location per read along with its confidence, while the all-mapping paradigm considers a comprehensive set of mapping locations per read. It goes without saying that read mapper and downstream analysis tools must agree on a common paradigm. Thus these paradigms are valid not only for read mappers but also for any downstream analysis tool, e.g. variant callers. Read mapping and variant calling are indeed tightly coupled steps within reference-based HTS pipelines.

### 5.2.3 Best-mapping

As said, best-mapping methods rely on a single mapping location per read. In order to maximize recall, best-mappers often adopt complex scoring schemes taking into account gaps and base quality values, and at the same time implement sophisticated heuristics to speed up the search. Best-mappers annotate any mapping location with its mapping quality. Subsequently, in order to maximize precision, variant calling tools decide whether to consider or discard reads not mapping confidently to any location. The GATK [DePristo *et al.*, 2011] and Samtools [Li *et al.*, 2009a] are popular best-mapping tools to call small variants. In section 5.3, I discuss how this paradigm systematically fails on reads belonging some critical genomic regions, thus is limited to the analysis of *high mapability* regions.

#### Mapping quality score

Mapping quality has been introduced in the tool MAQ [Li *et al.*, 2008]. The study considers short 30–40 bp reads, produced by early Illumina and ABI/SOLiD sequencing technologies, whose sequencing error rates were quite high. Given the short lengths and high error rates, a significant fraction of such reads can be aligned to multiple mapping locations, even considering only co-optimal Hamming distance locations. Since base callers output base call probabilities in Phred-scale along with the reads, Li *et al.* propose a novel probabilistic scoring scheme called *mapping quality*, encoding the probability that a read aligns correctly at a mapping location in the reference genome.

Mapping quality scores offer a way to prioritize the results produced in downstream data analyses. Li *et al.* [2008] write that “*it is possible to act conservatively by discarding reads that map ambiguously at some level, but this leaves no information in the repetitive regions and it also discards data, reducing coverage in an uneven fashion, which may complicate the calculation of coverage.*” For instance, the GATK HaplotypeCaller [DePristo *et al.*, 2011] annotates its variant calls with qualities whose value depends on mapping qualities, rather than removing data by applying a hard mapping quality cutoff. Below, I define the mapping quality score as done in [Li *et al.*, 2008].

Fix the alphabet  $\Sigma = \{A, C, G, T\}$ . Consider a known donor genome  $G$  over  $\Sigma$  and a read  $R$  sequenced at location  $l$  from the template  $G_{l \dots l+|R|}$ . The base calling error  $\epsilon_i$  from equation 5.1 represents the probability  $\epsilon_i$  of miscalling a base  $r_i$  instead of calling its corresponding base  $g_{l+i}$  in the donor genome. The probability  $\Pr[r_i | g_{l+i}]$  of observing the base  $r_i$  given the donor genome base  $g_{l+i}$ , is:

$$\Pr[r_i | g_{l+i}] = \begin{cases} 1 - \epsilon_i & \text{if } g_{l+i} = r_i \\ \frac{\epsilon_i}{|\Sigma| - 1} & \text{if } g_{l+i} \in \Sigma \setminus \{r_i\} \end{cases} \quad (5.2)$$

and assuming i.i.d. base calling errors, it follows that the probability  $\Pr[R | G, l]$  of observing the read  $R$ , given the donor genome template  $G_{l \dots l+|R|}$ , is:

$$\Pr[R | G, l] = \prod_{i=0}^{|R|-1} \Pr[r_i | g_{l+i}]. \quad (5.3)$$

By applying Bayes' theorem, Li *et al.* [2008] derive the posterior probability  $\Pr[l|G, R]$ , that location  $l$  in the reference genome  $G$  is the correct mapping location of read  $R$ . Assuming uniform coverage, each location  $l \in [0, |G| - |R|]$  has equal probability of being the origin of a read in the donor genome, thus the prior probability  $\Pr[l]$  is simply:

$$\Pr[l] = \frac{1}{|G| - |R| + 1} \quad (5.4)$$

Therefore, recalling  $\Pr[R|G, l]$  from equation 5.3, the posterior probability  $\Pr[l|G, R]$  equals the probability of the read  $R$  originating at location  $l$ , normalized over all possible locations in the reference genome:

$$\Pr[l|G, R] = \frac{\Pr[R|G, l]}{\sum_{i=0}^{|G|-|R|} \Pr[R|G, i]} \quad (5.5)$$

which in Phred-scale becomes:

$$Q[l|G, R] = -10 \log_{10}(1 - \Pr[l|G, R]). \quad (5.6)$$

Computing the exact mapping quality as in equation 5.6 requires aligning each read to all positions in the reference genome. On the one hand, this computation is practically infeasible. On the other hand, suboptimal locations not close to the optimum one contribute very little to the sum in equation 5.5. Therefore, read mapping programs approximate equation 5.5 using only the mapping locations they repute relevant.

Some objections can be raised against the above definition of mapping quality scores. First, the score is derived under the unlikely assumption of the reference genome being equal to the donor genome. In other words, equation 5.2 considers only errors due to base miscalls and disregards genetic variation; thus the risk is to prefer mapping locations supported by known low base qualities rather than by true but unknown SNVs. Second, mapping quality is nonetheless strongly correlated to mapping uniqueness, as discussed in section 5.3; it is easy to see that the probability of any location in equation 5.5 dilutes in presence of a large number of co-optimal or close suboptimal mapping locations. Therefore, in chapter 7, I give an alternative definition of mapping quality.

### 5.2.4 All-mapping

All-mapping analysis methods consider a comprehensive set of locations per read. Almost all read mappers in this category adopt edit distance and report all mapping locations within an error threshold, absolute or relative w.r.t. to the length of the reads. Variant calling algorithms based on all-mapping have the potential to detect a wider spectrum of genomic variation events than their best-mapping counterparts. For instance, variant callers based on the all-mapping paradigm detect CNVs [Alkan *et al.*, 2009], and SNVs in homologous regions [Simola and Kim, 2011].

### 5.3 Limits of high-throughput sequencing

A fraction of high-throughput sequencing reads cannot be mapped confidently due to repetitive elements. Which regions of a model organism's genome cannot be resequenced confidently by a high-throughput sequencing technology? And how accurate is downstream analysis on these low confidence regions? Two recent studies [Derrien *et al.*, 2012; Lee and Schatz, 2012] answer these questions. Below, I report their key ideas and most relevant findings.

#### 5.3.1 Genome mappability

Derrien *et al.* [2012] define *genome mappability* as a function of a genome for a fixed  $q$ -gram length, distance measure i.e. the Hamming or edit distance, and distance threshold  $k$ . Given a genomic sequence  $G$ , they define the  $(q, k)$ -frequency  $F_k^q(l)$  of the  $q$ -gram  $G_{l...l+q}$  at location  $l$  in  $G$  as the number of occurrences of the  $q$ -gram in  $G$  and its reverse complement  $\bar{G}$ . The  $(q, k)$ -mappability  $M_k^q(l)$  is the inverse  $(q, k)$ -frequency, i.e.  $M_k^q(l) = 1/F_k^q(l)$  with  $M_k^q : \mathbb{N} \rightarrow ]0, 1]$ . Note that  $M_k^q(l)$  can be seen as the prior probability that any read of length  $q$  originating at location  $l$  will be mapped correctly. The values of  $(q, k)$ -frequency and mappability obviously vary with the distance threshold  $k$ . Nonetheless, under any distance measure, it holds that the  $q$ -gram at location  $l$  is unique up to distance  $k$  iff  $M_k^q(l) = 1$  and repeated otherwise.

Unique mappability determines which fraction of a genome can be analyzed according to strategy (i) of Treangen and Salzberg [2011] (i.e. discarding non-unique reads, see section 5.2.2). Derrien *et al.* quantify the *unique mappability* of whole human, mouse, fly, and worm genomes. Mimicking typical Illumina read mapping setups, they consider  $q$ -grams of length 36, 50 and 75 bp, and Hamming distance 2. They find out that about 30 % of the whole human genome is not unique w.r.t. (36, 2)-mappability. At (75, 2)-mappability, 17 % of the human genome is not yet unique. This last result is slightly optimistic, as typical mapping setups call for up to 3–4 edit distance errors in order to map a significant fraction of the reads. Table 5.1 shows some results obtained from [Derrien *et al.*, 2012].

To estimate single-base resequencing accuracy, Derrien *et al.* consider the mappability

**Table 5.1:** Mappability of model genomes. Data obtained from [Derrien *et al.*, 2012].

	H.sapiens (hg19)	M.musculus (mm9)	D.mel (dm3)
Repeats content [%]	45.25	42.33	26.50
Uniqueome (36 bp, 2 msm) [%]	69.99	72.07	68.09
Uniqueome (50 bp, 2 msm) [%]	76.59	77.06	69.44
Uniqueome (75 bp, 2 msm) [%]	83.09	81.65	71.00

bility of all possible  $q$ -grams spanning any single genomic location. They define *pileup mappability*  $P_k^q$  at position  $l$  as the average mappability of all  $q$ -grams spanning position  $l$ :

$$P_k^q(l) = 1/q \sum_{i=l-q}^{l+q} M_k^q(i). \quad (5.7)$$

Derrien *et al.* [2012] find out in their own resequencing studies that “*low pileup-mappability regions are more prone to show a high value of heterozygosity than those with high mappability*”. Ideally, variant calling tools call a locus as heterozygous whenever the consensus alignment column at that locus contains two distinct bases. This situation tends to arise more frequently whenever the consensus alignment contains reads originating from similar yet distinct regions.

### 5.3.2 Genome mappability score

Genome mappability score (GMS) [Lee and Schatz, 2012], analogously to pileup mappability, estimates single-locus resequencing accuracy for a specific sequencing technology. Instead of considering the inverse  $q$ -gram frequency, Lee and Schatz use mapping quality (see section 5.2.3) to estimate the probability that a read originating at a given position can be mapped correctly. Subsequently, they derive the average mapping probability of all reads (in  $\mathbb{R}_l$ ) spanning a location  $l$  of the reference genome  $G$  as:

$$\Pr[l|G] = \sum_{R \in \mathbb{R}_l} \frac{\Pr[l|G, R]}{|\mathbb{R}_l|} \quad (5.8)$$

which in Phred-scale becomes:

$$Q[l|G] = \sum_{R \in \mathbb{R}_l} \frac{1 - 10^{-Q[l|G, R]/10}}{|\mathbb{R}_l|}. \quad (5.9)$$

Thus, fixed a genomic sequence  $G$ , they define the genome mappability score  $\text{GMS}(l)$  in percentual value:

$$\text{GMS}(l) = 100 \cdot Q[l|G] \quad (5.10)$$

Lee and Schatz proceed as follow to compute GMS. They first simulate reads from all genomic locations, having length and error profiles similar to those issue by actual sequencing technologies. Subsequently, they compute mapping quality scores by mapping all simulated reads with the best-mapper BWA [Li and Durbin, 2009]. Then, as just explained, they compute GMS at any location by averaging the quality scores. Finally, they define *low GMS* regions as those locations for which  $\text{GMS}(l) \leq 10$  and *high GMS* otherwise. Table 5.2 shows the performance of various sequencing technologies on the whole human genome (data obtained from [Lee and Schatz, 2012]).

**Table 5.2:** Human genome mappability score of various sequencing technologies. Data obtained from [Lee and Schatz, 2012].

Sequencing technology	Read length [bp]	Error rate [%] (msm, ins, del)	Low GMS [%]	High GMS [%]
Illumina-like	100	(0.10, 0.00, 0.00)	10.51	89.49
Ion Torrent-like	200	(0.04, 0.01, 0.95)	9.35	90.65
Roche/454-like	800	(0.18, 0.54, 0.36)	8.91	91.09
PacBio EC-like	2000	(0.33, 0.33, 0.33)	8.61	91.39

Lee and Schatz measure variant calling accuracy by GMS for the popular combination of best-mapping tools BWA and SAMtools [Li *et al.*, 2009a]. They simulate an Illumina-like resequencing study and feed it to such analysis pipeline. They find out that, at  $30\times$  sequencing coverage, accuracy approaches 100 % in high GMS regions, while it levels off to 25 % in low GMS regions. Their analysis “shows that most SNP detection errors are false negatives, and most of the missing variations are in regions with low GMS scores” [Lee and Schatz, 2012]. These are the limits of the analysis of high-throughput sequencing data.

## 5.4 Popular read mappers

Following the boom of NGS technologies, recent bioinformatics research has produced dozens of tools to perform read mapping. Two surveys [Li and Homer, 2010; Fonseca *et al.*, 2012] try to help bioinformaticians to find their way in the jungle of read mapping tools. The survey by Li and Homer first classifies read mapping algorithms by data structure: those based on hash tables and those based on suffix/prefix trees. However, the adopted data structure is often an implementation detail, indeed most algorithms covered in their survey could fit into both classes. The survey primarily considers the application of SNP calling; in the considered setup, tools enumerating a comprehensive set of locations always lag behind those designed to report only one location per read. The survey by Fonseca *et al.* instead catalogs read mappers by the features exposed to the user. It considers supported input-output formats, rate of errors and variation, number and type (i.e. local or semi-global read alignments) of mapping locations reported. After this exhaustive catalog, the survey concludes that the choice of a read mapper “involves application-specific requirements such as how well it works in conjunction with downstream analysis tools (i.e. variant callers)”. Read mapping and variant calling are indeed tightly coupled steps within reference-based HTS analysis pipelines.

As explained above, secondary and tertiary analysis methods are based on one of the two following paradigms: best-mapping and all-mapping. In the light of the above consideration, the most important feature of a read mapper is the number of mapping locations reported, followed by their type, while the other features are mostly of technical relevance. Most read mappers are specifically designed to fit one paradigm, while others



are versatile enough to work well in both cases.

The rest of this section presents most popular read mapping tools. Table 5.3 gives an overview of all these tools. Among them, BWA [Li and Durbin, 2009], Bowtie [Langmead *et al.*, 2009] and Bowtie 2 [Langmead and Salzberg, 2012], and Soap [Li *et al.*, 2009b] are prominent tools designed for best-mapping, while SHRiMP 2 [David *et al.*, 2011], mr(s)Fast [Alkan *et al.*, 2009; Hach *et al.*, 2010], RazerS [Weese *et al.*, 2009] and RazerS 3 [Weese *et al.*, 2012], and Hobbes 2 [Kim *et al.*, 2014] are designed for all-mapping. Grosso modo, most prominent best-mappers recursively enumerate substrings on a suffix/prefix tree of the reference genome via backtracking algorithms. Backtracking alone is impractical as its time complexity grows exponentially with the number of errors considered, hence best-mappers apply heuristics to reduce and prioritize the enumeration. Conversely, all-mappers are based on filtering algorithms for approximate string matching. They quickly determine, often with the help of an index, locations of the reference genome candidate to contain approximate occurrences, then verify them with conventional methods. Their efficiency is bound to filtration specificity and thus deteriorates with increasing error rates and genome lengths. GEM [Marco-Sola *et al.*, 2012] tries to fit both best and all-mapping paradigms. It speeds up best-mapping by stratifying mapping locations by edit distance and prioritizing filtration accordingly. Finally, Masai [Siragusa *et al.*, 2013a] and Yara are read mapping programs developed by myself. I present the engineering and evaluation of these tools in chapters 6 and 7.

### 5.4.1 Bowtie and Bowtie 2

Bowtie [Langmead *et al.*, 2009] is a mapper designed to have a small memory footprint and quickly report a few good mapping locations for early generation short Illumina and ABI/SOLiD reads. The tool achieves the former goal by indexing the reference genome with an FM-index and the latter one by performing a greedy backtracking on it. The greedy top-down traversal visits first the subtree yielding the least number of mismatches and stops after having found a candidate (not guaranteed to be optimal when  $k > 1$ ). In addition, Bowtie speeds up backtracking by applying *case pruning* [Mäkinen *et al.*, 2010], a simple application of the pigeonhole principle. However, this technique is mostly suited for  $k = 1$  and requires the index of the forward and reverse reference genome. Bowtie can be configured to search by strata, but the search time increases significantly while the traversal still misses a large fraction of the search space due to seeding heuristics.

Bowtie 2 [Langmead and Salzberg, 2012] has been designed to quickly report a couple of mapping locations for longer Illumina, Ion Torrent and Roche/454 reads, usually having lengths in the range from 100 bp to 400 bp. This tool uses an heuristic seed-and-extend approach, collecting seeds of fixed length, partially overlapping, and searching them exactly in the reference genome using an FM-index. Bowtie 2 randomly chooses candidate locations, to avoid uncompressing large suffix array intervals and executing many DP instances. The tool verifies candidate locations using a striped vectorial dynamic programming algorithm by Farrar [2007], implemented using SIMD instructions. Bowtie 2 can be configured to report semi-global or local alignments, scored using a tunable affine scoring scheme.

### 5.4.2 BWA

BWA [Li and Durbin, 2009] is designed to map Illumina reads and report a few best semi-global alignments. The program backtracks the FM-index of the reference genome with a *greedy breadth-first search*. The tool ranks nodes to be visited by edit distance score: the best node is popped from a priority queue and visited, its children are then inserted again in the queue. The traversal considers indels using a more involved 9-fold recursion. Li and Durbin speed up backtracking by adopting a more stringent pruning strategy [Mäkinen *et al.*, 2010] that nonetheless takes some preprocessing time and requires the index of the reverse reference genome. BWA performs paired-end alignments by trying to anchor both paired-end reads and verifying the corresponding mate, within an estimated insert size, using the classic DP-based algorithm by Smith and Waterman [1981]. Consequently, the program in paired-end mode aligns reads at a slower rate than in single-end mode. The program is not fully multi-threaded, therefore it scales poorly on modern multi-core machines.

### 5.4.3 Soap

Soap 2 [Li *et al.*, 2009b] has been designed to produce a very quick but shallow mapping of short Illumina reads, up to 2 mismatches and without indels. The tool performs backtracking using the so-called bi-directional (or 2-way) BWT [Belazzougui *et al.*, 2013]. Soap 2 supports paired-end mapping but at a slower alignment rate, it lacks native output in the *de-facto* standard SAM format, and it is not open source.

### 5.4.4 SHRiMP 2

The *SHort Read Mapping Program* (SHRiMP 2) [David *et al.*, 2011] is designed to map short Illumina and ABI/SOLiD reads. The tool achieves high accuracy at the expense of speed. SHRiMP 2 indices the reference genome using multiple gapped  $q$ -grams. At query time, it projects each read to identify candidate mapping locations, which are verified with a DP algorithm [Smith and Waterman, 1981]. The SHRiMP 2 project has been recently discontinued.

### 5.4.5 RazerS and RazerS 3

RazerS [Weese *et al.*, 2009] has been designed to report all mapping locations within a fixed Hamming or edit distance error rate. It is based on a full-sensitive  $q$ -gram counting filtration method (see section 4.3) combined with the edit distance verification algorithm by Myers [Myers, 1999]. On demand, the tool throttles filtration to be more specific at the expense of a controlled loss rate. Stronger filtration reduces the number of candidate locations and improves the overall speed of the program. All in all, the SWIFT filter is very slow while not highly specific.

RazerS 3 [Weese *et al.*, 2012] is a faster version featuring shared-memory parallelism, a banded version of Myers' algorithm, and a quicker filtration method based on exact seeds (see section 4.1). Such filtration method however turns out to be very weak on

mammal genomes. Because of this fact, RazerS 3 is one-two orders of magnitude slower than Bowtie 2 and BWA on such datasets.

All RazerS versions index the reads and scan the reference genome. One positive aspect of this strategy is that no preprocessing of the reference genome is required. However, other mapping strategies beyond all-mapping, e.g. mapping by strata, cannot be efficiently implemented. Moreover, these programs exhibit an high memory footprint as they remember the mapping locations of all input reads until the whole reference genome has been scanned.

#### 5.4.6 mrFast and mrsFast

The tools mrsFast [Hach *et al.*, 2010] and mrFast [Ahmadi *et al.*, 2012] are designed to map Illumina reads. They report all mapping locations within a fixed absolute number errors, respectively under the edit and Hamming distance. Similarly to RazerS 3, these two programs implement full-sensitive filtration using exact seeds (section 4.1). Their peculiarity is a cache-oblivious strategy to mitigate the high cost of verifying clusters of candidate locations. In addition, mrsFast computes the edit distance between one read and one mapping location in the reference genome with an antidiagonal-wise vectorial dynamic programming algorithm, implemented using SIMD instructions. These tools perform only all-mapping, produce files of impractical size and lack multi-threading support.

#### 5.4.7 Hobbes 2

Hobbes 2 [Kim *et al.*, 2014] is designed to identify all read mapping locations within a fixed Hamming or edit distance threshold. In order to improve filtering efficiency, the tool employs a novel technique of so-called *prefix q-grams* that enriches the reference genome *q*-gram index. However, this technique does not guarantee full-sensitivity.

#### 5.4.8 GEM

The GEM mapper [Marco-Sola *et al.*, 2012] is a flexible read aligner for Illumina and Ion Torrent reads. The tool can be configured either as an all-mapper, as a best/unique-mapper, or to search by strata; however, it supports the best-mapping paradigm only to some extent, as it does not annotate mapping locations with qualities.

GEM implements full-sensitive filtration with approximate seeds (see section 4.2). The program indexes the reference genome with an FM-index, tries to find an optimal filtration scheme per read, and verifies candidate locations using Myers' algorithm Myers [1999]. GEM maps paired-reads in two ways: either it maps both ends independently and then combines them, or maps one end and then verifies the other end using an online method. Unfortunately, the tool is not open source and provides obscure parameterization.

Table 5.3: Overview of popular read mappers.

Sequencer			Paradigm			Alignment			Index		
Illumina	Ion		Best	Strata	All	Type	Optimal	Method	Type	Reference	Reads
Bowtie	short	✗	✓	✓	✗	mismatches	✗	backtracking	FM-index	✓	✗
Bowtie 2	✓		✓	✗	✗	local	✗	exact seeds	FM-index	✓	✗
BWA	✓	✗	✓	✗	✗	indels	✗	backtracking	FM-index	✓	✗
Soap 2	short	✗	✓	✗	✗	mismatches	✗	backtracking	FM-index	✓	✗
RazerS	✓	✗	✗	✓	✓	indels	✓	$q$ -grams	$q$ -gram index	✗	✓
RazerS 3	✓	✗	✗	✓	✓	indels	✓	exact seeds	$q$ -gram index	✗	✓
SHRIMP 2	✓	✓	✓	✗	✓	local	✗	$q$ -grams	$q$ -gram index	✓	✗
mrsFast	short	✗	✗	✗	✓	mismatches	✓	exact seeds	$q$ -gram index	✓	✓
mrFast	✓	✗	✗	✗	✓	indels	✓	exact seeds	$q$ -gram index	✓	✓
Hobbes 2	✓	✗	✗	✗	✓	indels	✗	prefix $q$ -grams	$q$ -gram index	✓	✗
GEM	✓	✓	✓	✓	✓	indels	✓	apx seeds	FM-index	✓	✗
Masai	✓	✗	✓	✗	✓	indels	✓	apx seeds	generic	✓	✓
Yara	✓	✓	✓	✓	✓	indels	✓	apx seeds	FM-index	✓	✗

This chapter presents the engineering and evaluation of an efficient all-mapper for Illumina reads. When I started this project, in October 2011, the fastest all-mappers (mrFast and RazerS 3) were two order of magnitude slower than popular best-mappers (Bowtie and BWA). On the one hand, those all-mappers employed filtration based on exact seeds, which is efficient on short reference genomes but becomes too weak on mammal genomes; clearly, a stronger filtration method would have been beneficial. On the other hand, those best-mappers are based on heuristic backtracking, which is inadequate to map longer Illumina reads.

After a thorough literature review, I came out with a novel read mapping method combining seed-based filtration with backtracking, published in the peer-reviewed journal *Nucleic Acids Research* [Siragusa *et al.*, 2013a]. My method is packaged in a C++ tool nicknamed *Masai*, which stands for *multiple backtracking of approximate seeds on a suffix array index*. Masai is part of the SeqAn library, it is distributed under the BSD license and can be downloaded from <http://www.seqan.de/projects/masai>.

In the engineering section, I explain Masai's indexing, filtration and verification methods for all-mapping. In particular, the filtration method is based on approximate seeds: by employing approximate seeds instead of exact seeds, the tool obtains stronger, non-heuristic and quasi full-sensitive filtration for mammal reference genomes. Masai finds approximate seeds by backtracking the index of the reference genome. The tool speeds up the backtracking phase by searching all seeds simultaneously, with the help of an additional index and the multiple backtracking algorithm. Lastly, Masai implements also a quicker best-mapping method, though without mapping qualities.

In the evaluation section, I extensively compare Masai with popular read mappers, both on simulated and real datasets. Compared to the all-mappers mrFast and RazerS 3, Masai is an order of magnitude faster and has comparable sensitivity. In addition, Masai in best-mapping is 2–4 times faster and more accurate than Bowtie 2 [Langmead and Salzberg, 2012] and BWA [Li and Durbin, 2009]. Finally, I discuss the limitations of Masai that led me to engineer yet another read aligner.

## 6.1 Engineering

I first give an outline of the read mapping method implemented in Masai. Then, I explain each step in details, motivating relevant engineering choices that led me to the final implementation.

Masai requires an index capable of simulating a top-down traversal of the suffix trie of the reference genome. The tool gives to users the possibility to choose among various indices (see section 6.1.2). Similarly to all read mappers relying on an index of the reference genome, the tool indexes the reference genome only once, stores it on disk and reuses it for all subsequent read mapping jobs.

At mapping time, Masai requires two parameters to be provided: a maximum number of errors per read and a minimum seed length. Default parameters work well for actual Illumina reads, otherwise the user has to adequately parametrize the tool for optimal performance. Nonetheless, independently of the chosen parameterization, filtration is guaranteed to be quasi full-sensitive (see section 6.1.1).

Masai partitions all reads (and their reverse complements) into non-overlapping seeds and indexes them in a conceptual *trie*. Using the *multiple backtracking* algorithm explained in section 3.3.6, the tool backtracks simultaneously all indexed seeds in the suffix trie of the reference genome. The program verifies all candidate locations, reported by the multiple backtracking algorithm, performing seed extension with a banded version of *Myers bit-vector algorithm* [Myers, 1999] (details in section 6.1.3).

### 6.1.1 Filtration

My original intent was to improve the speed of the all-mapper RazerS [Weese *et al.*, 2009] while preserving full-sensitivity under the edit distance. RazerS was based on a  $q$ -gram filter; I was aware that gapped  $q$ -grams could have brought a huge speedup, but I could not see any straightforward generalization of gapped  $q$ -grams to the edit distance. At the same time, I experienced that weaker but quicker filtration using exact seeds is more advantageous than filtration using  $q$ -grams. Indeed, a typical Illumina read mapping setup requires only moderate error rates, in the range of 4–6 %. For instance, RazerS 3 [Weese *et al.*, 2012] went back to filtration with exact seeds (similarly to mrFast). Nonetheless, I wanted to improve filtration specificity of exact seeds, as the runtime of RazerS 3 on mammal genomes became dominated by verifications. I knew that to improve filtration specificity I had to increase the seed length.

While reviewing past literature in the field of approximate string matching, I rediscovered the works of Myers [1994] and Navarro and Baeza-Yates [2000] on approximate seeds, providing stronger filtration than exact seeds while preserving full-sensitivity under the edit distance. Their idea is to partition the pattern into fewer non-overlapping seeds, which obviously can be longer than exact seeds but have to be searched approximately. First, I slightly improved the filtration lemma of [Navarro and Baeza-Yates, 2000] to use approximate seeds with variable thresholds (see section 4.2). Then, as discussed in section 4.2, I chose to parameterize Masai’s filter by the seed length rather than by the number of seeds. Indeed, the minimum seed length provides a direct estimate of the expected number of candidate locations to verify and thus of filtration specificity. The resulting filter is thus flexible: by increasing the seed length, filtration becomes more specific at the expense of a higher filtration time. In practice, the optimal seed length depends on the reference genome as well as on read lengths and the absolute number of errors. In section 6.2.5, I experimentally evaluate filtration schemes with exact and

approximate seeds. When mapping current Illumina reads on short to medium length genomes, exact seeds are still more efficient than approximate seeds. Conversely, on larger genomes (e.g. mammalian genomes) 1-approximate seeds outperform exact seeds by an order of magnitude.

Following the work of Navarro and Baeza-Yates [2000], I decided to find approximate seeds by backtracking the suffix trie of the reference genome. In section 6.1.2, I recall the engineering work I did to implement efficient backtracking of approximate seeds. In order to achieve faster filtration, I opted to find approximate seeds only under the Hamming distance. For this reason, when resorting to approximate seeds, Masai does not attain strict full-sensitivity under the edit distance. Nonetheless, such implementation choice sacrifices only 0.1% sensitivity (see section 6.2).

### **Best-mapping**

Masai is a tool primarily designed to perform all-mapping rather than best-mapping. In best-mapping, Masai simply reports the first optimal mapping location encountered per read. Clearly, this policy makes sense if the edit distance is effective at identifying the original mapping locations. The evaluation of section 6.2.1 shows that Masai is competitive in best-mapping with tools using more complex scoring schemes. Nonetheless, Masai's best-mapping method is ad-hoc and limited. In best-mapping, the tool does not compute mapping qualities nor supports the paired-end and mate-pair protocols. The reader is thus referred to the next chapter (section 7.1) for the complete description of an efficient best-mapping method that, in standard scenarios, is an order of magnitude faster than all-mapping.

### **6.1.2 Indexing**

Initially, the SeqAn library provided me with only two indices capable of simulating a top-down suffix trie traversal: the enhanced suffix array (ESA) [Abouelhoda *et al.*, 2004] and the lazy suffix tree (LST) [Giegerich *et al.*, 1999]. To improve the efficiency of Masai, I implemented a generic top-down traversal for some additional indices, namely the suffix array (SA) [Manber and Myers, 1990], the  $q$ -gram index, and various specializations of the full-text minute index (FM-index) [Ferragina and Manzini, 2001]. Below I discuss the performance of these indices within Masai, while I refer the reader to chapter 3 for their extensive explanation.

#### **Indexing the reference genome**

I initially chose the ESA over the LST because of better construction times. Indeed, the ESA provides a linear time construction algorithm (an adaptation of the DC7 algorithm [Dementiev *et al.*, 2008] to multiple sequences [Weese, 2013] for the generalized SA, followed by the algorithms proposed in [Kasai *et al.*, 2001; Abouelhoda *et al.*, 2004]), while the LST construction algorithm takes quadratic time (using the radix sort based *wotd*-algorithm [Giegerich *et al.*, 1999]). The construction of the ESA of the *H. sapiens* reference genome (GRCh37) takes about 1.5 hours and the index consumes 39 GB of memory. Apart

from that, both ESA and LST implementations require 13 bytes per base pair and exhibit comparable query speed.

At this point, Masai required high-end hardware to process large reference genomes. Therefore, thinking about a space-time trade-off, I designed a generic suffix trie top-down traversal for the SA (see section 3.1.1). The SA consumes only 5 bytes per base pair but is theoretically slower than the ESA, as it adds a logarithmic factor to query times. However, with surprise, I found out that within Masai the SA had equal or better performance than the ESA (see table 6.4). Ultimately, this change brought down the memory footprint of the index from 39 GB to 15 GB but preserved query speed.

I tried to further improve query speed by removing the logarithmic factor introduced by the SA. Therefore, to cut the most expensive binary searches, I put a  $q$ -gram index on top of the SA and extended my generic suffix trie top-down traversal accordingly (see section 3.1.3). Yet, the  $q$ -gram index did not bring significant speedup to the application; indeed, the lookup table turned out to be useful when searching patterns one by one, but not when coupled with the multiple backtracking algorithm.

Finally, I explored additional space-time trade-offs. Starting from the implementation of [Singer, 2013], I realized a generalized FM-index based on a wavelet tree [Grossi *et al.*, 2003]. This initial FM-index consumed about 1.5 bytes per base pair with a SA sampling of 10 %. Thus the memory footprint of the index went down to 4.5 GB, but Masai became almost twice as slow (see table 6.4).

To sum up, I preferred the SA as it provides a good compromise between query speed and memory consumption. Nevertheless, Masai leaves to the user the possibility of choosing among the aforementioned data structures. Table 6.4 summarizes the runtime of the program with various indices.

## Indexing the reads

In order to improve index query speed, I designed and implemented the algorithms presented in sections 3.3.5 and 3.3.6. These algorithms search simultaneously many exact or approximate seeds, achieving a speedup of 2–5 times over their conventional counterparts. As this multiple string matching algorithm requires a trie of the seeds, I also engineered an efficient trie implementation. A short explanation can be found in section 3.1.4

Building the SA via quicksort turned out to be faster than building the LST via radix sort but, within Masai, the more involved LST data structure paid off in terms of query time (see section 3.3.2). Indeed, the LST stores all trie nodes and thus provides node traversal in constant time, while the SA explicitly stores only the leaves and then derives internal nodes via binary search. As the memory footprint of the trie is negligible within this application, I chose the LST to perform multiple backtracking of approximate seeds.

When performing multiple backtracking of exact seeds, the LST construction time dominates the overall filtration time (see section 3.3.5). Therefore, I decided to resort to the  $q$ -gram index to emulate a trie in this case: I build a partial  $q$ -gram index efficiently and in linear time by bucket sort, again considering only the first suffix of each seed in the collection. Such index represents a trie truncated at depth  $q$  (which I fixed to 12 in



Masai). Truncation is only a minor concern: at depth  $q$  the search continues separately on each active node using the conventional binary search algorithms (see sections 3.3.3 and 3.3.4).

### 6.1.3 Verification

To verify candidate locations reported by the filtration algorithm, Masai employs a banded version of Myers bit-vector algorithm [Myers, 1999] implemented for the tool RazerS 3 by Weese *et al.* and available in the SeqAn library. This banded version computes only a diagonal band of the DP matrix and is thus faster than the original algorithm by Myers. However, Masai uses this algorithm in a different way than RazerS 3. While RazerS 3 performs one semi-global alignment to verify one parallelogram surrounding any seed occurrence, Masai performs two global alignments on both ends of any seed occurrence. Given a seed occurring with  $e$  errors, Masai first performs seed extension on the left side within an error threshold of  $k - e$  errors. Only if the seed extension on the left side succeeds, the tool performs a seed extension on the right side within the remaining error threshold. Moreover, the tool first computes the *longest common prefix* on each side of the seed and lets the global alignment algorithm start from the first mismatching positions. This approach is up to two times faster than the one implemented by RazerS 3 (data not shown).

## 6.2 Evaluation

In order to evaluate Masai, I propose three experiments: (i) the Rabema benchmark on simulated data, (ii) variant detection on simulated data, and (iii) performance on real data. This evaluation focuses on the capability of the mappers to retrieve the location of a single read without the help of its paired-end, which can of course disambiguate some mapping locations. As references, I use whole genomes of *E. coli* (NCBI NC\_000913.2), *C. elegans* (WormBase WS195), *D. melanogaster* (FlyBase release 5.42), and *H. sapiens* (GRCh37.p2).

I compare Masai in all-mapping with RazerS 3, Hobbes, mrFAST and SHRiMP 2, while in best-mapping with Bowtie 2, BWA and Soap 2. Masai, RazerS 3, Hobbes and mrFAST use edit distance, while Bowtie 2, BWA, Soap 2 and SHRiMP 2 rely on scoring schemes taking into account base quality values. When relevant, I configured some read mappers with the appropriate absolute number of errors (Masai, mrFAST, Hobbes, Soap 2) or error rate (RazerS 3). In section A.1, I give the exact parameterization of the read mappers considered in this evaluation.

### 6.2.1 Rabema benchmark on simulated data

I first consider the Rabema benchmark [Holtgrewe *et al.*, 2011] (v1.1) for a thorough evaluation and comparison of read mapping sensitivity. I consider the benchmark categories *all* and *best*, in addition to *precision* and *recall*. In the categories *all* and *best*, a

read mapper has to find for each read respectively all or one of the best edit distance mapping locations. The categories precision and recall require a read mapper to find the *original* location of each simulated read, which is a measure independent of the used scoring model, e.g. edit distance or quality based. A simulated read is mapped *correctly* if the mapper reports its original location, and it is mapped *uniquely* if the mapper reports only one location. Rabema defines *recall* to be the fraction of reads which were correctly mapped and *precision* the fraction of uniquely mapped reads that were mapped correctly.

Similarly to [Langmead and Salzberg, 2012], I used the read simulator Mason [Holtgrewe, 2010] with default profile settings to simulate, from each whole genome, 100 k reads of length 100 bp having sequencing errors distributed like in a typical Illumina run. I performed the benchmark for an error rate of 5 %, which corresponds to edit distance 5 for reads of length 100 bp. Therefore, I built a Rabema gold standard for each dataset by running RazerS 3 in full-sensitive mode up to edit distance 5. I further classified mapping locations in each category by their edit distance.

For a more fair and thorough comparison, I also consider BWA and Bowtie 2 in all-mapping (Soap 2 cannot be configured accordingly). To this extent, I parametrized these tools to be highly sensitive and output all found mapping locations. Since BWA and Bowtie 2 were not designed to be used in this way, they spent much more time than proper all-mappers, i.e. up to 3 hours in a run compared to several minutes. However, the aim of this experiment is to investigate read mapping sensitivity, therefore I do not report any running times. Table 6.1 shows the results on *H. sapiens*.

### All-mapping

As expected, RazerS 3 shows full-sensitivity. In contrast, mrFAST loses a minimal percentage of mapping locations. Overall, Masai does not lose more than 0.1 % of all mapping locations. In particular, Masai is full-sensitive for low-error locations and loses only a small percentage of high-error locations, i.e. its loss is limited to 0.1 % and 1.4 % of mapping locations at edit distance 4 and 5.

Conversely, BWA and Bowtie 2 miss 35 % and 45 % of all mapping locations at edit distance 5 and their recall values as all-mappers do not substantially increase. Likewise, SHRiMP 2 is not able to enumerate all mapping locations, although its recall values are good. Again, Hobbes has the worst performance.

As mentioned in section 6.1.1, Masai is not full-sensitive whenever approximate seeds are used, e.g. on *H. sapiens*. Indeed, Masai loses 0.1 % overall sensitivity in respect to RazerS 3. In general, RazerS 3 should be used when full-sensitivity is required, i.e. in read mapping benchmarks. However, these results show that Masai can replace RazerS 3 or mrFAST in practical all-mapping setups.

### Best-mapping

Masai shows the best recall values, not losing more than 2.3 % recall on edit distance 5. Conversely, the recall values of BWA and Bowtie 2 drop significantly with increasing edit distance and lose up to 15.4 % and 11.5 % on edit distance 5. As expected, Soap 2 turns out to be inadequate for mapping reads of length 100 bp at this error rates. Precision

**Table 6.1:** Rabema benchmark results on  $100k \times 100bp$  Illumina-like reads. Rabema scores are given in percent (average fraction of edit distance locations reported per read). Large numbers show total scores in each Rabema category and small numbers show the category scores separately for reads with  $\begin{pmatrix} 0 & 1 & 2 \\ 3 & 4 & 5 \end{pmatrix}$  errors.

		All locations				Best locations				Recall				Precision			
All-mapping	Masai	99.9	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0
	mrFAST	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0
	RazerS 3	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0
	Bowtie 2	95.7	100.0	99.9	99.4	99.2	100.0	100.0	99.0	98.5	99.7	99.6	98.3	99.8	100.0	99.9	99.9
	BWA	95.9	100.0	99.9	99.5	98.8	100.0	100.0	99.8	97.8	99.0	99.0	98.7	98.1	93.2	97.6	98.4
Best-mapping	Masai	93.3	99.2	98.7	97.9	100.0	100.0	100.0	100.0	97.7	97.9	97.8	97.8	97.8	97.9	97.8	97.8
	Bowtie 2	92.0	99.2	98.7	96.8	98.1	100.0	100.0	97.6	95.9	98.0	97.7	95.6	96.6	98.0	97.7	96.0
	BWA	92.2	99.2	98.7	97.8	98.8	100.0	100.0	99.8	96.4	97.9	97.7	97.3	97.5	97.9	97.7	97.4
	Soap 2	65.9	99.2	95.5	91.3	71.4	100.0	96.8	93.2	69.9	98.1	94.6	91.2	97.7	98.1	97.7	97.7
			8.7	0.7	0.0		9.2	0.8	0.0		11.9	1.4	0.4		94.9	84.1	91.7

values have less variance than recall values. Masai shows the best precision values with 97.8 %, followed by Soap 2 with 97.7 %, and BWA with 97.5 %. Interestingly, Bowtie 2 shows the worst precision values, losing up to 5.6 % on edit distance 5.

## 6.2.2 Variant detection on simulated data

The second experiment analyzes the theoretical performance of Masai and other read mappers in variant detection pipelines. Similarly to [David *et al.*, 2011], this experiment considers simulated reads containing sequencing errors, SNPs and indels. Each simulated read has an edit distance of at most 5 to its genomic origin, and it is grouped according to the number of contained SNPs and indels, where class  $(s, i)$  consists of all reads with  $s$  SNPs and  $i$  indels. The experiment considers a read to be mapped *correctly* if a mapping location is reported within 10 bp of its genomic origin; it considers a read to map *uniquely* if only one location is reported by the mapper. For each class, the experiment defines *recall* to be the fraction of reads which were correctly mapped and *precision* the fraction of uniquely mapped reads that were mapped correctly.

I simulated 5 million Illumina-like reads of length 100 bp from the whole human genome using Mason. I mapped the reads with each tool and measured its sensitivity in each class. Table 6.2 shows the results.

### All-mapping

Looking at all-mapping results, Masai shows 100 % precision and recall in all classes, except for classes (2,0) and (1,1) where it loses only 0.1 % and 0.7 % recall. Masai is therefore roughly comparable to the full-sensitive read mappers RazerS 3 and mrFAST.

**Table 6.2:** Variant detection results on  $5\text{ M} \times 100\text{ bp}$  Illumina-like reads. The table shows percentages of found origins (recall) and fraction of unique reads mapped to their origin (precision) classed by reads with  $s$  SNPs and  $i$  indels ( $s, i$ ).

		(0,0)		(2,0)		(4,0)		(1,1)		(1,2)		(0,3)	
		Prec.	Recl.	Prec.	Recl.	Prec.	Recl.	Prec.	Recl.	Prec.	Recl.	Prec.	Recl.
All-mapping	Masai	100.0	100.0	100.0	99.9	100.0	100.0	100.0	99.3	100.0	100.0	100.0	100.0
	RazerS 3	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0
	Hobbes	99.9	99.9	99.9	99.9	100.0	100.0	100.0	99.8	100.0	93.6	99.6	90.5
	mrFAST	100.0	99.9	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0
	SHRiMP 2	100.0	99.4	100.0	99.7	100.0	99.7	100.0	99.5	100.0	99.2	100.0	99.6
Best-mapping	Masai	98.2	98.2	97.6	97.5	96.8	96.8	97.8	97.2	97.9	97.9	97.2	97.2
	Bowtie 2	97.6	97.3	94.6	92.0	92.6	82.5	95.3	93.3	93.5	92.3	96.1	95.4
	BWA	98.2	97.9	97.6	95.3	94.9	85.1	97.4	90.9	97.1	80.3	96.3	66.5
	Soap 2	98.1	82.9	97.4	31.0	0.0	0.0	90.6	6.2	0.0	0.0	0.0	0.0

SHRiMP 2 shows 100 % precision in all classes but loses between 0.3 % and 0.8 % recall in each class. Hobbes has the lowest performance among all-mappers: it appears to have problems with indels, indeed it loses 9.5 % recall in class (0,3).

### Best-mapping

Masai shows the highest precision and recall in all best-mapping classes. In particular, Masai does not lose more than 3.2 % recall in class (4,0), whereas Bowtie 2 and BWA lose respectively 17.5 % and 14.9 % and Soap 2 is not able to map any read. The recall values of Bowtie 2, BWA and Soap 2 are negatively correlated with the amount of genomic variation. For instance, in the Rabema benchmark, Bowtie 2 loses respectively 7.2 % and 11.5 % of mapping locations at distance 4 and 5, but in class (4,0) of this experiment it loses 17.5 % recall. A similar trend is observable for BWA and Soap 2. The low performance of Soap 2 is also due to its limitation to at most 2 mismatches and no support for indels.

### 6.2.3 Performance on real data

The last experiment focuses on comparing read mappers performance on real data. I mapped the first  $10\text{ M} \times 100\text{ bp}$  reads from an Illumina lane of *E. coli* (ERR022075, Genome Analyzer IIx), *D. melanogaster* (SRR497711, HiSeq 2000), *C. elegans* (SRR065390, Genome Analyzer II), and *H. sapiens* (ERR012100, Genome Analyzer II). Whenever possible, I configured the tools to map the reads within edit distance 5. I measured mapping times on a cluster of nodes with 72 GB RAM and 2 Intel Xeon X5650 processors running Linux 3.2.0. For an accurate running time comparison, I ran the tools using a single thread and used local disks for I/O. I measured running times and peak memory consumptions.

I cannot measure precision and recall values as real reads have unknown origins. Therefore, for this evaluation, I count also the Rabema *best* locations. As explained in section 6.2.1, the Rabema best category assigns a point for a read if the mapper reports at least one mapping location with the optimal (i.e. best) number of errors; final Rabema best scores are normalized by the number of reads. All results are shown in table 6.3.

### All-mapping

On the *H. sapiens* dataset, Masai is 11.9 times faster than RazerS 3, 14.6 times faster than mrFAST, and 7.6 times faster than Hobbes. The current version of Hobbes constantly crashes and maps only half of the reads. SHRiMP 2 is not able to map the *H. sapiens* dataset within 4 days. On the *C. elegans* dataset, Masai is 2.0 times faster than RazerS 3, 10.9 times faster than Hobbes, 6.3 times faster than mrFAST and 50.1 times faster than SHRiMP 2. Hobbes constantly crashes and maps less reads than all other mappers in this category.

### Best-mapping

On the *H. sapiens* dataset, Masai is 2.6 times faster than Bowtie 2, 3.6 times faster than BWA but 2.1 times slower than Soap 2. On the *C. elegans* dataset, Masai is 7.7 times faster than Bowtie 2, 8.2 times faster than BWA and 1.5 times faster than Soap 2. Soap 2 is not able to map a consistent fraction of reads because of its limitation to 2 mismatches. Bowtie 2 misses respectively 22.0 % and 20.7 % of reads mappable at edit distance 5 on the *C. elegans* and *H. sapiens* datasets.

## 6.2.4 Performance with different indices

Table 6.4 shows the runtime and memory consumption of Masai using different indices. As discussed in section 6.1.2, Masai/SA is always faster than Masai/ESA, except for all-mapping on the *H. sapiens* dataset. Masai/FM-index shows significantly lower memory consumption only on the *H. sapiens* dataset, but with a consistent runtime penalty.

## 6.2.5 Filtration efficiency

This experiment assesses the contribution of approximate seeds and multiple backtracking on runtime results. To this intent, I perform all-mapping with Masai on all dataset of section 6.2.3, using either exact or approximate seeds in combination with either single or multiple backtracking. As this experiment focuses on filtration efficiency, I do not consider the time spent performing seed extensions and I/O, i.e. loading the reference genome and its index, loading the reads, and writing the results.

Table 6.5 shows the results. Column *time* reports filtration times, i.e. the time spent to index the seeds (in case of multiple backtracking) and to perform backtracking. Column *candid.* reports the number of candidate locations reported by the filter for which seed extension is subsequently performed. In bold, I report the optimal filtering scheme used to parameterize Masai in the experiments of section 6.2.3.

**Table 6.3:** Performance on real data using  $10\text{ M} \times 100\text{ bp}$  Illumina reads. Rabema best column: in large are shown the percentage of reads mapped with the minimal number of errors (up to 5%) and in small the percentage of reads that were mapped with  $\begin{pmatrix} 0 & 1\% & 2\% \\ 3\% & 4\% & 5\% \end{pmatrix}$  errors. Remarks: SHRiMP 2 is not able to map the *H. sapiens* dataset within 4 days; Hobbes systematically crashes in all but the *D. melanogaster* dataset.

		ERR012100 <i>H. sapiens</i>			SRR065390 <i>C. elegans</i>		
		Time [min:s]	Memory [GB]	Best locations [%]	Time [min:s]	Memory [GB]	Best locations [%]
All-mapping	Masai	307:16	19.66	100.0 100.0 100.0 100.0 100.0 100.0 99.5	10:49	2.76	100.0 100.0 100.0 100.0 100.0 100.0 100.0
	RazerS 3	3653:03	16.89	100.0 100.0 100.0 100.0 100.0 100.0 100.0	21:18	11.22	100.0 100.0 100.0 100.0 100.0 100.0 100.0
	Hobbes	2319:27	70.00	59.0 59.2 58.7 57.5 56.9 56.7 56.3	117:46	3.79	89.8 91.0 80.6 86.5 88.3 88.5 85.2
	mrFAST	4462:25	0.91	100.0 100.0 100.0 100.0 100.0 100.0 97.5	67:41	0.85	100.0 100.0 100.0 99.9 99.9 99.9 99.5
	SHRiMP 2	–	–	–	541:20	2.67	98.5 99.6 96.8 91.8 87.6 81.9 74.8
Best-mapping	Masai	22:35	19.25	100.0 100.0 100.0 100.0 99.9 99.9 99.5	3:10	2.87	100.0 100.0 100.0 100.0 100.0 100.0 100.0
	Bowtie 2	57:41	3.11	99.4 100.0 99.7 96.0 92.9 87.9 79.3	24:14	0.13	99.2 100.0 99.3 93.4 88.6 84.0 78.0
	BWA	80:58	4.37	99.5 100.0 99.5 98.0 93.4 88.9 84.4	25:53	0.32	99.3 100.0 99.1 95.6 89.7 85.9 82.3
	Soap 2	11:11	5.23	95.7 100.0 94.9 86.5 0.3 0.2 0.2	4:37	0.73	96.0 100.0 96.6 92.4 0.3 0.04 0.02
		SRR497711 <i>D. melanogaster</i>			ERR022075 <i>E. coli</i>		
		Time [min:s]	Memory [GB]	Best locations [%]	Time [min:s]	Memory [GB]	Best locations [%]
All-mapping	Masai	7:34	2.87	100.0 100.0 100.0 100.0 100.0 100.0 100.0	1:33	2.22	100.0 100.0 100.0 100.0 100.0 100.0 100.0
	RazerS 3	10:54	7.71	100.0 100.0 100.0 100.0 100.0 100.0 100.0	1:46	5.57	100.0 100.0 100.0 100.0 100.0 100.0 100.0
	Hobbes	42:05	2.44	99.9 100.0 100.0 100.0 100.0 99.3 96.6	9:14	0.68	95.1 95.1 95.2 95.1 95.1 95.3 94.8
	mrFAST	37:38	0.88	99.9 100.0 100.0 100.0 100.0 100.0 96.9	4:34	0.67	100.0 100.0 100.0 100.0 100.0 100.0 100.0
	SHRiMP 2	225:03	3.02	99.7 100.0 100.0 99.7 98.7 96.3 92.6	41:40	0.95	99.8 100.0 99.7 98.7 97.1 94.5 91.2
Best-mapping	Masai	4:52	2.98	100.0 100.0 100.0 100.0 100.0 99.8 99.0	0:44	2.33	100.0 100.0 100.0 100.0 100.0 99.5 97.5
	Bowtie 2	21:11	0.15	99.5 97.6 94.9 89.9	12:11	0.03	99.7 94.4 90.7 85.8
	BWA	30:46	0.35	98.5 100.0 99.6 98.4 90.7 82.1 73.5	10:13	0.16	99.7 100.0 99.5 97.5 94.4 91.8 89.5
	Soap 2	5:42	0.76	90.3 100.0 96.2 89.4 0.09 0.02 0.02	2:19	0.59	97.6 100.0 99.1 96.8 0.5 0.1 0.08

On *E. coli*, *D. melanogaster* and *C. elegans*, approximate seeds reduce the number of candidates respectively by 2.1 times, 9.9 times, and 4.3 times. Nevertheless I still prefer exact seeds as filtration dominates the total runtime. Multiple backtracking on exact seeds compared to single backtracking speeds up filtration by 2.9 times on *E. coli*, and

**Table 6.4:** Masai performance with different indices. The lowest runtimes are in bold.

		ERR022075 E. coli		SRR497711 D. melanogaster		SRR065390 C. elegans		ERR012100 H. sapiens	
		Time [min:s]	Memory [GB]	Time [min:s]	Memory [GB]	Time [min:s]	Memory [GB]	Time [min:s]	Memory [GB]
All-mapping	SA	<b>1:33</b>	2.22	<b>7:34</b>	2.87	<b>10:49</b>	2.76	307:16	19.66
	Esa	1:42	2.26	8:02	3.76	11:13	3.50	<b>297:13</b>	42.18
	FM-index	2:38	2.20	17:15	2.44	21:12	2.40	480:23	8.94
Best-mapping	SA	<b>0:44</b>	2.33	<b>4:52</b>	2.98	<b>3:10</b>	2.87	<b>22:35</b>	19.25
	Esa	0:46	2.37	5:01	3.88	3:15	3.61	26:24	41.78
	FM-index	0:52	2.31	10:17	2.55	5:29	2.51	42:26	8.57

3.8 times on *D. melanogaster* and *C. elegans*. Without the contribution of multiple backtracking, Masai would not be faster than RazerS 3, the second fastest all-mapper.

Approximate seeds become effective on *H. sapiens*, where they reduce the number of candidates by 10.8 times. On *H. sapiens*, seed extensions largely dominate the total runtime, therefore I prefer approximate seeds. Multiple backtracking on approximate seeds provides a speed-up of 3.2 times over single backtracking. The combination of the two methods makes Masai an order of magnitude faster than any other all-mapper.

### 6.3 Discussion

Masai consists of two important algorithmic methods: approximate seeds and multiple backtracking. Approximate seeds are of paramount importance to obtain very specific, yet full-sensitive filtration; their adoption speeds up Masai by one order of magnitude. Multiple backtracking further speeds up the filtration phase by 3–5 times on a

**Table 6.5:** Masai all-mapping filtration efficiency results. Filtering times include seeds indexing times. The best filtering schemes are in bold.

		ERR022075 E. coli		SRR497711 D. melanogaster		SRR065390 C. elegans		ERR012100 H. sapiens	
		Timee [min:s]	Candid. [M]	Timee [min:s]	Candid. [M]	Timee [min:s]	Candid. [M]	Timee [min:s]	Candid. [M]
Seeding	Backtracking								
Exact	Single	3:55	69.17	8:15	1020.28	8:25	1065.70	55:54	294943.86
Exact	Multiple	<b>1:20</b>	<b>69.17</b>	<b>2:11</b>	<b>1020.28</b>	<b>2:11</b>	<b>1065.70</b>	41:52	294943.86
1-Apx	Single	38:42	33.08	100:18	102.78	102:02	246.65	165:45	27396.01
1-Apx	Multiple	9:00	33.08	20:48	102.78	21:33	246.65	<b>52:15</b>	<b>27396.01</b>

(enhanced) suffix array index; this technique makes Masai twice as fast. In addition, Masai implements a best-mapping method that finds one optimal mapping location and is an order of magnitude faster than all-mapping.

Is the edit distance adequate for best-mapping? Both Rabema benchmark and variant detection results show that Masai has constantly better accuracy than other best-mappers relying on more complex scoring schemes. In particular, the Rabema benchmark results show that Rabema best values are tightly bound to recall scores. Hence, the edit distance is a pertinent and adequate scoring scheme for best-mapping. Vice versa, best-mappers using scoring schemes based on quality values show a generalized and substantial loss of mapping accuracy. This is likely due to the heuristics on which these tools rely. To sum up, it is better to stick to edit distance and guarantee full-sensitivity rather than to adopt an involved scoring scheme and explore the alignment space heuristically, hence partially.

How many mapping locations do heuristic best-mappers miss? By looking at precision and recall values on simulated data, or at Rabema best values on real data, it can be deduced that Bowtie 2, BWA and Soap 2 miss up to 20 % of reads mappable at 5 % error rate. Yet, it is not evident how these results affects variant calling pipelines.

Summing up, Masai in all-mapping is an order of magnitude faster and thus a valid alternative to tools like RazerS 3 and mrFast. Computational requirements of all-mapping are now close to those of best-mapping. Indeed, Masai in all-mapping is only 4 times slower than BWA in best-mapping, despite reporting two orders of magnitude more mapping locations. Masai in best-mapping is 2–4 times faster and more accurate than Bowtie 2 [Langmead and Salzberg, 2012] and BWA [Li and Durbin, 2009]. The achieved speedup is huge when RazerS 3 is used for best-mapping: in this scenario, Masai is roughly 200 times faster!

Despite these good results, Masai is not being widely used. This is mainly because the tool lacks some commonly requested features, including: direct support of the paired-end and mate-pair protocols, computation of mapping qualities, parallelization via multi-threading, low memory footprint, and automatic parameterization. Because of my initial inexperience and unclear or wrong design goals, I neglected these features while engineering Masai. The next chapter introduces *Yara*, a tool implementing these features.



*Yara* is a non-heuristic read mapper capable of quickly reporting all co-optimal or sub-optimal mapping locations within a given error rate. *Yara* works with Illumina or Ion Torrent reads, supports both paired-end and mate-pair protocols, computes accurate mapping qualities, offers parallelization via multi-threading, has a low memory footprint thanks to the FM-index, and does not require ad-hoc parameterization.

## 7.1 Engineering

### 7.1.1 Stratified mapping

*Yara* is based on the concept of *stratified all-mapping*. The all-mapping methods discussed so far consider a set of relevant mapping locations per read. Yet, this definition leaves open what relevant means. In all-mapping under the edit distance, the user defines relevant mapping locations by imposing a distance threshold. Despite being sound, this definition does not work well in practice. On the one hand a very low threshold leaves a consistent fraction of the reads unmapped, on the other hand a moderate threshold produces a deluge of mapping locations for some reads. In practice, at 5 %, error rate, Illumina reads map on average to hundreds of mapping locations on the human genome. It is questionable whether all these locations are relevant for the downstream analysis. Thus, a finer definition of all-mapping relevance is necessary in practice.

Stratification of mapping locations yields an equally sound yet practical definition of all-mapping under the edit distance. The *e*-stratum

$$\mathbb{S}_e = \{(i, j, e) : d_E(G_{i...j}, R) = e\} \quad (7.1)$$

denotes the set of all mapping locations of a read *R* at edit distance *e* from the reference genome *G*. According to the above definition, conventional all-mapping under the edit distance defines the set

$$\mathbb{S} = \mathbb{S}_0 \cup \mathbb{S}_1 \cup \dots \cup \mathbb{S}_k \quad (7.2)$$

as relevant mapping locations within an *absolute* error threshold *k*. Stratified all-mapping refines this definition by considering only mapping locations being co-optimal, or suboptimal up to a certain degree *s*. Formally, if the distance of any optimal mapping location for read *R* is

$$e^* = \min \{e \in [0, k - s] : \mathbb{S}_e \neq \emptyset\} \quad (7.3)$$

stratified all-mapping considers mapping locations

$$\mathbb{S} = \mathbb{S}_{e^*} \cup \dots \cup \mathbb{S}_{e^*+s} \quad (7.4)$$

within a *relative* suboptimality error threshold  $s$  to be relevant.

### Efficient filtration by strata

Yara significantly improves the runtime of stratified all-mapping over conventional all-mapping. Obviously, the most straightforward way to achieve stratified all-mapping consists of performing conventional all-mapping and subsequently filtering out any irrelevant mapping location. For instance, RazerS3 implements this method and gets no speedup. Another naïve method consists of performing up to  $k$  rounds of conventional all-mapping, using filtration schemes with thresholds incrementing from 0 to  $e^* + s$ . It is easy to see that also this method performs redundant computation: the total work for any read that maps at distance greater or equal to  $k$  corresponds to the sum of all  $k$  filtration schemes.

The following stratified all-mapping method, implemented in Yara, guarantees not to perform more work than conventional all-mapping. Indeed, the key idea is to simply reduce any filtration scheme full-sensitive within distance  $k$  to be full-sensitive within distance  $e^* + s$ . Given a filtration scheme  $T = (t_0, \dots, t_{s^*-1})$  full-sensitive within distance  $k$ , any subset consisting of  $s^* \leq s$  seeds with thresholds  $\mathbf{t}^* = (t_0^*, \dots, t_{s^*-1}^*)$  is full-sensitive within distance  $e^* + s$  if it satisfies

$$s^* + \sum_{i=0}^{s^*-1} t_i^* > e^* + s. \quad (7.5)$$

**Example 7.1.** Let  $k = 5$  be the absolute threshold and  $s = 0$  the relative threshold, defining only co-optimal locations to be relevant. A read  $R$  maps at distance 1, i.e.  $|\mathbb{S}_0| = 0$  and  $|\mathbb{S}_1| > 0$ , thus  $e^* = 1$ . Given the filtration scheme  $T = (1, 1, 1)$ , any subset full-sensitive up to distance  $e^* + s = 1$  finds all relevant mapping locations. All full-sensitive subsets of  $T$  are  $(0, 0, -)$ ,  $(0, -, 0)$ ,  $(-, 0, 0)$ ,  $(1, -, -)$ ,  $(-, 1, -)$ ,  $(-, -, 1)$ , where a  $-$  at position  $i$  indicates that the  $i$ -th seed is unnecessary. Thus, the verification of all candidates, either produced by *any* two exact seeds or by one 1-approximate seed, yields all mapping locations in the 1-stratum of  $R$ .

### Greedy verification strategy

In addition, Yara implements a simple greedy strategy to minimize the number of verifications necessary to find all relevant stratified mapping locations. As candidate locations can be verified in any order, Yara chooses an ordering of the seeds that minimizes the expected number of verifications. The tool first finds all seeds and ranks them by number of candidate locations produced. Then it processes all candidate locations, from the least to the most frequent seed, until it explores  $s$  strata from the first non-empty one.

## Best-mapping

Yara performs best-mapping by means of stratified filtration. Best-mapping requires one primary mapping location along with its confidence. Under the edit distance, without any further assumptions, any co-optimal location is equally likely to be correct. Thus, Yara performs stratified all-mapping with a relative threshold  $s = 0$ , picks one random co-optimal location, and subsequently estimates its mapping quality using all found mapping locations (see section 7.1.5). Thanks to this method, Yara in best-mapping is an order of magnitude faster than in all-mapping.

### 7.1.2 Adaptive filtration

Specific yet rapid filtration is fundamental in the design of an efficient read mapping tool. Read mappers like RazerS3 [Weese *et al.*, 2012] and mrFast [Ahmadi *et al.*, 2012] are designed around naïve filtration with exact seeds. This filtration method is always very quick, however it is not specific enough on long, repetitive reference genomes like the human genome. Masai [Siragusa *et al.*, 2013a] circumvents this problem by enforcing a minimum seed length, whose optimal value must be tuned for a specific reference genome, and eventually by resorting to approximate seeds in order to guarantee full-sensitivity. This filtration method speeds up Masai by an order of magnitude but has some drawbacks: it needs external parametrization, lacks flexibility and is suboptimal in practice.

Yara applies an adaptive filtration scheme *per read* because, under any fixed filtration scheme, the number of verifications per read is not uniform: within a typical human genome resequencing, most reads produce very few verifications and are easily mappable, while few other reads are problematic and often not even confidently mappable to one single location. Consequently, any fixed filtration scheme turns out to be too weak for some reads yet too strong for others, thus suboptimal in practice. An adaptive filtration scheme per read improves filtration efficiency by optimizing the ratio between filtration speed and specificity. Yara thus automatically chooses an adaptive filtration scheme per read, without requiring manual parameterization by the user.

Adaptive filtration works as follows. Yara initially applies filtration with exact seeds to all reads. The tool counts the number of verifications to be performed for each read, thus decides if it is worth proceeding with the verification phase or alternatively applying a stronger filtration scheme. This decision depends on fine-tuned internal verification thresholds. Under standard Illumina setups, exact seeds provide efficient filtration for up to 70–80 % of the reads; on the remaining reads, a filtration scheme using 1- or 2-approximate seeds works better. Thus, Yara starts with the quickest filtration scheme and becomes more specific whenever it pays off to do so.

### 7.1.3 Indexing

Yara uses an efficient FM-index specialization for the DNA alphabet, based on interleaved rank dictionaries (see section 3.2.2). This FM-index exhibits a fourfold speedup over the first FM-index implementation bundled with Masai. Surprisingly, this FM-index is faster than any other index, both in exact and approximate search (see section 3.3). Moreover,

this index consumes only 1.23 bytes per base pair with a SA sampled at 10 %, thus its memory footprint for the human genome is 3.7 GB. Thus, the FM-index always provides the most convenient suffix trie implementation, both in terms of speed and memory footprint.

Yara does not use the multiple search algorithms of section 3.3.6 to search seeds on the FM-index. As shown in section 3.3, on the FM-index it is always faster to search exact queries in a naïve way and approximate queries after sorting them in lexicographical order. This fact considerably simplifies the implementation and allows its fine-grained parallelization.

### 7.1.4 Paired-end and mate-pair protocols

Paired-end and mate-pair protocols are the sequencing protocols of choice of Illumina instruments. As reads are sequenced in pairs from the two ends of the same DNA fragment, they are expected to map closely in the reference genome. The added information of an expected DNA insert size allows to disambiguate the original location of read pairs more confidently than in the single-end protocol. Nonetheless, the lack of any proper pair of mapping locations signals a potential structural variation, e.g. a long indel or an inversion. Therefore, a read mapper should report equally important unpaired mapping locations.

In the paired-end or mate-pair workflows, Yara maps paired reads independently, exactly as in the single-end workflow, and then reports all relevant mapping locations per read. However, in addition to the single-end workflow, Yara implements a finer strategy to choose primary mapping locations. For any reads pair, among all pairs of co-optimal mapping locations, the tool selects the one with minimal deviation from the expected insert size. Since Yara outputs all relevant mapping locations, the choice of primary locations can be always corrected a posteriori.

### 7.1.5 Mapping qualities

Yara computes accurate mapping qualities using the number of found mapping locations stratified by error rate. In section 7.1.1 I defined the set  $\mathbb{S}$  of all relevant mapping locations as

$$\mathbb{S} = \bigcup_{e=0}^k \mathbb{S}_e. \quad (7.6)$$

I now denote by  $z_e$  the number of mapping locations in the stratum  $\mathbb{S}_e$ , i.e.  $z_e = |\mathbb{S}_e|$ ; I associate to  $\mathbb{S}$  a canonical partition function  $Z$

$$Z = \sum_{e=0}^k z_e w(e), \quad (7.7)$$

where the function  $w : \mathbb{N} \rightarrow \mathbb{R}$  assigns a weight to each (sub)optimal stratum. The function  $w$  expresses how well edit distance correlates to the distance of correct map-

ping locations. If all correct mapping locations occur at minimum edit distance  $e^*$ , then  $w(e^*) = 1$  and  $w(e) = 0$  for all  $e > e^*$ .

### Single-end reads

For a fixed single-end read, I denote with  $C$  a random variable assuming the value of its correct mapping location and I assume that

$$\Pr[C \in \mathbb{S}] = 1. \quad (7.8)$$

I estimate the probability of  $C$  being part of the stratum  $\mathbb{S}_e$  as

$$\Pr[C \in \mathbb{S}_e] = \frac{1}{Z} \cdot z_e w(e) \quad (7.9)$$

and the probability of  $C$  being equal to any found mapping location  $(i, j, e)$  as

$$\Pr[C = (i, j, e)] = \frac{1}{Z} \cdot w(e). \quad (7.10)$$

Hence, equation 7.10 gives the probability  $p$  that a random mapping location drawn from any stratum is correct. Such probability  $p$  is further encoded as the mapping quality  $Q = -10 \log_{10}(1 - p)$ .

### Paired-end reads

A paired-end read whose mate is unmapped is treated as a single-end read, thus the above probabilities still hold. Given a read end, I denote by  $\mathbb{S}'_e$  the  $e$ -stratum of its mate, by  $Z'$  the canonical partition function of its mate as in equation 7.7, and by  $C'$  its correct mapping location. I partition each stratum  $\mathbb{S}_e$  in two subsets  $\mathbb{S}_e^P$  and  $\mathbb{S}_e^I$ , where  $\mathbb{S}_e^P$  contains all mapping locations in  $\mathbb{S}_e$  being properly paired to some mapping location in  $\mathbb{S}'$ , and  $\mathbb{S}_e^I = \mathbb{S}_e^P \setminus \mathbb{S}_e$ . I distinguish the cases of a read end having some properly paired mapping location or not.

If a read end has no properly paired mapping location,  $\mathbb{S}_e^I = \mathbb{S}_e$ . I estimate the probability of  $C$  being equal to any found mapping location  $(i, j, e)$  as

$$\Pr[C = (i, j, e) | \mathbb{S}_e^P = \emptyset] = \Pr[C = (i, j, e)] \cdot \Pr[C' \notin \mathbb{S}'_e] \quad (7.11)$$

and since

$$\Pr[C' \notin \mathbb{S}'_e] = 1 - \frac{1}{Z'} \cdot z_{e'} w(e'), \quad (7.12)$$

it follows that

$$\Pr[C = (i, j, e) | \mathbb{S}_e^P = \emptyset] = \frac{1}{Z} \cdot w(e) - \frac{1}{ZZ'} \cdot z_e w(e) w(e'). \quad (7.13)$$

Now I consider the case of a read end having some properly paired mapping location. In this case, Yara picks a random location from  $\mathbb{S}_e^P$ , under the assumption that any location

in  $\mathbb{S}_e^P$  is more likely to be correct than those in  $\mathbb{S}_e^I$ . Therefore, I assign different weights to  $\mathbb{S}_e^P$  and  $\mathbb{S}_e^I$ . I assume that

$$\Pr[C \in \mathbb{S}^P] = \Pr[C' \in \mathbb{S}_e^{P'}] = p' \quad (7.14)$$

and

$$\Pr[C \in \mathbb{S}^I] = 1 - p'. \quad (7.15)$$

Hence the probability of  $C$  being part of the properly paired locations in stratum  $\mathbb{S}_e$  is

$$\Pr[C \in \mathbb{S}_e^P] = \frac{p' z_e^P w(e)}{p' \cdot Z^P + (1 - p') \cdot Z^I} \quad (7.16)$$

and the probability that any of these locations is correct as

$$\Pr[C = (i, j, e) \in \mathbb{S}_e^P] = \frac{p' w(e)}{p' \cdot Z^P + (1 - p') \cdot Z^I}. \quad (7.17)$$

## 7.2 Evaluation

The evaluation consists of three experiments, all performed on the human reference genome (GRCh37). The first experiment assesses all-mapping sensitivity by applying the Rabema benchmark on real data. The second experiment evaluates best-mapping accuracy on simulated data. The last experiment assesses the performance of best-mappers within a variant calling pipeline applied to real data.

All experiments consider also read mapping throughputs and memory consumptions. Read mapping throughput is measured in *giga base pairs per hour* (Gbp/h). The Illumina HiSeq 2500 in a six days run produces<sup>1</sup> up to 800 Gbp as  $2 \times 100$  bp paired-end reads. Under this measure the maximum throughput of the Illumina HiSeq 2500 is 5.56 Gbp/h.

In addition to Yara, I consider the following state-of-the-art tools: Bowtie 2, BWA, GEM, RazerS 3, and Hobbes 2. The tools Bowtie 2 and BWA are suited for best-mapping, while RazerS 3 and Hobbes 2 for all-mapping. GEM, in addition to all-mapping, can be used for best-mapping only to some extent, as it does not compute mapping qualities.

### 7.2.1 Experimental setup

#### Read mappers parametrization

In appendix A.2, I give the exact parametrization of each read mapper considered in the evaluation. I parameterized Yara in three different ways, denoted as [s=0], [s=1], and [all]. Yara [s=0] thus reports only the co-optimal stratum as in equations 7.4, Yara [s=1] up to the first sub-optimal stratum, and Yara [all] all sub-optimal strata as in equation 7.2. Whenever possible, I configured the tools with the appropriate error rate (Yara, GEM, RazerS 3) or absolute number of errors (Hobbes 2). When processing paired-end reads, I provided the tools with appropriate insert size information.

<sup>1</sup> According to the specifications at [http://res.illumina.com/documents/products/datasheets/datasheet\\_hiseq2500.pdf](http://res.illumina.com/documents/products/datasheets/datasheet_hiseq2500.pdf) for high output run mode with dual flow cell.

## Infrastructure

All tools run on a desktop computer running Linux 3.10.11, equipped with one Intel® Core i7-4770K CPU @ 3.50 GHz, 32 GB RAM and a 2 TB HDD @ 7200 RPM. For maximum throughput, all tools run using eight threads. For accurate running time comparisons, I disabled Intel Turbo Boost; therefore, all measured running times might be slightly higher than the actual ones.

### 7.2.2 Rabema benchmark on real data

The first experiment uses the Rabema benchmark to evaluate all-mappers sensitivity on real data. The *Rabema benchmark* [Holtgrewe *et al.*, 2011] (v1.2) measures the sensitivity of read mappers in finding *relevant* mapping locations of genomic reads. I subdivide this experiment in two categories: *suboptimal* and *co-optimal*. In the category *suboptimal*, Rabema counts as relevant, for each read, all suboptimal mapping locations within a maximal edit distance error rate, i.e. all strata; in the category *all-best* it considers just *co-optimal* mapping locations, i.e. only the best stratum. Rabema computes the *sensitivity* of each tool as the fraction of relevant mapping locations found per read. For a thorough evaluation, Rabema classes mapping locations by their *error rate*, then computes sensitivity within each error rate class. The benchmark reports percentual scores normalized by the number of reads.

The data used in this experiment is a publicly released sequencing run (SRA/ENA id: ERR161544) performed at the Beijing Genome Institute; the genomic DNA used in this study came from an anonymous male Han Chinese individual who has no known genetic diseases. This dataset consists of  $2 \times 100$  bp whole genome sequencing reads produced by an Illumina HiSeq 2000 instrument. For practical reasons, in the category *co-optimal* I consider the first 10 M reads, while in the category *suboptimal* only the first 1 M reads.

An error rate of 5 %, is sufficient to map almost all reads in the dataset. Therefore, I built a Rabema gold standard by running RazerS 3 in full-sensitive mode within 5 % error rate. Subsequently, I provided the reads as unpaired to each tool, as the Rabema benchmark is not meaningful for paired-end reads.

#### Co-optimal mapping locations

Results are shown in table 7.1 (left panel). Yara [s=0] is the most sensitive tool in finding all co-optimal locations; it is full-sensitive up to 3 % error rate. GEM is not full-sensitive even though it claims to be so; it loses 5.6 % of normalized locations at 5 % error rate. Bowtie 2 and BWA are not designed for this task; indeed, they lose a significant number of co-optimal locations. In addition, Yara [s=0] has the highest throughput, being 1.7 times faster than GEM, 2.3 times faster than Bowtie 2 and 3.2 times faster than BWA.

#### Suboptimal mapping locations

Results are shown in table 7.1 (right panel). RazerS 3 is the only full-sensitive tool in finding all suboptimal locations. Hobbes 2 loses only a few points at 4-5 % error rate. Yara

**Table 7.1:** Rabema benchmark results on whole human genome 100 bp Illumina HiSeq 2000 reads (SRA/ENA id: ERR161544). The left panel shows the results of finding all co-optimal mapping locations of the first 10 M reads; the right panel shows the results of finding all suboptimal mapping locations of only the first 1 M reads. Big numbers show total Rabema scores, while small numbers show marginal scores for the mapping locations at  $\begin{pmatrix} 0 & 1 & 2 \\ 3 & 4 & 5 \end{pmatrix}$  % error rate.

	Co-opt. locations	Throughput	Memory		Subopt. locations	Throughput	Memory
	[%]	[Gbp/h]	[GB]		[%]	[Gbp/h]	[GB]
Yara [s=0]	100.0	13.59	4.88	Yara [all]	99.8	1.31	5.44
GEM	99.9	7.90	4.31	GEM	95.2	1.82	4.58
Bowtie 2	95.9	5.99	3.25	Hobbes 2	99.9	0.39	14.59
BWA	96.0	4.24	4.47	RazerS 3	100.0	0.06	22.80

[all] is full-sensitive up to 3 % error rate, nonetheless it loses only 2.8 % of normalized locations at 5 % error rate. However, looking at throughputs, Yara [all] is 3.6 times faster than Hobbes 2 and 21.9 times faster than RazerS 3. In addition, Yara [all] has a significantly lower memory footprint than its two competitors. GEM is again not full-sensitive; it loses more than 50 % of normalized locations at 5 % error rate. This could be due to a misconfiguration of the tool; nonetheless, I took the same parameterization used in the supplemental material of [Marco-Sola *et al.*, 2012].

### 7.2.3 Accuracy on simulated data

The second experiment evaluates the ability of best-mappers to find the *original* location of simulated reads and to estimate their mapping quality. Contrarily to the Rabema benchmark (section 7.2.2), this experiment relies on simulated data and considers only one *primary mapping location* per read. I consider only tools suited for this specific task, i.e. Bowtie 2 and BWA, in addition to Yara; GEM does not compute mapping qualities, while RazerS 3 and Hobbes 2 are two orders of magnitude slower than Yara, thus I discard them *a priori*.

For each tool, the accuracy benchmark counts each read as *correctly mapped* if its *primary* mapping location has been reported within 10 bp of the simulated location, or *incorrectly mapped* otherwise. Subsequently, the benchmark *stratifies* (i.e. sorts) primary mapping locations by mapping quality, s.t. mapping locations estimated to be correct by the mapper precede those estimated to be incorrect. Finally, the benchmark cumulates the counts of correctly and incorrectly mapped locations and plots them as *receiver operating characteristic* (ROC) curves.

The dataset consists of 10 M Illumina-like  $2 \times 100$  bp paired-end reads, simulated from the whole human reference genome using Mason [Holtgrewe, 2010]; the mean insert size is  $INS = 300$  and the standard deviation is  $ERR = 20$ . To assess the improvements due to the knowledge of the insert size, the experiment considers twice the same simulated reads, first as unpaired and then as paired-end including insert size information.



**Figure 7.1:** Accuracy on 100 bp Illumina-like single-end reads.

### Single-end reads

Accuracy results are shown in figure 7.1. The ROC curves show incorrect locations on a logarithmic scale. Yara [s=1] always dominates BWA, which in turn always dominates Bowtie 2. It is interesting to see that Yara [s=0] is not able to compute accurate mapping qualities without the knowledge of mapping locations from sub-optimal strata.

Performances are shown in table 7.2. Yara [s=0] is 3.6 times faster than BWA and 4.5 times faster than Bowtie 2. Bowtie 2 uses 1.5 GB less memory than Yara [s=0], probably thanks to a more compact FM-index implementation; however, memory footprint is not a limiting factor. Yara [s=1] is 3.0 times slower than Yara [s=0], but still faster than all other tools.

### Paired-end reads

Accuracy results are shown in figure 7.2. Again, the ROC curves show incorrect locations on a logarithmic scale. Yara [s=1] still dominates BWA on the most mappable reads; however, compared to single-end reads, BWA on paired-end reads is closer to Yara. Bowtie 2 shows the worst accuracy.

**Table 7.2:** Performance on 100 bp Illumina-like single-end reads.

	Yara [s=0]	Yara [s=1]	Bowtie 2	BWA
Throughput [Gbp/h]	26.28	8.88	5.84	7.37
Memory [GB]	4.88	4.93	3.25	4.48

**Figure 7.2:** Accuracy on  $2 \times 100$  bp Illumina-like paired-end reads.

Performances are shown in table 7.3. Bowtie 2 on paired-end reads has a higher throughput than on single-end reads. Conversely, BWA on paired-end reads is considerably slower than on single-end reads. Yara [s=0] is still 3.4 times faster than Bowtie 2 and 4.0 times faster than BWA.

#### 7.2.4 Variant calling on real data

The last experiment uses well-characterized datasets to estimate both true-positive and false-negative rates induced by various mappers within a best-mapping pipeline calling variants on real data. Each pipeline consists of one best-mapper, whose output is sorted by genomic position using *Samtools* [Li and Durbin, 2009], and subsequently given to the widely used *GATK Unified Genotyper* (UG) [DePristo *et al.*, 2011] to call both SNVs and INDELs. The evaluation consists of comparing the variants called by each pipeline to a set of high-confidence, single-nucleotide variant (SNV) and INDEL calls provided by the *Genome in a Bottle* (GIAB) consortium [Zook *et al.*, 2014].

The GIAB consortium provides a set of calls (NIST v2.18) for the pilot genome NA12878. In this experiment, each pipeline has to call variants from a whole exome sequencing (WES) run of the same NA12878 individual (SRA/ENA id: SRR1611178). This run, pro-

**Table 7.3:** Performance on  $2 \times 100$  bp Illumina-like paired-end reads.

	Yara [s=0]	Yara [s=1]	Bowtie 2	BWA
Throughput [Gbp/h]	23.04	8.37	6.78	5.70
Memory [GB]	5.00	5.11	3.29	4.66

**Figure 7.3:** SNVs calling accuracy on a  $150 \times$  coverage WES run consisting of  $2 \times 100$  bp Illumina HiSeq 2000 reads (SRA/ENA id: SRR1611178).



duced at the Icahn School of Medicine at Mount Sinai [Linderman *et al.*, 2014], consists of  $2 \times 100$  bp Illumina HiSeq 2000 reads and has a mean coverage of  $150 \times$ .

According to the GIAB guidelines, the evaluation first decomposes complex variants using the tool *vcfallelicprimitives* [Danecek *et al.*, 2011] and subsequently compares called variants to the set of GIAB trusted variants using the tool *USeq VCFComparator* [Nix *et al.*, 2008]. This last tool counts the number of correct and incorrect variants, stratifies them by variant quality, and reports their cumulated counts. The output of the evaluation is plotted as ROC curves.

### SNVs calling

Figure 7.3 shows the SNVs calling accuracy results. The ROC curve show incorrect SNVs on a logarithmic scale, to highlight high-quality calls. The plot shows that Yara induces a significantly higher rate of high-quality calls and up to 3 times less incorrect SNVs than BWA. There is no appreciable difference between Yara [s=0] and [s=1]. The poor performance of GEM could be partially due to the fact that it does not annotate mapping locations with qualities.

### INDELs calling

Figure 7.4 shows the INDELs calling accuracy results. Again, the ROC curves show incorrect INDELs on a logarithmic scale. Bowtie 2 dominates both BWA and Yara, possibly thanks to the fact that it computes local alignments rather than semi-global alignments. Yara [s=0] and [s=1] produce exactly the same set of calls: they induce a lower rate of correct INDEL calls and at the same time a higher rate of incorrect calls than Bowtie 2 and BWA. Surprisingly, GEM induces almost only incorrect INDEL calls.

**Figure 7.4:** INDELs calling accuracy on a  $150 \times$  coverage WES run consisting of  $2 \times 100$  bp Illumina HiSeq 2000 reads (SRA/ENA id: SRR1611178).



## Performance

Table 7.4 shows tools performances. Yara [s=0] is 2.0 times faster than GEM, 3.2 times faster than Bowtie 2, and 4.7 times faster than BWA. Compared to simulated WGS data (section 7.2.3), all tools exhibit a higher throughput on real WES data. This is due to the fact that the mappability of the whole human exome is higher than the average mappability of the whole human genome. This also explains why Yara [s=1] is now only 1.3 times slower than Yara [s=0].

## 7.3 Discussion

Stratified all-mapping is a well-defined and practical method for read mapping. The knowledge of mapping strata allows to accurately compute mapping qualities and pair read ends. To show the advantages of this method in real applications, I considered the genotyping of a human individual. Nonetheless, other HTS applications would equally benefit from this method.

**Table 7.4:** Throughput results on a  $150 \times$  coverage WES run consisting of  $2 \times 100$  bp Illumina HiSeq 2000 reads (SRA/ENA id: SRR1611178).

	Yara [s=0]	Yara [s=1]	Bowtie 2	BWA	GEM
Throughput [Gbp/h]	28.83	22.82	9.03	6.11	14.36
Memory [GB]	5.00	5.00	3.29	4.65	4.33

One such application is metagenomics classification: HTS reads coming from DNA samples recovered from e.g. soil or water are screened against a database consisting of bacterial reference genomes. The presence of several homologous genomes lowers the mappability of the database and requires the knowledge of multiple mapping locations per read. Actual read mappers are much slower or not sensitive enough for this task.

Another similar application is the resolution of alternative loci in the new human reference genome (GRCh38). This task could be simply performed by inferring the most likely alternative locus, e.g. from the density of mapped reads, and then re-assigning multi-reads to that locus. With de-facto standard mappers BWA and Bowtie 2, multi-reads would have to be re-mapped.

The weak spot of Yara is INDELs accuracy. Yara is unable to find medium to long INDELs (i.e. > 5 bp) because the reads spanning those INDELs remain unmapped. The edit distance does not support long INDELs over random errors, thus increasing the error rate only increases the rate of incorrect INDELs. However, in the paired-end protocol, the other end of such unmapped reads is usually mapped. As the insert size is known and relatively narrow, it is possible to spot these unmapped reads by searching them *online* within a neighborhood of their mapped end. This strategy is implemented by the tools BWA and Bowtie 2.

Performance-wise, Yara can still be improved. The verification time can be drastically reduced: David Weese and me vectorized the banded Myers' algorithm such that it verifies multiple patterns simultaneously; a prototypical implementation on a modern processor (256 bit AVX instructions) shows a speed-up of 4–6 times over the conventional algorithm. The filtration time can be almost halved by indexing both the forward and reverse complemented strands of the reference genome; this would be paid with a total memory footprint of about 8GB, which is nonetheless acceptable. Finally, the estimation of mapping qualities in Yara [s=0] can use matches from sub-optimal strata, which are found but not reported; this could render Yara [s=0] almost as good as Yara [s=1].

Yara is part of the SeqAn library [Döring *et al.*, 2008], it is distributed under the BSD license and can be downloaded from <http://www.seqan.de/projects/yara>. A manuscript for this work is in preparation and will be submitted to a peer-reviewed journal.



---

# A Read mappers parameterization

## A.1 Masai evaluation

In the following, I give the exact parameterization of each read mapper considered in the evaluation of section 6.2.

**Masai** Version 0.5 was used. In order to use Masai as an all-mapper, I passed the argument `-all`, otherwise the argument `-any-best` is used by default. I set the maximal edit distance using the parameter `-e`. I configured the seed length with the parameter `-seed-length`; on *E. coli*, *D. melanogaster* and *C. elegans* I chose a seed length of 16, while on *H. sapiens* I chose a seed length of 33. I selected the SAM output format with `-os` and enabled CIGAR output with `-oc`.

**Bowtie 2** Version 2.0.0-beta6 was used. I used the parameter `-end-to-end` to enforce semi-global read alignments. For the Rabema experiment I used the parameter `-k 100`.

**BWA** Version 0.6.1-r104 was used. For the Rabema experiment I passed the parameter `-N` to `aln` and `-n 100` to `samse`.

**Soap 2** Version 2.1 was used.

**RazerS3** Version 3.1 was used. I mapped with indels using the pigeonhole filter (default) and set the error rate through the parameter `-i`, e.g. `-i 95` to map within an error rate of 5 %. I selected the native or SAM output format with `-of 0` or `-of 4`.

**Hobbes** Version 1.3 was used. I built the index using the recommended  $q$ -gram length 11. Since I focus on edit distance, I used the 16 bit bit-vector version. I enabled indels with `-indels` and set maximal edit distance using the parameter `-v`. For resource measurement I used the output without CIGAR, for analyzing the results I enabled CIGAR output using `-cigar`.

**mrFAST** Version 2.1.0.6 was used. I set maximal edit distance using the parameter `-e`.

**SHRiMP 2** Version 2.2.2 was used.

## A.2 Yara evaluation

In the following, I give the exact parameterization of each read mapper considered in the evaluation of section 7.2. Below, MIN and MAX are placeholders for minimal and maximal insert size, while INS is the mean insert size and ERR its allowed deviation, i.e.  $INS = (MIN + MAX) / 2$ ,  $ERR = (MAX - MIN) / 2$ .

**Yara** Version 0.9.4 was used. By default, the tool runs as a best-mapper; I set the number of strata with the parameter `-s`. To perform all-mapping, I passed the argument `-all`. I set the error rate using the parameter `-e`. In paired-end mode, the parameters used were `-library-length INS -library-error ERR`. The number of threads was set with the parameter `-t`.

**GEM** Version 1.376 was used. I set the error rate using the parameters `-m` and `-e`, then I disabled adaptive mapping using the parameter `-quality-format ignore`. In best-mapping, to analyze only the best stratum, I passed the argument `-s 0`; in all-mapping, to analyze all strata, I passed `-d all -D all -s all -max-big-indel-length 0`. In single-end mode, I passed the parameter `-expect-single-end-reads`; in paired-end mode, I passed `-paired-end-alignment`, along with `-min-insert-size MIN -max-insert-size MAX`, and `-map-both-ends` to select the workflow mapping both reads independently. The number of threads was selected using the parameter `-t`.

**Bowtie 2** Version 2.2.1 was used. I used the parameter `-end-to-end` to enforce semi-global read alignments. In paired-end mode, I used the parameters `-minins MIN -maxins MAX`. The number of threads was selected using the parameter `-p`.

**BWA** Version 0.7.7-r441 was used. I used the parameter `-t` to select the number of threads in the `aln` step; the `sampe` and `samse` steps were performed using one thread since BWA does not offer any parallelization here.

**Hobbes 2** Version 2.0 was used. I built the index using the recommended  $q$ -gram length 11. I enabled edit distance with `-indels` and set the distance threshold using the parameter `-v`. In paired-end mode, I used the parameters `-pe -min MIN -max MAX`. Multi-threading was enabled using `-p`.

**RazerS 3** Version 3.3 was used. I set the error rate through the parameter `-i`, e.g. `-i 95` to map within an error rate of 5 %. I passed the option `-rr 100` to set the recognition rate to 100 % and `-m 10000000` to output all mapping locations per read. In paired-end mode, the parameters used were `-library-length INS -library-error ERR`. The number of threads was set with the `-tc` parameter.



## Personal information

**Name** Enrico Siragusa  
**Date of birth** 28 March 1985  
**Citizenship** Italian  
**Address** Meininger Str. 10  
10823 Berlin, Germany  
**Mobile phone** +49 (0)160 5168591  
**Email address** enrico.siragusa@fu-berlin.de



## Education

**2010–2015** **PhD in Bioinformatics**, *Free University of Berlin*. IMPRS for Computational Biology and Scientific Computing.  
**2008–2009** **MSc in Bioinformatics and Genome Research**, *Bielefeld University*. EU exchange program.  
**2007–2010** **MSc in Computer Science**, *UMLV Paris-Est*. Mention «Bien».  
**2006–2007** **BSc in Computer Science and Maths**, *UMLV Paris-Est*. Degree achieved via the international project «Vinci» sponsored by Italian and French Ministry of Education.  
**2003–2006** **BSc in Computer Science**, *Palermo University*. Grade «110/110 cum Laude».

## Professional experience

**2015– now** **Research associate**, *Free University of Berlin*. Berlin, Germany.  
**2014–2015** **Research associate**, *Berlin Institute of Health*. Berlin, Germany.  
**Apr–Jul 2013** **CUDA Software Engineer**, *NVIDIA Research*. Berlin, Germany.  
**2009–2010** **Software Developer**, *EffiCity*. Paris, France.  
**2007–2008** **Software Developer**, *eoRezo*. Paris, France.

## Teaching experience

- WS 2011** Data Structures and Algorithms for Bioinformatics — practical lectures and lab. *FU Berlin, 19700.*
- SS 2010** Project Management with SeqAn. *FU Berlin, 19588a.*

## Summer schools

- Aug 2013** Otto Warburg Summer School on NGS and its Impact on Genetics.
- Sep 2011** Otto Warburg Summer School on Evolutionary Genomics.

## Technical skills

- Languages** C++, OpenMP, Python, Perl, R, Matlab, Java, JavaScript, Bash,  $\text{\LaTeX}$ , PGF/TikZ.
- Libraries** STL, Boost, Thrust, NumPy, OpenCV, GTK+, Tk, Django, ExtJS, JQuery.
- Tools** Xcode & Instruments, Eclipse, CMake, GNU Toolchain, Valgrind, Flex & Bison, PerlXS, Subversion, Git, Mercurial.
- Methods** Generic Programming, Object Oriented Programming, Design Patterns.
- Databases** MySQL, PostgreSQL and PostGIS, MongoDB.

## Language skills

- Italian** Mother tongue
- English** Fluent
- French** Fluent
- Spanisch** Fluent
- German** Intermediate

## Transferable skills

### *Workshops provided by the Dahlem Research School*

Scientific writing in the life sciences. Teaching and learning at university level. Advanced scientific presentation. Confident speaking with impact. Speed reading.

## Presentations and talks

Massively parallel methods for NGS read mapping. *Otto Warburg Summer School 2013 on NGS and its Impact on Genetics*.

CUDA-accelerated SeqAn: effective sequence analysis on GPUs. In *GPU Acceleration of Bioinformatics Pipeline, NVIDIA Corporation, BoF Session. ISMB/ECCB 2013*.

Scalable string similarity search/join with approximate seeds and multiple backtracking. *EDBT/ICDT 2013 Workshop on Scalable String Similarity Search/Join*.

## Publications

Weese, D. and Siragusa, E. (submitted, 2015). *Algorithms for next-generation sequencing data*, chapter Full-Text Indexes for High-Throughput Sequencing. John Wiley & Sons

Smedley, D., Jacobson, J., Jäger, M., Köhler, S., Holtgrewe, M., Schubach, M., Siragusa, E., Zemojtel, T., Buske, O., Bone, W., Haendel, M., and Robinson, P. N. (accepted, 2015). Next-generation diagnostics and disease gene discovery with the exomiser. *Nat. Protoc*.

Wandelt, S., Deng, D., Gerdjikov, S., Mishra, S., Mitankin, P., Patil, M., Siragusa, E., Tiskin, A., Wang, W., Wang, J., *et al.* (2014). State-of-the-art in string similarity search and join. *ACM SIGMOD Record*, **43**(1), pages 64–76

Siragusa, E., Weese, D., and Reinert, K. (2013b). Scalable string similarity search/join with approximate seeds and multiple backtracking. In *Proceedings of the Joint EDBT/ICDT 2013 Workshops*, pages 370–374. ACM

Siragusa, E., Weese, D., and Reinert, K. (2013a). Fast and accurate read mapping with approximate seeds and multiple backtracking. *Nucleic Acids Res*.

Giancarlo, R., Siragusa, A., Siragusa, E., and Utro, F. (2007). A basic analysis toolkit for biological sequences. *Algorithms for Molecular Biology*, **2**(1), page 10



# C Declaration

I declare that this thesis is my own work and has not been submitted in any form for another degree or diploma at any university or other institute of tertiary education. Information derived from the published and unpublished work of others has been acknowledged in the text and a list of references is given.

---

Enrico Siragusa  
April 21, 2015



## BIBLIOGRAPHY

- Abouelhoda, M., Kurtz, S., and Ohlebusch, E. (2004). Replacing suffix trees with enhanced suffix arrays. *Journal of Discrete Algorithms*, **2**(1), pages 53–86.
- Ahmadi, A., Behm, A., Honnalli, N., Li, C., Weng, L., and Xie, X. (2012). Hobbes: optimized gram-based methods for efficient read alignment. *Nucleic Acids Res.*, **40**(6), page e41.
- Alkan, C., Kidd, J. M., Marques-Bonet, T., Aksay, G., Antonacci, F., Hormozdiari, F., Kitzman, J. O., Baker, C., Malig, M., Mutlu, O., Sahinalp, S. C., Gibbs, R. A., and Eichler, E. E. (2009). Personalized copy number and segmental duplication maps using next-generation sequencing. *Nat. Genet.*, **41**(10), pages 1061–1067.
- Apostolico, A. (1985). The myriad virtues of subword trees. In *Combinatorial algorithms on words*, pages 85–96. Springer.
- Arlazarov, V., Dinic, E., Kronrod, M., and Faradzev, I. (1970). On economical construction of the transitive closure of a directed graph. *Dokl. Akad. Nauk.*, **11**, page 194.
- Baeza-Yates, R. A. and Perleberg, C. H. (1992). Fast and practical approximate string matching. In *Combinatorial Pattern Matching*, pages 185–192. Springer.
- Bailey, J. A., Yavor, A. M., Massa, H. F., Trask, B. J., and Eichler, E. E. (2001). Segmental duplications: organization and impact within the current human genome project assembly. *Genome research*, **11**(6), pages 1005–1017.
- Bauer, M. J., Cox, A. J., and Rosone, G. (2013). Lightweight algorithms for constructing and inverting the bwt of string collections. *Theoretical Computer Science*, **483**, pages 134–148.
- Belazzougui, D., Cunial, F., Kärkkäinen, J., and Mäkinen, V. (2013). Versatile succinct representations of the bidirectional burrows-wheeler transform. In *Algorithms-ESA 2013*, pages 133–144. Springer.
- Bentley, D. R., Balasubramanian, S., Swerdlow, H. P., Smith, G. P., Milton, J., Brown, C. G., Hall, K. P., Evers, D. J., Barnes, C. L., Bignell, H. R., Boutell, J. M., Bryant, J., Carter, R. J., Keira Cheetham, R., Cox, A. J., Ellis, D. J., Flatbush, M. R., Gormley, N. A., Humphray, S. J., Irving, L. J., Karbelashvili, M. S., Kirk, S. M., Li, H., Liu, X., Maisinger, K. S., Murray, L. J., Obradovic, B., Ost, T., Parkinson, M. L., Pratt, M. R., Rasolonjatovo, I. M. J., Reed, M. T., Rigatti, R., Rodighiero, C., Ross, M. T., Sabot, A., Sankar, S. V., Scally, A., Schroth, G. P., Smith, M. E., Smith, V. P., Spiridou, A., Torrance, P. E., Tzonev, S. S., Vermaas, E. H., Walter,

- K., Wu, X., Zhang, L., Alam, M. D., Anastasi, C., Aniebo, I. C., Bailey, D. M. D., Bancarz, I. R., Banerjee, S., Barbour, S. G., Baybayan, P. A., Benoit, V. A., Benson, K. F., Bevis, C., Black, P. J., Boodhun, A., Brennan, J. S., Bridgham, J. A., Brown, R. C., Brown, A. A., Buermann, D. H., Bundu, A. A., Burrows, J. C., Carter, N. P., Castillo, N., Chiara E. Catenazzi, M., Chang, S., Neil Cooley, R., Crake, N. R., Dada, O. O., Diakoumakos, K. D., Dominguez-Fernandez, B., Earnshaw, D. J., Egbujor, U. C., Elmore, D. W., Etchin, S. S., Ewan, M. R., Fedurco, M., Fraser, L. J., Fuentes Fajardo, K. V., Scott Furey, W., George, D., Gietzen, K. J., Goddard, C. P., Golda, G. S., Granieri, P. A., Green, D. E., Gustafson, D. L., Hansen, N. F., Harnish, K., Haudenschild, C. D., Heyer, N. I., Hims, M. M., Ho, J. T., Horgan, A. M., Hoschler, K., Hurwitz, S., Ivanov, D. V., Johnson, M. Q., James, T., Huw Jones, T. A., Kang, G.-D., Kerelska, T. H., Kersey, A. D., Khrebtukova, I., Kindwall, A. P., Kingsbury, Z., Kokko-Gonzales, P. I., Kumar, A., Laurent, M. A., Lawley, C. T., Lee, S. E., Lee, X., Liao, A. K., Loch, J. A., Lok, M., Luo, S., Mammen, R. M., Martin, J. W., McCauley, P. G., McNitt, P., Mehta, P., Moon, K. W., Mullens, J. W., Newington, T., Ning, Z., Ling Ng, B., Novo, S. M., O'Neill, M. J., Osborne, M. A., Osnowski, A., Ostadan, O., Paraschos, L. L., Pickering, L., Pike, A. C., Pike, A. C., Chris Pinkard, D., Pliskin, D. P., Podhasky, J., Quijano, V. J., Raczy, C., Rae, V. H., Rawlings, S. R., Chiva Rodriguez, A., Roe, P. M., Rogers, J., Rogert Bacigalupo, M. C., Romanov, N., Romieu, A., Roth, R. K., Rourke, N. J., Ruediger, S. T., Rusman, E., Sanches-Kuiper, R. M., Schenker, M. R., Seoane, J. M., Shaw, R. J., Shiver, M. K., Short, S. W., Sizto, N. L., Sluis, J. P., Smith, M. A., Ernest Sohna, J., Spence, E. J., Stevens, K., Sutton, N., Szajkowski, L., Tregidgo, C. L., Turcatti, G., Vandevondele, S., Verhovsky, Y., Virk, S. M., Wakelin, S., Walcott, G. C., Wang, J., Worsley, G. J., Yan, J., Yau, L., Zuerlein, M., Rogers, J., Mullikin, J. C., Hurler, M. E., McCooke, N. J., West, J. S., Oaks, F. L., Lundberg, P. L., Klennerman, D., Durbin, R., and Smith, A. J. (2008). Accurate whole human genome sequencing using reversible terminator chemistry. *Nature*, **456**(7218), pages 53–59.
- Brown, D. G. (2008). A survey of seeding for sequence alignment. *Bioinformatics algorithms: techniques and applications*, pages 126–152.
- Burkhardt, S. and Kärkkäinen, J. (2001). Better filtering with gapped q-grams. In *Proc. of the 12th Annual Symposium on Combinatorial Pattern Matching, CPM '01*, pages 73–85. Springer.
- Burkhardt, S., Crauser, A., Ferragina, P., Lenhof, H.-P., Rivals, E., and Vingron, M. (1999). q-gram based database searching using a suffix array (QUASAR). In *Proc. of the 3rd Annual International Conference on Research in Computational Molecular Biology, RECOMB '99*, pages 77–83. ACM Press.
- Burrows, M. and Wheeler, D. J. (1994). A block-sorting lossless data compression algorithm. Technical Report 124, Digital SRC Research Report.
- Califano, A. and Rigoutsos, I. (1993). Flash: A fast look-up algorithm for string homology. In *Computer Vision and Pattern Recognition, 1993. Proceedings CVPR'93., 1993 IEEE Computer Society Conference on*, pages 353–359. IEEE.
- Consortium, I. H. G. S. (2001). Initial sequencing and analysis of the human genome. *Nature*, **409**(6822), pages 860–921.



- Crochemore, M., Grossi, R., Kärkkäinen, J., and Landau, G. M. (2013). A constant-space comparison-based algorithm for computing the burrows–wheeler transform. In *Combinatorial Pattern Matching*, pages 74–82. Springer.
- Danecek, P., Auton, A., Abecasis, G., Albers, C. A., Banks, E., DePristo, M. A., Handsaker, R. E., Lunter, G., Marth, G. T., Sherry, S. T., *et al.* (2011). The variant call format and vcftools. *Bioinformatics*, **27**(15), pages 2156–2158.
- David, M., Dzamba, M., Lister, D., Ilie, L., and Brudno, M. (2011). SHRiMP2: sensitive yet practical short read mapping. *Bioinformatics*, **27**(7), pages 1011–1012.
- Dehal, P. and Boore, J. L. (2005). Two rounds of whole genome duplication in the ancestral vertebrate. *PLoS biology*, **3**(10), page e314.
- Dementiev, R., Kärkkäinen, J., Mehnert, J., and Sanders, P. (2008). Better external memory suffix array construction. *J. Exp. Algorithmics*, **12**, pages 3.4:1–3.4:24.
- DePristo, M. A., Banks, E., Poplin, R., Garimella, K. V., Maguire, J. R., Hartl, C., Philippakis, A. A., del Angel, G., Rivas, M. A., Hanna, M., *et al.* (2011). A framework for variation discovery and genotyping using next-generation dna sequencing data. *Nature genetics*, **43**(5), pages 491–498.
- Derrien, T., Estellé, J., Marco Sola, S., Knowles, D. G., Raineri, E., Guigó, R., and Ribeca, P. (2012). Fast computation and applications of genome mappability. *PLoS ONE*, **7**(1), page e30377.
- Döring, A., Weese, D., Rausch, T., and Reinert, K. (2008). SeqAn an efficient, generic C++ library for sequence analysis. *BMC Bioinformatics*, **9**, page 11.
- Elias, P. (1975). Universal codeword sets and representations of the integers. *Information Theory, IEEE Transactions on*, **21**(2), pages 194–203.
- Ewing, B. and Green, P. (1998). Base-calling of automated sequencer traces using phred. ii. error probabilities. *Genome research*, **8**(3), pages 186–194.
- Ewing, B., Hillier, L., Wendl, M. C., and Green, P. (1998). Base-calling of automated sequencer traces using phred. i. accuracy assessment. *Genome research*, **8**(3), pages 175–185.
- Faro, S. and Lecroq, T. (2013). The exact online string matching problem: a review of the most recent results. *ACM Computing Surveys (CSUR)*, **45**(2), page 13.
- Farrar, M. (2007). Striped smith–waterman speeds database searches six times over other simd implementations. *Bioinformatics*, **23**(2), pages 156–161.
- Ferragina, P. and Manzini, G. (2000). Opportunistic data structures with applications. In *Foundations of Computer Science, 2000. Proceedings. 41st Annual Symposium on*, pages 390–398. IEEE.

- Ferragina, P. and Manzini, G. (2001). An experimental study of an opportunistic index. In *SODA*, pages 269–278.
- Fonseca, N. A., Rung, J., Brazma, A., and Marioni, J. C. (2012). Tools for mapping high-throughput sequencing data. *Bioinformatics*, **28**(24), pages 3169–3177.
- Galil, Z. and Giancarlo, R. (1988). Data structures and algorithms for approximate string matching. *Journal of Complexity*, **4**(1), pages 33–72.
- Gallant, J., Maier, D., and Astorer, J. (1980). On finding minimal length superstrings. *Journal of Computer and System Sciences*, **20**(1), pages 50–58.
- Giancarlo, R., Siragusa, A., Siragusa, E., and Utro, F. (2007). A basic analysis toolkit for biological sequences. *Algorithms for Molecular Biology*, **2**(1), page 10.
- Giegerich, R., Kurtz, S., and Stoye, J. (1999). Efficient implementation of lazy suffix trees. In *Algorithm Engineering*, pages 30–42. Springer.
- Giegerich, R., Kurtz, S., and Stoye, J. (2003). Efficient implementation of lazy suffix trees. *Softw., Pract. Exper.*, pages 1035–1049.
- Grossi, R., Gupta, A., and Vitter, J. S. (2003). High-order entropy-compressed text indexes. In *Proc. of the 14th annual ACM-SIAM symposium on Discrete algorithms, SODA '03*, pages 841–850, Philadelphia, PA, USA. Society for Industrial and Applied Mathematics.
- Gusfield, D. (1997). *Algorithms on strings, trees, and sequences: Computer science and computational biology*. Cambridge University Press, New York, NY, USA.
- Hach, F., Hormozdiari, F., Alkan, C., Hormozdiari, F., Birol, I., Eichler, E. E., and Sahinalp, S. C. (2010). mrsFAST: a cache-oblivious algorithm for short-read mapping. *Nat. Methods*, **7**(8), pages 576–577.
- Hamming, R. W. (1950). Error detecting and error correcting codes. *Syst. Tech. J.*, **29**, pages 147–160.
- Holtgrewe, M. (2010). Mason – a read simulator for second generation sequencing data. Technical Report TR-B-10-06, Institut für Mathematik und Informatik, Freie Universität Berlin.
- Holtgrewe, M., Emde, A.-K., Weese, D., and Reinert, K. (2011). A novel and well-defined benchmarking method for second generation read mapping. *BMC Bioinformatics*, **12**, page 210.
- Intel (2011). *Intel® 64 and IA-32 Architectures Optimization Reference Manual*. Intel Corporation.
- Jacobson, G. (1989). Space-efficient static trees and graphs. In *Foundations of Computer Science, 1989, 30th Annual Symposium on*, pages 549–554. IEEE.

- Jokinen, P. and Ukkonen, E. (1991). Two algorithms for approximate string matching in static texts. In *Mathematical Foundations of Computer Science 1991*, pages 240–248. Springer.
- Kärkkäinen, J. and Sanders, P. (2003). Simple linear work suffix array construction. *ICALP*, pages 943–955.
- Karp, R. M., Luby, M., and Madras, N. (1989). Monte-carlo approximation algorithms for enumeration problems. *Journal of algorithms*, **10**(3), pages 429–448.
- Kasai, T., Lee, G., Arimura, H., Arikawa, S., and Park, K. (2001). Linear-time longest-common-prefix computation in suffix arrays and its applications. In *CPM*, pages 181–192.
- Kehr, B., Weese, D., and Reinert, K. (2011). Stellar: fast and exact local alignments. *BMC Bioinf.*, **12**(Suppl 9), page S15.
- Kim, J., Li, C., and Xie, X. (2014). Improving read mapping using additional prefix grams. *BMC bioinformatics*, **15**(1), page 42.
- Knuth, D. (1973). *The Art of Computer Programming. Volume 3, Addison-Wesley*.
- Kucherov, G., Noé, L., and Roytberg, M. (2005). Multiseed lossless filtration. *IEEE/ACM Transactions on Computational Biology and Bioinformatics (TCBB)*, **2**(1), pages 51–61.
- Kurtz, S. (1999). Reducing the space requirement of suffix trees. *Software-Practice and Experience*, **29**(13), pages 1149–71.
- Langmead, B. and Salzberg, S. L. (2012). Fast gapped-read alignment with Bowtie 2. *Nat. Methods*, **9**(4), pages 357–359.
- Langmead, B., Trapnell, C., Pop, M., and Salzberg, S. L. (2009). Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol.*, **10**(3), page R25.
- Lee, H. and Schatz, M. C. (2012). Genomic dark matter: the reliability of short read mapping illustrated by the genome mappability score. *Bioinformatics*, **28**(16), pages 2097–2105.
- Levenshtein, V. I. (1966). Binary codes capable of correcting deletions, insertions and reversals. *Soviet Physics – Doklady*, **10**, pages 707–710.
- Li, H. (2012). Exploring single-sample snp and indel calling with whole-genome de novo assembly. *Bioinformatics*, **28**(14), pages 1838–1844.
- Li, H. and Durbin, R. (2009). Fast and accurate short read alignment with burrows-wheeler transform. *Bioinformatics*, **25**(14), pages 1754–1760.

- Li, H. and Homer, N. (2010). A survey of sequence alignment algorithms for next-generation sequencing. *Brief Bioinform.*, **11**(5), pages 473–483.
- Li, H., Ruan, J., and Durbin, R. (2008). Mapping short dna sequencing reads and calling variants using mapping quality scores. *Genome research*, **18**(11), pages 1851–1858.
- Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., Durbin, R., and 1000 Genome Project Data Processing Subgroup (2009a). The sequence alignment/map format and SAMtools. *Bioinformatics*, **25**(16), pages 2078–2079.
- Li, R., Yu, C., Li, Y., Lam, T.-W., Yiu, S.-M., Kristiansen, K., and Wang, J. (2009b). SOAP2: an improved ultrafast tool for short read alignment. *Bioinformatics*, **25**(15), pages 1966–1967.
- Linderman, M. D., Brandt, T., Edelmann, L., Jabado, O., Kasai, Y., Kornreich, R., Mahajan, M., Shah, H., Kasarskis, A., and Schadt, E. E. (2014). Analytical validation of whole exome and whole genome sequencing for clinical applications. *BMC medical genomics*, **7**(1), page 20.
- Maier, D. and Storer, J. A. (1977). A note on the complexity of the superstring problem. *Computer Science Laboratory, Report*, (233).
- Mäkinen, V., Välimäki, N., Laaksonen, A., and Katainen, R. (2010). Unified view of backward backtracking in short read mapping. In T. Elomaa, H. Mannila, and P. Orponen, editors, *Algorithms and Applications*, volume 6060 of *Lecture Notes in Computer Science*, pages 182–195. Springer Berlin Heidelberg.
- Manber, U. and Myers, G. (1990). Suffix arrays: a new method for on-line string searches. In *SODA*, pages 319–327.
- Mantaci, S., Restivo, A., Rosone, G., and Sciortino, M. (2007). An extension of the burrows–wheeler transform. *Theoretical Computer Science*, **387**(3), pages 298–312.
- Marco-Sola, S., Sammeth, M., Guigó, R., and Ribeca, P. (2012). The gem mapper: fast, accurate and versatile alignment by filtration. *Nature methods*, **9**(12), pages 1185–1188.
- Meyne, J., Baker, R. J., Hobart, H. H., Hsu, T., Ryder, O. A., Ward, O. G., Wiley, J. E., Wurster-Hill, D. H., Yates, T. L., and Moyzis, R. K. (1990). Distribution of non-telomeric sites of the (ttaggg) *n* telomeric sequence in vertebrate chromosomes. *Chromosoma*, **99**(1), pages 3–10.
- Morrison, D. R. (1968). Patricia – practical algorithm to retrieve information coded in alphanumeric. *J. ACM*, **15**(4), pages 514–534.
- Mortazavi, A., Williams, B., McCue, K., Schaeffer, L., and Wold, B. (2008). Mapping and quantifying mammalian transcriptomes by RNA-seq. *Nat. Methods*, **5**(7), pages 621–628.

- Myers, E. W. (1994). A sublinear algorithm for approximate keyword searching. *Algorithmica*, **12**(4-5), pages 345–374.
- Myers, E. W. (2005). The fragment assembly string graph. *Bioinformatics*, **21**(suppl 2), pages ii79–ii85.
- Myers, G. (1999). A fast bit-vector algorithm for approximate string matching based on dynamic programming. *J. ACM*, **46**(3), pages 395–415.
- Navarro, G. (2001). A guided tour to approximate string matching. *ACM Comput. Surv.*, **33**(1), pages 31–88.
- Navarro, G. and Baeza-Yates, R. A. (2000). A hybrid indexing method for approximate string matching. *Journal of Discrete Algorithms*, **1**(1), pages 205–239.
- Navarro, G. and Mäkinen, V. (2007). Compressed full-text indexes. *ACM Computing Surveys (CSUR)*, **39**(1), page 2.
- Navarro, G., Baeza-Yates, R. A., Sutinen, E., and Tarhio, J. (2001). Indexing methods for approximate string matching. *IEEE Data Eng. Bull.*, **24**(4), pages 19–27.
- Needleman, S. B. and Wunsch, C. D. (1970). A general method applicable to the search for similarities in the amino acid sequence of two proteins. *J. Mol. Biol.*, **48**, pages 443–453.
- Nemhauser, G. L., Wolsey, L. A., and Fisher, M. L. (1978). An analysis of approximations for maximizing submodular set functions—i. *Mathematical Programming*, **14**(1), pages 265–294.
- Nicolas, F. and Rivals, E. (2005). Hardness of optimal spaced seed design. In *Combinatorial Pattern Matching*, pages 144–155. Springer.
- Nix, D. A., Courdy, S. J., and Boucher, K. M. (2008). Empirical methods for controlling false positives and estimating confidence in chip-seq peaks. *BMC bioinformatics*, **9**(1), page 523.
- Pevzner, P. A., Tang, H., and Waterman, M. S. (2001). An eulerian path approach to dna fragment assembly. *Proceedings of the National Academy of Sciences*, **98**(17), pages 9748–9753.
- Rasmussen, K. R., Stoye, J., and Myers, E. W. (2006). Efficient q-gram filters for finding all  $\epsilon$ -matches over a given length. *J. Comput. Biol.*, **13**(2), pages 296–308.
- Rusk, N. (2010). Torrents of sequence. *Nature Methods*, **8**(1), pages 44–44.
- Samonte, R. V. and Eichler, E. E. (2002). Segmental duplications and the evolution of the primate genome. *Nature Reviews Genetics*, **3**(1), pages 65–72.
- Sanger, F., Nicklen, S., and Coulson, A. R. (1977). DNA sequencing with chain-terminating inhibitors. *PNAS*, **74**(12), pages 5463–5467.

- Schürmann, K.-B. and Stoye, J. (2007). An incomplex algorithm for fast suffix array construction. *Software: Practice and Experience*, **37**(3), pages 309–329.
- Simola, D. F. and Kim, J. (2011). Sniper: improved snp discovery by multiply mapping deep sequenced reads. *Genome Biol*, **12**(6), page R55.
- Singer, J. (2013). *A Wavelet Tree Based FM-Index for Biological Sequences in SeqAn*. Master's thesis, Freie Universität Berlin.
- Siragusa, E., Weese, D., and Reinert, K. (2013a). Fast and accurate read mapping with approximate seeds and multiple backtracking. *Nucleic Acids Res*.
- Siragusa, E., Weese, D., and Reinert, K. (2013b). Scalable string similarity search/join with approximate seeds and multiple backtracking. In *Proceedings of the Joint EDBT/ICDT 2013 Workshops*, pages 370–374. ACM.
- Smedley, D., Jacobson, J., Jäger, M., Köhler, S., Holtgrewe, M., Schubach, M., Siragusa, E., Zemojtel, T., Buske, O., Bone, W., Haendel, M., and Robinson, P. N. (accepted, 2015). Next-generation diagnostics and disease gene discovery with the exomiser. *Nat. Protoc*.
- Smit, A. F. (1996). The origin of interspersed repeats in the human genome. *Current opinion in genetics & development*, **6**(6), pages 743–748.
- Smith, T. F. and Waterman, M. S. (1981). Identification of Common Molecular Subsequences. *J. Mol. Biol.*, **147**, pages 195–197.
- Treangen, T. J. and Salzberg, S. L. (2011). Repetitive dna and next-generation sequencing: computational challenges and solutions. *Nature Reviews Genetics*, **13**(1), pages 36–46.
- Turner, J. S. (1989). Approximation algorithms for the shortest common superstring problem. *Information and computation*, **83**(1), pages 1–20.
- Ukkonen, E. (1993). Approximate string-matching over suffix trees. In *CPM*, pages 228–242.
- Vazirani, V. V. (2001). *Approximation algorithms*. springer.
- Venter, J. C., Adams, M. D., Myers, E. W., Li, P. W., Mural, R. J., Sutton, G. G., Smith, H. O., Yandell, M., Evans, C. A., Holt, R. A., Gocayne, J. D., Amanatides, P., Ballew, R. M., Huse, D. H., Wortman, J. R., Zhang, Q., Kodira, C. D., Zheng, X. H., Chen, L., Skupski, M., Subramanian, G., Thomas, P. D., Zhang, J., Gabor Miklos, G. L., Nelson, C., Broder, S., Clark, A. G., Nadeau, J., McKusick, V. A., Zinder, N., Levine, A. J., Roberts, R. J., Simon, M., Slayman, C., Hunkapiller, M., Bolanos, R., Delcher, A., Dew, I., Fasulo, D., Flanigan, M., Florea, L., Halpern, A., Hannenhalli, S., Kravitz, S., Levy, S., Mobarry, C., Reinert, K., Remington, K., Abu-Threideh, J., Beasley, E., Biddick, K., Bonazzi, V., Brandon, R., Cargill, M., Chandramouliswaran, I., Charlab, R., Chaturvedi, K., Deng, Z., Di Francesco, V., Dunn, P., Eilbeck, K., Evangelista, C., Gabrielian, A. E., Gan, W., Ge, W., Gong, F., Gu,

- Z., Guan, P., Heiman, T. J., Higgins, M. E., Ji, R. R., Ke, Z., Ketchum, K. A., Lai, Z., Lei, Y., Li, Z., Li, J., Liang, Y., Lin, X., Lu, F., Merkulov, G. V., Milshina, N., Moore, H. M., Naik, A. K., Narayan, V. A., Neelam, B., Nusskern, D., Rusch, D. B., Salzberg, S., Shao, W., Shue, B., Sun, J., Wang, Z., Wang, A., Wang, X., Wang, J., Wei, M., Wides, R., Xiao, C., Yan, C., Yao, A., Ye, J., Zhan, M., Zhang, W., Zhang, H., Zhao, Q., Zheng, L., Zhong, F., Zhong, W., Zhu, S., Zhao, S., Gilbert, D., Baumhueter, S., Spier, G., Carter, C., Cravchik, A., Woodage, T., Ali, F., An, H., Awe, A., Baldwin, D., Baden, H., Barnstead, M., Barrow, I., Beeson, K., Busam, D., Carver, A., Center, A., Cheng, M. L., Curry, L., Danaher, S., Davenport, L., Desilets, R., Dietz, S., Dodson, K., Doup, L., Ferriera, S., Garg, N., Gluecksmann, A., Hart, B., Haynes, J., Haynes, C., Heiner, C., Hladun, S., Hostin, D., Houck, J., Howland, T., Ibegwam, C., Johnson, J., Kalush, F., Kline, L., Koduru, S., Love, A., Mann, F., May, D., McCawley, S., McIntosh, T., McMullen, I., Moy, M., Moy, L., Murphy, B., Nelson, K., Pfannkoch, C., Pratts, E., Puri, V., Qureshi, H., Reardon, M., Rodriguez, R., Rogers, Y. H., Rombiad, D., Ruhfel, B., Scott, R., Sitter, C., Smallwood, M., Stewart, E., Strong, R., Suh, E., Thomas, R., Tint, N. N., Tse, S., Vech, C., Wang, G., Wetter, J., Williams, S., Williams, M., Windsor, S., Winn-Deen, E., Wolfe, K., Zaveri, J., Zaveri, K., Abril, J. F., Guigó, R., Campbell, M. J., Sjolander, K. V., Karlak, B., Kejariwal, A., Mi, H., Lazareva, B., Hatton, T., Narechania, A., Diemer, K., Muruganujan, A., Guo, N., Sato, S., Bafna, V., Istrail, S., Lippert, R., Schwartz, R., Walenz, B., Yooseph, S., Allen, D., Basu, A., Baxendale, J., Blick, L., Caminha, M., Carnes-Stine, J., Caulk, P., Chiang, Y. H., Coyne, M., Dahlke, C., Mays, A., Dombroski, M., Donnelly, M., Ely, D., Esparham, S., Fosler, C., Gire, H., Glanowski, S., Glasser, K., Glodek, A., Gorokhov, M., Graham, K., Gropman, B., Harris, M., Heil, J., Henderson, S., Hoover, J., Jennings, D., Jordan, C., Jordan, J., Kasha, J., Kagan, L., Kraft, C., Levitsky, A., Lewis, M., Liu, X., Lopez, J., Ma, D., Majoros, W., McDaniel, J., Murphy, S., Newman, M., Nguyen, T., Nguyen, N., Nodell, M., Pan, S., Peck, J., Peterson, M., Rowe, W., Sanders, R., Scott, J., Simpson, M., Smith, T., Sprague, A., Stockwell, T., Turner, R., Venter, E., Wang, M., Wen, M., Wu, D., Wu, M., Xia, A., Zandieh, A., and Zhu, X. (2001). The sequence of the human genome. *Science*, **291**, pages 1304–1351.
- Wandelt, S., Deng, D., Gerdjikov, S., Mishra, S., Mitankin, P., Patil, M., Siragusa, E., Tiskin, A., Wang, W., Wang, J., *et al.* (2014). State-of-the-art in string similarity search and join. *ACM SIGMOD Record*, **43**(1), pages 64–76.
- Wang, Z., Weber, J. L., Zhong, G., and Tanksley, S. (1994). Survey of plant short tandem dna repeats. *Theoretical and applied genetics*, **88**(1), pages 1–6.
- Weese, D. (2013). *Indices and Applications in High-Throughput Sequencing*. Ph.D. thesis, Freie Universität Berlin.
- Weese, D. and Siragusa, E. (submitted, 2015). *Algorithms for next-generation sequencing data*, chapter Full-Text Indexes for High-Throughput Sequencing. John Wiley & Sons.
- Weese, D., Emde, A.-K., Rausch, T., Döring, A., and Reinert, K. (2009). RazerS—fast read mapping with sensitivity control. *Genome Res.*, **19**(9), pages 1646–1654.
- Weese, D., Holtgrewe, M., and Reinert, K. (2012). RazerS 3: Faster, fully sensitive read mapping. *Bioinformatics*. 10.1093/bioinformatics/bts505.

- Weese, D., Schulz, M. H., Holtgrewe, M., and Richard, H. (2013). Fiona: a versatile and automatic strategy for read error correction. *to appear*.
- Weiner, P. (1973). Linear pattern matching algorithms. In *SWAT (FOCS)*, pages 1–11. IEEE.
- Wilkes, M. V. (1995). The memory wall and the cmos end-point. *ACM SIGARCH Computer Architecture News*, **23**(4), pages 4–6.
- Wolfe, K. H. and Shields, D. C. (1997). Molecular evidence for an ancient duplication of the entire yeast genome. *Nature*, **387**(6634), pages 708–712.
- Wooster, R., Cleton-Jansen, A.-M., Collins, N., Mangion, J., Cornelis, R., Cooper, C., Gusteron, B., Ponder, B., Von Deimling, A., Wiestler, O., *et al.* (1994). Instability of short tandem repeats (microsatellites) in human cancers. *Nature genetics*, **6**(2), pages 152–156.
- Zook, J. M., Chapman, B., Wang, J., Mittelman, D., Hofmann, O., Hide, W., and Salit, M. (2014). Integrating human sequence data sets provides a resource of benchmark snp and indel genotype calls. *Nature biotechnology*.



## LIST OF FIGURES

2.1	Edit transcript and alignment . . . . .	9
2.2	$k$ -Differences occurrence . . . . .	12
2.3	Suffix trie and suffix tree . . . . .	14
2.4	Generalized suffix trie . . . . .	15
3.1	Suffix array and generalized suffix array . . . . .	20
3.2	$q$ -Gram index . . . . .	24
3.3	Burrows-Wheeler transform . . . . .	27
3.4	Permutations $LF$ and $\Psi$ . . . . .	28
3.5	Inversion of the Burrows-Wheeler transform . . . . .	29
3.6	Binary rank dictionaries . . . . .	31
3.7	DNA rank dictionary . . . . .	32
3.8	Wavelet tree . . . . .	33
3.9	Depth-first traversal runtime . . . . .	37
3.10	Exact string matching runtime . . . . .	38
3.11	$k$ -Mismatches runtime . . . . .	39
3.12	Multiple exact string matching runtime . . . . .	42
3.13	Multiple $k$ -mismatches runtime . . . . .	42
3.14	Multiple $k$ -mismatches speed-up on SA . . . . .	43
3.15	Multiple $k$ -mismatches speed-up on FM-index . . . . .	43
4.1	Filtration with exact seeds . . . . .	46
4.2	Filtration with approximate seeds . . . . .	48
4.3	Filtration with contiguous $q$ -grams . . . . .	50
4.4	Parallelogram buckets . . . . .	51
4.5	Filtration with gapped $q$ -grams . . . . .	52
4.6	Filtration with multiple gapped $q$ -grams . . . . .	60
4.7	Filters runtime on $k$ -mismatches . . . . .	61
4.8	Filters runtime on $k$ -differences . . . . .	61
4.9	Filtration versus verification time on $k$ -mismatches . . . . .	63
4.10	Filters specificity on $k$ -mismatches . . . . .	63
7.1	Yara accuracy on Illumina-like single-end reads . . . . .	101
7.2	Yara accuracy on Illumina-like paired-end reads . . . . .	102
7.3	Yara SNVs calling accuracy . . . . .	103
7.4	Yara INDELs calling accuracy . . . . .	104



## LIST OF TABLES

2.1	Classification of text locations by filtering methods . . . . .	16
3.1	Index construction times and memory footprints . . . . .	36
4.1	Best $Q$ -grams found by Burkhardt and Kärkkäinen . . . . .	56
4.2	Optimal threshold computation results . . . . .	56
4.3	Measurement of filtering methods efficiency . . . . .	62
5.1	Mappability of model genomes . . . . .	74
5.2	Human genome mappability score . . . . .	76
5.3	Overview of popular read mappers . . . . .	80
6.1	Masai results in the Rabema benchmark . . . . .	87
6.2	Masai variant detection results . . . . .	88
6.3	Masai performance on real data . . . . .	90
6.4	Masai performance with different indices . . . . .	91
6.5	Masai filtration efficiency results . . . . .	91
7.1	Yara results in the Rabema benchmark . . . . .	100
7.2	Yara performance on Illumina-like single-end reads . . . . .	101
7.3	Yara performance on Illumina-like paired-end reads . . . . .	102
7.4	Yara throughput on whole exome sequencing data – Illumina HiSeq 2000 . . . . .	104

