# 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

Datum des Disputation: **XX.XX.201X**

Gutachter:
**Prof. Dr. Knut Reinert**, *Freie Universität Berlin, Deutschland*
**Prof. Dr. XXX XXX**, *XXX, XXX*

# Abstract

Bla bla bla.

# CONTENTS

# Part I

# APPROXIMATE STRING MATCHING

# Algorithms for Edit Distance Verification

# Part II

# FULL-TEXT INDICES

# Classic Full-Text Indices

# Compressed Full-Text Indices

# Online Approximate String Matching

4.1 Banded Myers' bit-vector algorithm

4.2 Increased bit-parallelism using SIMD instructions

# Part III

# APPLICATIONS

# Read Mapping

# String Similarity Search / Join

# Counting filters

## A.1   Ungapped seeds

## A.2   Gapped seeds

### A.2.1   Underlying principle

### A.2.2   Threshold computation

### A.2.3   Sensitivity computation

# APPENDIX B

# 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 Information derived from the published and unpublished work of others has been acknowledged in the text and a list of references is given.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 5, 2013

# LIST OF FIGURES

# LIST OF TABLES

# LIST OF ALGORITHMS

# LIST OF NOTATIONS