# Dissertation Title

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 201X

# Abstract

Bla bla bla.

# CONTENTS

# Part I

# APPROXIMATE STRING MATCHING

# Background

## 1.1 Introduction

### 1.1.1 Problem definition

### 1.1.2 Bioinformatics applications

## 1.2 Overview of existing methods

### 1.2.1 Online methods

### 1.2.2 Indexed methods

### 1.2.3 Filtering methods

## 1.3 Related problems

### 1.3.1 Local similarity search

### 1.3.2 Dictionary search

### 1.3.3 Overlaps computation

# Online Methods

# Indexed Methods

## 3.1 Classic Full-Text Indices

### 3.1.1 Trie

### 3.1.2 Suffix trie and suffix tree

### 3.1.3 Suffix array

### 3.1.4 $q$-Gram index

## 3.2 Compressed Full-Text Indices

### 3.2.1 Burrows-Wheeler transform

### 3.2.2 FM-index

### 3.2.3 Rank dictionaries

## 3.3 Backtracking

### 3.3.1 Pruning methods

### 3.3.2 Multiple backtracking

# Filtering Methods

## 4.1 $q$-Gram filters

### 4.1.1 Exact seeds

### 4.1.2 Gapped seeds

## 4.2 Factor filters

### 4.2.1 Exact seeds

### 4.2.2 Approximate seeds

## 4.3 Suffix filters

**Part II**

**APPLICATIONS**

# Read Mapping

## 5.1 Related work

### 5.1.1 Best mappers

### 5.1.2 All mappers

## 5.2 The Masai mapper

### 5.2.1 Single-end mapping

### 5.2.2 Paired-end mapping

## 5.3 The Masai 2 mapper

### 5.3.1 Single-end mapping

### 5.3.2 Paired-end mapping

### 5.3.3 Parallelization

### 5.3.4 Hardware acceleration

## 5.4 Experimental results

### 5.4.1 Comparison of filtration strategies

### 5.4.2 Rabema benchmark results

### 5.4.3 Variant detection results

### 5.4.4 Runtime results

## 5.5 Discussion

# String Similarity Search / Join

APPENDIX

# A

# 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 Information derived from the published and unpublished work of others has been acknowledged in the text and a list of references is given.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
October 6, 2013

# LIST OF FIGURES

# LIST OF TABLES

# LIST OF ALGORITHMS

# LIST OF NOTATIONS