

Xinhui Cai Thesis

**Computational identification of within-host diversity of  
SARS-CoV-2 and its benefits for mRNA vaccine design**

**Xinhui Cai**

School of Computer Science  
Faculty of Engineering and Information Technology  
University of Technology Sydney

Dec 28<sup>th</sup>, 2022

# CERTIFICATE OF ORIGINAL AUTHORSHIP

I, Xinhui Cai, declare that this thesis is submitted in fulfillment of the requirements for the award of Master of Analytics, in the School of Computer Science, Faculty of Engineering and Information Technology 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.

This research is supported by the Australian Government Research Training Program.

Signature: Production Note:  
Signature removed prior to publication.

Date: 28/12/2022

# Acknowledgment

Through out the time of this research I have received a great amount of help and support from many people, without them this research would be impossible. The first person I'd like to thank is my principle supervisor, prof.Jinyan Li. During my master studies, prof.Li helps me form the idea of identifying co-existing strains from the early SARS-CoV-2 sample and provides me valuable feedback and guidance at each stage of my research and helping me with the scholarship which reduces my financial burden considerably. I would also like to thank the other members in the team for their support during my research, especially Tian Lan and Pengyao Ping. Tian made a great contribution in obtaining the protein docking information and Pengyao develops an important part of the double-model error correction method. Dr.Xuan Zhang developed the InsEC which is a source of inspiration for the strain identification algorithm I designed, she also has given me some useful advice from time to time. My parents support me firmly during my research which I appreciate very much. The sequencing read information of SARS-CoV-2 reference genome sequence provided by Dr.Si Haorui from Prof.Zhengli Shi's lab is very important to my research and I appreciate it very much.

I would like to express my gratitude to my co-supervisor Prof.Brian Oliver for providing very helpful insights for both of my papers and playing an important role in forming some of the ideas presented in the papers. His knowledge in virology and immune system has been very helpful.

The Australian Research Council has provided financial support to this research through the scholarship I received. I would like to thank the organization for its recognition and generosity.

I also want to thank the panel of my candidate assessments, Prof.Wei Liu, Prof.Longbing Cao, Prof.Jinyan Li and Prof.Brian Oliver, who have shared their valuable opinions on each stage of my research.

Finally, for people special to me and support me throughout my study and life, I am very grateful and happy to have them.

Xinhui Cai

December 2022 in Sydney

## Abstract

Since the emergence of COVID-19 in Wuhan in late 2019, the world has been largely affected. Despite extensive research about the virus itself (SARS-CoV-2) causing the pandemic, some questions regarding its emergence are still unaccounted for. This paper contains an investigation and analysis of the sequencing reads and the SARS-CoV-2 reference genome sequence assembled from those reads. By retracing the steps of assembling the reference sequence, a probability that multiple strains of SARS-CoV-2 co-existed inside the patient's body is found. The assembly tool, MEGAHIT, applies an assembly process that tends to ignore multiple routes to form contigs. Therefore, we design a workflow that could rectify this potential issue and identify multiple strains from the COVID-19 patient sample. It involves error correction, extracting relevant reads, strain identification, phylogenetic study, and protein structure analysis. The results indicate that more than one strain of SARS-CoV-2 could be produced from the sample that was used to produce the reference sequence. Their binding affinity and phylogenetic relationships with the published SARS-CoV-2 reference genome, SARS-CoV, and some other variants of SARS-CoV-2 are also revealed. The discovered strains show some possible structural differences that affect the protein binding affinity between the spike protein and human ACE2. Consequently, this workflow for identifying within-host diversity highlights the existence of co-existing strains with distinct nucleotide sequences, emphasizing the importance of considering these variations when designing mRNA vaccines.

**Keywords:** SARS-CoV-2 strains; error correction; Illumina short reads; phylogenetic study;mRNA vaccine design

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Background . . . . .	1
1.2	Research Objectives . . . . .	2
1.3	Research Significance . . . . .	2
1.4	Motivation . . . . .	3
1.5	Thesis Structure . . . . .	5
1.6	Research Contribution . . . . .	6
<b>2</b>	<b>Literature review</b>	<b>7</b>
2.1	Error Correction Algorithm . . . . .	7
2.1.1	Instance-based Error Correction . . . . .	7
2.1.2	karect . . . . .	9
2.2	mRNA vaccines design and properties . . . . .	10
2.2.1	mRNA vaccines properties . . . . .	10
2.2.2	Comparison of mRNA vaccines with other types of vaccines . . . . .	11
2.3	mRNA vaccines designed for COVID-19 . . . . .	13
2.3.1	mRNA vaccines designed for preventing other diseases . . . . .	17
2.4	Existed Within-host Diversity Identification Tools . . . . .	21
2.5	RNA viruses' quasi-spices and within-host diversity . . . . .	22
<b>3</b>	<b>Computational discovery of intra-host coexisting strains of the SARS-CoV-2 reference strain: Motivation and methods</b>	<b>24</b>
3.1	Data Accessibility . . . . .	24
3.2	Workflow . . . . .	24
3.2.1	Non-human reads extraction . . . . .	26
3.2.2	Error correction for the non-human reads . . . . .	26
3.2.3	New strain identification . . . . .	27

3.2.4	Phylogenetic Study . . . . .	29
3.2.5	Protein binding affinity investigation . . . . .	30
<b>4</b>	<b>Computational discovery of intra-host coexisting strains of the SARS-CoV-2 reference strain: Results</b>	<b>34</b>
4.1	MEGAHIT Analysis . . . . .	34
4.2	Discovered Strain Differences . . . . .	35
4.3	Comparison with Other Tools . . . . .	36
4.4	Phylogenetic Relationship Analysis . . . . .	38
4.5	Sequences Comparison . . . . .	44
4.6	Protein Binding Affinity Analysis . . . . .	64
4.7	Case Study: Co-existing Strain Identification of SARS-CoV-2 in Wastewater Sample from California . . . . .	66
4.8	Foot-and-Mouth Disease Virus Animal Samples . . . . .	79
<b>5</b>	<b>Conclusions and future research perspectives</b>	<b>82</b>
<b>6</b>	<b>Current Publications and Future Plan</b>	<b>85</b>
<b>7</b>	<b>Appendix</b>	<b>86</b>
7.1	Supplemental Material . . . . .	98
7.1.1	SRR11092062 Discovered Strain 1 Spike Protein Sequence . . . . .	99
7.1.2	Difference of SRR11092062 Discovered Strain 1 Spike Protein Sequence and SARS-CoV-2 Spike Protein Sequence . . . . .	99
7.1.3	SRR11092062 Discovered Strain 2 Spike Protein Sequence . . . . .	99
7.1.4	Difference of Discovered Strain 2 Spike Protein Sequence and SARS-CoV-2 Spike Protein Sequence . . . . .	100
7.1.5	SRR11092062 Discovered Strain 1 Nucleotide Sequence . . . . .	100
7.1.6	SRR11092062 Discovered Strain 2 Nucleotide Sequence . . . . .	106
7.1.7	SRR12596175 Assembled Contig Spike Protein Sequence . . . . .	116
7.1.8	SRR12596175 Discovered Strain 1 Possible Spike Protein Sequence .	117
7.1.9	SRR12596175 Discovered Strain 2 Possible Spike Protein Sequence .	118
7.1.10	FMDV sample sequence differences . . . . .	118
7.1.11	SRR12596175 Assembled Contig Nucleotide Sequence . . . . .	144
7.1.12	SRR12596175 Discovered Strain 1 Nucleotide Sequence . . . . .	155
7.1.13	SRR12596175 Discovered Strain 2 Nucleotide Sequence . . . . .	165