School of Computer Science

University of Technology Sydney

# Unsupervised learning for high performance compression of genomic data collections

A thesis submitted in partial fulfillment of
the requirements for the degree of
**Doctor of Philosophy**

by

## Tao Tang

June 2021

# CERTIFICATE OF ORIGINAL AUTHORSHIP

I, Tao Tang declare that this thesis, is submitted in fulfilment of the requirements for the award of Doctor of Philosophy, in the School of Computer Science at the University of Technology Sydney.

This thesis is wholly my own work unless otherwise referenced or acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

This document has not been submitted for qualifications at any other academic institution.

Signature of Candidate
Production Note: Signature
removed prior to publication

Date 11/03/2021

# Acknowledgments

# Contents

# List of Figures

# List of Tables

# List of Publications

Below is the list of journal papers associated with my PhD research:

**Journal Papers Published**

- **Tang, T.**, Liu, Y., Zhang, B., Su, B., & Li, J. (2019). Sketch distance-based clustering of chromosomes for large genome database compression. *BMC genomics*, 20(10), 1-9.

- **Tang, T.**, & Li, J. (2021). Transformation of FASTA files into feature vectors for unsupervised compression of short reads databases. *Journal of Bioinformatics and Computational Biology*, 2050048-2050048.

- **Tang, T.**, Hutvagner. G., Wang, W. & Li, J. Assembly-improved compression of genomic short reads via error correction: survey and advancement. Under review.

# Abstract

The advanced next-generation sequencing (NGS) technologies have launched a new era of all fields of genetics. However, the vast quantity of data generated by NGS technologies also proposed great challenges to data storage, transmission and analysis. In this thesis, we focus on the compression of multiple data collections of short reads and assembled genome, we also explore the relationship between compression, error correction and de novo assembly of short reads data. First, we introduce an efficient clustering-based reference selection algorithm for the compression of genome databases. This method clusters the genomes into subsets of highly similar genomes using MinHash sketch distance, then applies a two-level compression based on the clustering result. The compression ratio gain of our approach can reach up to 20-30% in most cases for the datasets from NCBI, the 1000 Human Genomes Project and the 3000 Rice Genomes Project.

Furthermore, we propose a new clustering-based method for the compression of short reads datasets. Our approach transforms each file into a feature vector for clustering, then compresses the files in the same group together to increase the total number of detected overlappings during compression. The experiments show that our method achieves 20%-30% improvements in compression ratio than the previous one-by-one compression.

Finally, we review the relationship between reference-free compression, MSA based error correction and de novo assembly of short reads data. We demonstrate that high quality error correction can significantly reduce the number of mismatched nucleotides during reference-free compression and

hence improve the final compression ratio. The experiment results verify our estimation and show that the same error correction also has a positive effect on de novo assembly in most cases. In addition, we also propose a path graph based method for compression of short reads datasets.