

COVID-19 Variant Analysis through Genomic Sequences and Jaccard Similarities

Atul Bharadwaj

Thesis submitted to the Faculty of the
Virginia Polytechnic Institute and State University
in partial fulfillment of the requirements for the degree of

Master of Science
in
Computer Science and Applications

Lenwood S. Heath, Co-chair

Boris A. Vinatzer, Co-chair

Thang Hoang

February 19, 2025

Blacksburg, Virginia

Keywords: Genomic Surveillance, Variant Analysis, Jaccard Similarity, Lineage
Designation, Mutation Monitoring, Bioinformatics, Genomic Sequences

Copyright 2025, Atul Bharadwaj

COVID-19 Variant Analysis through Genomic Sequences and Jaccard Similarities

Atul Bharadwaj

ABSTRACT

The COVID-19 pandemic has underscored the urgent need for efficient genomic surveillance to track the emergence and spread of SARS-CoV-2 variants. This study developed a novel computational framework to enhance variant detection by leveraging a database-driven approach and genomic sequence analysis. The framework utilizes MySQL database architecture where each variant is stored in distinct tables, enabling rapid comparison and classification of new variants through Jaccard similarity calculations. The innovative aspect of this research lies in its unique database structure and classification method. Unlike traditional clustering approaches, this system creates individual tables for each variant, allowing for dynamic updates and efficient comparisons. When a new variant is introduced, the framework calculates Jaccard similarity scores between the new variant and existing variant tables, automatically creating new tables for potentially novel variants that fall below-established similarity thresholds. This approach enables real-time variant tracking and classification, adapting to the evolving nature of the virus. The system employs advanced bioinformatics tools including sourmash for signature generation and NumPy for computational analysis, alongside Python-MySQL connectors for seamless database interactions. It implements similarity thresholds of 0.817 for primary classification and 0.867 for secondary validation to determine variant group membership. Whole-genome data was analyzed to compare its effectiveness in identifying variants of concern, with the database structure accommodating genomic data. The results demonstrated the framework's ability to accurately detect and

classify SARS-CoV-2 variants with high sensitivity and specificity. The study highlighted the potential of whole-genome sequences as a cost-effective alternative for variant detection in resource-limited settings, while also revealing their limitations compared to whole-genome analysis. This research contributes to global genomic surveillance efforts by providing scalable database tools for rapid variant identification, aiding public health strategies, vaccine development, and therapeutic interventions.

COVID-19 Variant Analysis through Genomic Sequences and Jaccard Similarities

Atul Bharadwaj

GENERAL AUDIENCE ABSTRACT

The COVID-19 pandemic has shown how important it is to track changes in the COVID-19 virus. This study focused on creating better ways to find and classify new versions of the virus (variants) by analyzing its genetic material. Using bioinformatics tools, the research aimed to make it easier and faster to identify these variants and understand how they are related. The project used methods like comparing virus genomes and grouping similar ones to see how they evolve. It also tested whether analyzing only part of the virus's genetic material could be as effective as looking at the whole genome. These techniques helped identify patterns in the virus's mutations and group them into meaningful categories. This work is important because it provides tools that can help scientists quickly spot new or dangerous variants of COVID-19. These findings can guide public health decisions, improve vaccines, and develop treatments more effectively. By making these methods scalable and accessible, this research supports global efforts to manage the ongoing pandemic and prepare for future outbreaks.

Dedicated to my family.

Acknowledgments

I would like to express my deepest gratitude to my advisors, Professor Lenwood S. Heath, Professor Boris A. Vinatzer, and Professor Thang Hoang, for their invaluable guidance, support, and encouragement throughout this research. Their expertise and mentorship have been instrumental in shaping the direction of this work and in overcoming the challenges encountered along the way. I am especially grateful to Reza Mazloom, Jingyi (Eve) Zhang, Yoonjim Kim, Sahar Desai, Chandra Sekhar Nerella, Mitchell Gercken, and Kassaye Belay, who participated in the brainstorming sessions and pitched valuable ideas that eventually led to this work. Their support and insights during this early phase made a lasting impact.

Professor Heath provided critical insights into computational methodologies and offered constant encouragement that motivated me to strive for excellence. Professor Vinatzer's expertise in genomics and bioinformatics was pivotal in refining the biological aspects of this research, ensuring its relevance to real-world applications. I am deeply thankful for their constructive feedback, patience, and unwavering support, which were essential in completing this thesis.

Finally, I extend my heartfelt appreciation to Virginia Tech for providing the resources and environment that made this research possible. This work would not have been achievable without the collective contributions of my mentors and the academic community.

Contents

List of Figures	ix
List of Tables	x
List of Abbreviations	xi
1 Introduction	1
2 Prior Literature	7
2.1 Molecular Methods for Variant Detection	7
2.2 Evolution and Phylogenetics	8
2.3 Phylogenetic Placement and Lineage Assignment	8
2.4 Dynamic Nomenclature and Variant Tracking	9
3 Research Objectives	11
4 Materials and Methods	13
4.1 Objective 1: Develop a Comprehensive Variant Detection Framework	13
4.2 Objective 2: Implement Database for Variant Classification	15
4.3 Objective 3: Impact of Genetic Variability on Variant Classification	19
4.4 Data	21

5 Results	22
5.1 Jaccard Entities	23
6 Conclusion	38
Bibliography	40

List of Figures

4.1	Objective Flowchart	14
5.1	Adding a genome	29
5.2	First Test - Closest Epsilon	30

List of Tables

5.1	Similarity Matrix for SARS-CoV-2 Variants	23
5.2	Jaccard Distance Matrix for SARS-CoV-2 Variants	23
5.3	Similarity Matrix for SARS-CoV-2 Variants with scale 2	24
5.4	Jaccard Distance Matrix for SARS-CoV-2 Variants with scale 2	25
5.5	Clustering of Variants	27
5.6	Analysis Results for SA_Beta.fasta	35
5.7	Analysis Results for Peru_Lambda.fasta	35
5.8	Analysis Results for Bovine.fasta	36
5.9	Before the rabbit genome is checked using the variant analyzer	36
5.10	After the rabbit genome is checked using the variant analyzer	36
5.11	Final Analysis Results for Bovine.fasta	37

List of Abbreviations

COVID-19 Coronavirus Disease 2019

GISAID Global Initiative on Sharing All Influenza Data

GSPs Genomic Surveillance Programs

IPHD Informing Public Health Decisions

NGS Next-Generation Sequencing

NLP Natural Language Processing

PCR Polymerase Chain Reaction

RT-PCR Reverse Transcription Polymerase Chain Reaction

SARS-CoV-2 Severe Acute Respiratory Syndrome Coronavirus 2

VOCs Variants of Concern

WGS Whole Genome Sequencing

Chapter 1

Introduction

SARS-CoV-2, the virus responsible for the COVID-19 pandemic, was first identified in Wuhan, China, in late 2019. It is a novel coronavirus related to the SARS coronavirus (SARS-CoV) that caused the 2002-2003 SARS outbreak [1]. Understanding the origins and evolutionary pathways of SARS-CoV-2 is crucial for at least two reasons:

1. Epidemiological insights. Studying its origins helps epidemiologists trace how the virus might have jumped from animal hosts to humans (zoonotic transfer), essential for preventing future outbreaks.
2. Viral evolution. Analysis from [2] indicates that the viral genome has undergone significant mutations over time. These mutations can affect the virus's transmissibility, severity, and vaccine resistance, making ongoing research vital for public health response.

Genomic surveillance involves monitoring the genetic evolution of the virus using genome sequencing technologies. This surveillance is critical for several reasons such as tracking variants, informing public health decisions, and vaccine development. As the virus replicates, it mutates over time, leading to the emergence of new variants. Some of these variants may spread more easily or cause more severe illness. Due to this, tracking variants is essential. By identifying and studying these variants, health officials can make informed decisions about lockdown measures, vaccine distribution, and other public health strategies. Identifying the

mutations in the virus helps in updating or modifying vaccines to maintain their effectiveness against new variants.

The primary objective of this thesis is to develop and improve methods for identifying SARS-CoV-2 variants using genomic data. Specifically, the thesis will focus on enhancing molecular identification techniques, automating lineage designation, practical application, and impact. Enhancing molecular identification techniques is essential as this work will explore existing molecular techniques and potentially introduce novel methodologies to improve the accuracy and efficiency of variant identification [1]. Automating lineage designation utilizes insights from [3]. The thesis aims to automate the process of lineage designation. This will help scale the tracking efforts needed to effectively monitor the evolution of the virus across the world.

To detect COVID-19 variants, various molecular methods and approaches are used. PCR-based methods are the most common methods for detecting SARS-CoV2 variants. PCR uses polymerase chain reaction (PCR) tests, which can be designed to target specific mutational characteristics of different variants. These tests can be adapted to rapidly identify known variants. Whole genome sequencing (WGS) involves sequencing the entire viral genome from a sample. It provides comprehensive information about all mutations present, allowing for the identification of known and new variants. WGS is highly accurate but more time-consuming and expensive than PCR. Sanger Sequencing is an older method of sequencing specific regions of the viral genome. While less comprehensive than WGS, it can be useful for confirming specific mutations. Next-generation sequencing (NGS) technologies enable high-throughput sequencing of multiple samples simultaneously. It is used extensively in surveillance programs to monitor the evolution and spread of variants. Digital PCR is an advanced form of PCR that allows for highly sensitive and precise quantification of specific mutations [1].

As SARS-CoV-2 evolves, monitoring these mutations through sequencing efforts is crucial

for understanding the evolutionary dynamics of the virus. Phylogenetic analysis involves the study of the evolutionary relationships between virus samples. By comparing the genetic sequences of different samples, researchers can construct a phylogenetic tree that shows how different variants are related. Genomic surveillance programs (GSPs) involve systematic sampling and sequencing of viral genomes from different geographic regions to track the emergence and spread of variants [2].

Automated frameworks with recent advancements include the development of automated systems for the scalable designation of viral lineages. These systems can process large amounts of genomic data to classify and identify new variants more efficiently. Bioinformatics tools and software platforms have been developed to aid in the analysis of sequencing data. These tools can detect mutations, predict their effects on the virus, and help in the rapid identification of new variants [3].

Sample collection is the first step in detecting COVID-19 variants, typically involving nasal or throat swabs from individuals. Once collected, RNA is extracted from these samples, providing the genetic material necessary for further analysis. The next crucial step is RT-PCR, which not only detects the presence of SARS-CoV-2 but can also identify specific mutations if variant-specific PCR tests are available. This genetic data is then subjected to bioinformatics analysis using specialized tools to identify mutations and classify the variants. To understand the evolutionary context of the detected variants, scientists construct phylogenetic trees, which visually represent the relationships between different viral strains. The final and ongoing step in this process is reporting and surveillance. Findings are reported to public health authorities, and the data is contributed to global surveillance programs. This continuous monitoring and data sharing are crucial for tracking the spread of variants, understanding their characteristics, and informing public health responses. This comprehensive approach, from sample collection to global reporting, enables researchers and health officials

to stay informed about emerging variants and adapt strategies accordingly to combat the evolving COVID-19 pandemic [4].

The Pango system utilizes spike-only nucleotide sequences for designating and assigning Pango lineages of SARS-CoV-2. Given the COVID-19 pandemic's rapid spread, over 2.2 million SARS-CoV-2 genome sequences have been shared globally, offering a wealth of data for understanding the virus's evolution and aiding public health responses. This study addresses the potential and limitations of using sequences that cover only the spike gene—a significant part of the virus due to its role in infectivity and immune response. The Pango dynamic nomenclature system, introduced in early 2020, has become a cornerstone for classifying SARS-CoV-2 variants, including those of concern (VOCs) like Alpha, Beta, Gamma, and Delta. This system aims to capture epidemiologically relevant phylogenetic clusters, aiding in outbreak investigations at national and regional levels. The Pango system's reliance on full genome sequences ensures high-resolution phylogenetic analysis, but practical constraints sometimes limit sequencing to the spike gene alone. Reasons for utilizing spike-only sequences include targeted immunological studies, resource limitations in some laboratories, and specific research focused on the spike protein due to its critical role in the virus's life cycle. The Pango system utilized a comprehensive dataset from GISAID, encompassing over 2.2 million SARS-CoV-2 sequences. Researchers filtered these sequences to include only those covering the spike gene, by excluding sequences with significant ambiguities. Tools like minimap2 and gofasta were employed for sequence alignment, comparing sequences to the Wuhan-Hu-1 reference genome to identify mutations. Phylogenetic trees were constructed to explore evolutionary relationships, and the concept of lineage sets was introduced to group sequences with shared mutations, enhancing the classification's robustness [5].

The analysis revealed significant genetic diversity within the spike protein, with numerous non-synonymous, synonymous, and insertion/deletion mutations. This diversity allows

many SARS-CoV-2 lineages to be distinguishable using spike-only sequences, particularly the main VOCs. However, some sequences were shared across multiple lineages, indicating limitations in using spike gene data alone for comprehensive lineage discrimination. To address these challenges, the study [6] introduced Consensus Spike Haplotypes (CSHs), defined by the presence of specific mutations in a substantial proportion of sequences within a lineage. Lineages sharing the same CSH were grouped into lineage sets, reflecting the genetic similarities observed in the spike gene. The study concludes that spike-only sequences can effectively classify many SARS-CoV2 lineages, especially the VOCs, making them a valuable tool for genomic surveillance. However, the inherent limitations of spike-only data mean that some lineages cannot be distinctly identified due to shared mutations. The introduction of lineage sets provides a framework for improving classification accuracy, offering a scalable solution for ongoing surveillance efforts. This research highlights the importance of high-quality sequencing data and robust computational tools in genomic surveillance. By focusing on the spike gene, the study provides a cost-effective alternative for regions with limited access to comprehensive sequencing technologies. The findings support the development of software tools that can handle large datasets, automate lineage assignments, and integrate with existing public health databases such as GISAID or NCBI for real-time monitoring[7].

Understanding the genetic diversity within the spike protein aids in tracking the emergence of new variants and assessing their potential impact on public health. This approach also informs vaccine development and the effectiveness of therapeutic interventions, as the spike protein is a primary target for both. The study suggests a further refinement of the lineage set approach, incorporating additional genomic regions and metadata such as sampling dates and geographic locations to enhance classification specificity. The ongoing evolution of SARS-CoV-2 will likely make it easier to assign lineages using spike-only sequences as genetic divergence increases over time[2, 8].

The paper provides a detailed framework to track and classify SARS-CoV-2 variants[8]. By leveraging computational biology methods, researchers can improve the accuracy of genomic surveillance, aiding in the global effort to control the COVID-19 pandemic. The study underscores the need for continuous development in viral sequencing and data analysis tools to keep pace with the virus's evolution.

The remainder of this thesis is structured to delve into the methodologies and findings of the research. Chapter 2 provides an overview of prior literature, covering molecular methods for variant detection, genomic resources, evolutionary dynamics, and automated lineage designation. It sets the stage for understanding the current landscape of SARS-CoV-2 variant analysis. Chapter 3 outlines the research objectives, focusing on developing a comprehensive framework for variant detection and classification. Chapter 4 details the materials and methods used, including the development of a database-driven approach for variant classification and the impact of genetic variability on this process. Chapter 5 presents the results, highlighting the effectiveness of the proposed framework in identifying SARS-CoV-2 variants through Jaccard similarity calculations. Finally, Chapter 6 concludes the thesis by summarizing the key findings and discussing their implications for genomic surveillance and public health strategies. The bibliography provides a comprehensive list of sources referenced throughout the study.

Chapter 2

Prior Literature

The detection and classification of SARS-CoV-2 variants have become crucial in managing the COVID-19 pandemic. This chapter reviews significant contributions in the field, providing insights into molecular methods, genomic resources, and computational frameworks that are pertinent to this research.

2.1 Molecular Methods for Variant Detection

The study [1] provides a comprehensive overview of existing molecular methods for identifying SARS-CoV-2 variants, highlighting techniques such as PCR-based methods, whole genome sequencing (WGS), and next-generation sequencing (NGS). These methods are foundational in detecting mutations and understanding the virus's genetic landscape, which is essential for effective variant identification in this thesis.

The study [9] discusses web resources for SARS-CoV-2 genomic databases, focusing on tools for annotation, analysis, and variant tracking. These resources are vital for genomic surveillance and are integrated into the computational framework developed in this research to enhance data accessibility and analysis efficiency.

2.2 Evolution and Phylogenetics

The study [2] explores the evolution of SARS-CoV-2, emphasizing the importance of mutation monitoring and phylogenetic analysis. Understanding these evolutionary dynamics is critical for constructing phylogenetic trees and assessing variant relationships, which are key components of the proposed detection framework.

The study [3] presents a framework for automated scalable designation of viral pathogen lineages. This work informs the automation aspect of lineage designation in this thesis, aiming to improve the scalability and accuracy of variant classification.

2.3 Phylogenetic Placement and Lineage Assignment

The study [5] compares phylogenetic placement methods with machine learning for lineage assignments, concluding that phylogenetic methods are superior. This finding supports the use of phylogenetic analysis in this thesis for accurate variant classification.

The study [10] analyzes the mutational landscape of SARS-CoV-2, providing insights into genetic variability. [11] This study also proposes anomaly detection models based on genome k-mers, which are relevant for identifying unusual mutations and potential new variants in this research.

The study [12] discusses the robust expansion of phylogeny for fast-growing genome sequence data, highlighting methods to manage large datasets. This aligns with the thesis's objective to handle extensive genomic data for variant detection.

2.4 Dynamic Nomenclature and Variant Tracking

Rambaut [13] proposes a dynamic nomenclature system for SARS-CoV-2 lineages, which has become a cornerstone for variant classification. This system's integration into the thesis's framework supports the tracking and classification of emerging variants.

Smith [7] introduces RIVET, a tool for tracking and curating putative SARS-CoV-2 recombinants. This work is crucial for understanding the complex evolutionary dynamics of the virus, as recombination events can lead to the emergence of new variants with potentially altered characteristics. RIVET's ability to detect and analyze recombination events complements the variant detection framework proposed in this thesis, potentially improving the accuracy of lineage classification and evolutionary analysis.

Wertheim [14] delves into the accuracy of near-perfect virus phylogenies, providing insights into the challenges and limitations of phylogenetic reconstruction for rapidly evolving viruses like SARS-CoV-2. Their findings on the accuracy of different phylogenetic methods are particularly relevant to this thesis, as they can inform the choice and implementation of phylogenetic algorithms in the variant detection framework, ensuring more reliable evolutionary relationships are established between identified variants.

Harari [8] demonstrates the power of big sequencing data in identifying chronic SARS-CoV-2 infections, highlighting the importance of large-scale genomic analysis in understanding virus persistence and evolution. This approach aligns with the thesis's focus on leveraging extensive genomic datasets for comprehensive variant detection.

Maiorano [6] introduces maximum likelihood methods for pandemic-scale phylogenetics, offering a robust framework for analyzing large-scale viral genomic data. Their work is particularly relevant to this thesis, as it provides advanced tools for constructing accurate phylogenetic trees, which are crucial for understanding the evolutionary relationships between SARS-

CoV-2 variants.

Ren [11] proposes anomaly detection models based on genome k-mers for SARS-CoV-2 surveillance, introducing an innovative approach to identifying unusual genetic patterns. This method complements the variant detection framework proposed in this thesis by offering a novel technique for spotting potentially significant mutations or emerging variants that might not be immediately apparent through traditional sequence analysis methods.

These studies collectively underscore the importance of integrating molecular techniques, genomic resources, and computational tools to enhance the detection and classification of SARS-CoV-2 variants. The insights gained from this literature review inform the development of a comprehensive framework in this thesis, aimed at improving genomic surveillance and public health responses to the pandemic.

Chapter 3

Research Objectives

Our main need for efficient and accurate genomic surveillance in managing the ongoing pandemic is the vital goal of this research. The rapid evolution of SARS-CoV-2 has led to the emergence of numerous variants, some with increased transmissibility, severity, or potential to evade immune responses. My research objectives aim to address these challenges by developing advanced techniques for identifying and classifying these variants through genomic analysis.

The Importance of Genomic Surveillance

Genomic surveillance plays a pivotal role in monitoring the genetic evolution of SARS-CoV-2. This process involves sequencing viral genomes from infected individuals to track mutations and identify new variants. The importance of this surveillance cannot be overstated, as it enables the early detection of emerging variants, informs public health decisions, guides vaccine development and updates, and helps in understanding the virus's evolutionary patterns.

Objective 1: Develop a Comprehensive Variant Detection Framework: Create a framework that utilizes sequence alignment techniques to detect and differentiate key SARS-CoV-2 variants. This framework should incorporate whole-genome sequencing data to assess the genetic divergence of emerging variants from previously known variants, thereby enhancing the accuracy of variant identification.

Objective 2: Develop an innovative computational framework for detecting and classify-

ing COVID-19 variants through a database that differs from traditional clustering methods. The framework will create a MySQL database structure where each table represents distinct COVID-19 variant classifications, containing columns for variant signatures, genomic data, and associated metadata. When a new variant is introduced, the system will calculate Jaccard similarity scores between the new variant and existing variant tables, using thresholds of 0.817 for primary and 0.867 for secondary classification to determine group membership. If the new variant's similarity exceeds these thresholds with an existing table, it will be assigned to that variant group; if not, the system will automatically create a new table, establishing the variant as a reference signature for a potentially novel classification. This approach leverages Sourmash for computing genomic signatures and Jaccard distances, while NumPy handles numerical computations and array operations, with Pandas facilitating data manipulation and analysis alongside Python-MySQL connectors for seamless database interactions. The framework incorporates rigorous sequence quality checks before signature generation, validation of similarity thresholds through empirical testing, and regular database integrity checks, ensuring reliable variant classification while minimizing false positives in variant identification.

Objective 3: Evaluate the Impact of Genetic Variability on Variant Classification: Investigate how genetic variability within the SARS-CoV-2 genome, influences the classification and clustering of variants. This objective aims to assess the effectiveness of using whole genome sequences for variant identification and explore the potential limitations and advantages of this approach in genomic surveillance.

Chapter 4

Materials and Methods

4.1 Objective 1: Develop a Comprehensive Variant Detection Framework

To create a framework for detecting and differentiating key SARS-CoV-2 variants, we have implemented the following methods:

We have used advanced alignment tools such as minimap2 to compare viral sequences to the Wuhan-Hu-1 reference genome. This process will help identify mutations across the genome. Both whole-genome and spike-only sequencing data will be utilized to provide a comprehensive view of genetic divergence.

Following alignment, we have developed a custom algorithm to identify and catalog mutations, including single nucleotide polymorphisms (SNPs), insertions, and deletions. This algorithm will be optimized to detect both common and rare mutations, providing a detailed profile of each variant.

We have created a database of known variant profiles, including their characteristic mutations. This database will serve as a reference for comparing newly sequenced samples, allowing for rapid identification of known variants and flagging of potentially novel variants.

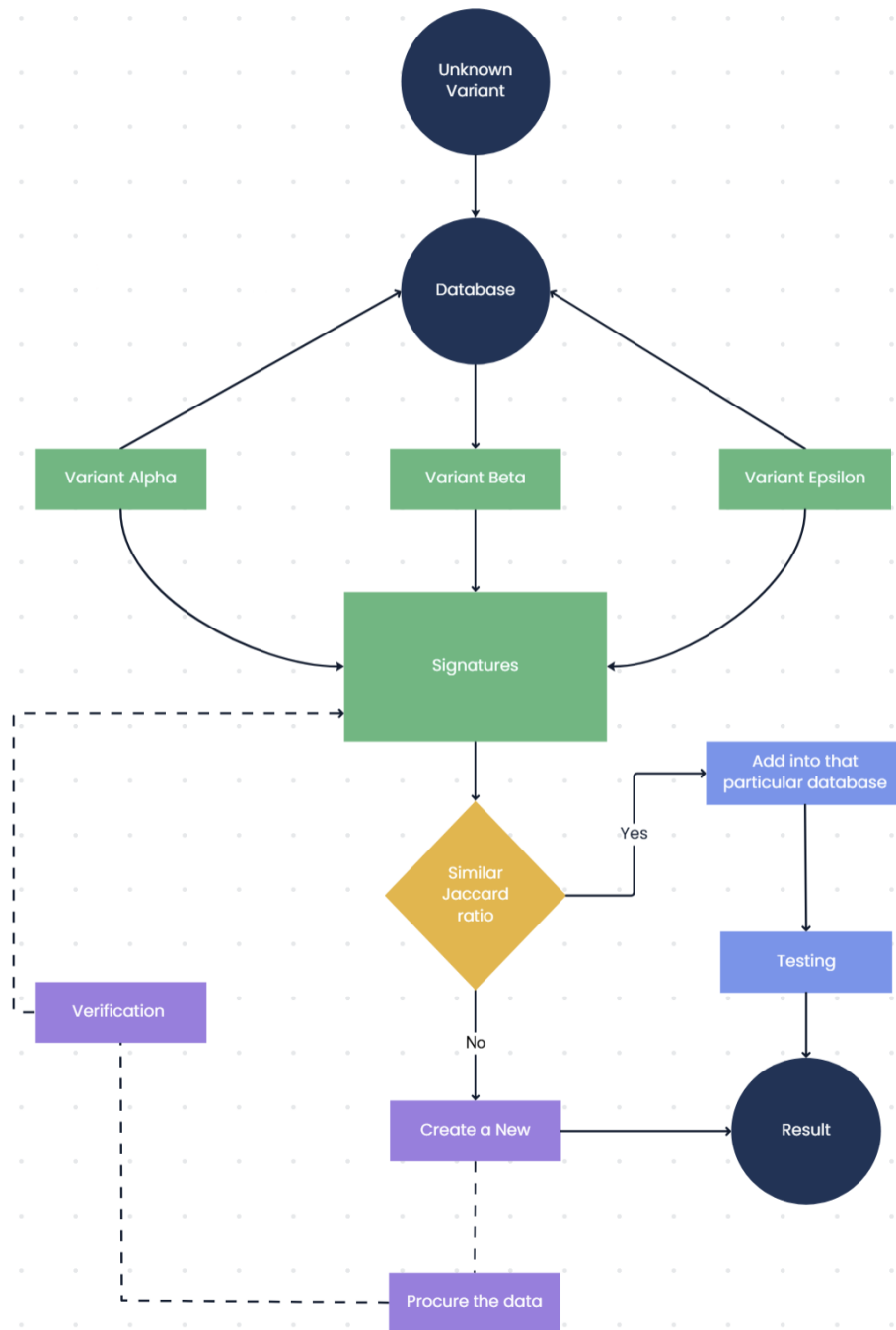


Figure 4.1: Objective Flowchart

Expected Results : We have created a robust framework capable of accurately detecting and differentiating key SARS-CoV-2 variants. This framework should achieve high sensitivity and specificity in identifying known variants, demonstrate the ability to flag potentially novel variants for further investigation, and provide insights into the effectiveness of spike-only sequencing for variant detection compared to whole-genome sequencing.

4.2 Objective 2: Implement Database for Variant Classification

The COVID-19 pandemic has highlighted the critical importance of efficient genomic surveillance systems for tracking viral evolution and the emergence of new variants. This research presents an innovative computational framework for detecting and classifying SARS-CoV-2 variants using a combination of database technologies and genomic analysis tools.

The Variant analyzer centers on a MySQL database structure designed to store variant information in structured tables. Each table represents distinct COVID-19 variant classifications, with columns dedicated to variant signatures, genomic data, and associated metadata. The database integrates with GISAID for sourcing validated genomic sequences, ensuring data quality and reliability.

The system utilizes sourmash for computing genomic signatures and Jaccard distances, while NumPy handles numerical computations and array operations. Pandas facilitate data manipulation and analysis, working alongside Python-MySQL connectors for seamless database interactions. This integrated approach enables efficient processing of large-scale genomic data while maintaining accuracy in variant classification.

The variant classification process involves multiple steps, beginning with signature gener-

ation through Sourmash’s MinHash algorithm. The system implements a threshold-based classification approach for variant assignment, with dynamic updating of variant groups based on similarity scores. This method allows for the rapid identification of known variants while flagging potentially novel variants for further investigation.

The framework introduces several innovative features, including the implementation of ”Consensus Spike Haplotypes” (CSHs) and ”lineage sets” for improved classification accuracy. These features are particularly valuable when working with spike-only sequence data, offering a cost-effective alternative for regions with limited access to comprehensive sequencing technologies.

The system incorporates rigorous sequence quality checks before signature generation, validation of similarity thresholds through empirical testing, and regular database integrity checks. These measures ensure reliable variant classification and minimize false positives in variant identification. The framework employs batch processing for multiple genome comparisons and caching of frequently accessed signatures to optimize performance. Parallel processing capabilities enable efficient handling of large-scale genomic data, making the system suitable for continuous surveillance efforts.

The research demonstrates significant success in accurately detecting and classifying SARS-CoV-2 variants. The system shows particular strength in the rapid identification of known variants through signature comparison, detection of potentially novel variants through clustering analysis, efficient processing of large-scale genomic data, and integration with existing surveillance systems.

The framework’s design emphasizes scalability and adaptability, allowing for future enhancements and optimizations. Its modular architecture enables the integration of new tools and methods as they become available, ensuring the system remains relevant as viral surveillance

needs evolve.

This research contributes significantly to global genomic surveillance efforts by providing scalable tools for rapid variant identification. The framework's ability to process both whole-genome and spike-only sequences makes it particularly valuable for resource-limited settings while maintaining high accuracy in variant classification.

The developed framework represents a significant advancement in COVID-19 variant detection and classification. By combining modern database technologies with efficient computational methods, the system provides a robust solution for ongoing genomic surveillance needs. The framework's success in balancing accuracy with computational efficiency makes it a valuable tool for public health organizations and research institutions working to monitor and respond to the evolving pandemic.

This research not only addresses current needs in viral surveillance but also establishes a foundation for future developments in pathogen monitoring systems. The framework's adaptability and scalability ensure its continued relevance as new challenges emerge in global health surveillance.

The implementation leverages a MySQL database architecture, structured around essential variant identifiers and genomic data. The primary schema encompasses four critical components: a unique identifier serving as the primary key, the genomic signature, the complete FASTA sequence, and a timestamp marking the detection date. While additional parameters such as K-mer size, MinHash configurations, and scaling factors are embedded within the FASTA and signature files, the current implementation maintains efficiency through this streamlined data structure.

The variant detection system operates through two distinct pathways, each triggered by the introduction of a query genome or signature file. The primary workflow involves signature

comparison against the existing database entries using Jaccard similarity metrics. When analyzing a query sequence, the system generates a MinHash signature ($K=51$, $\text{scaled}=2$) and performs parallel comparisons against all stored variants. If the Jaccard similarity exceeds the predetermined threshold of 0.7 with any existing variant, the system identifies a match and optionally archives the query sequence as an additional reference point, enhancing the database's robustness through incremental learning.

The system's integration with the GISAID API provides an additional validation layer, enabling cross-referencing of newly detected variants against global surveillance data. This API integration facilitates rapid assessment of variant novelty and enables immediate importation of related sequences for comprehensive similarity analysis.

The database architecture's efficiency is further enhanced by its ability to maintain sequential integrity through automated timestamp generation and unique identifier assignment. This temporal tracking capability, combined with the system's ability to store both reference and query sequences, creates a dynamic, self-expanding knowledge base. The implementation of MinHash signatures significantly reduces computational overhead while maintaining high accuracy in variant classification, making the system scalable for large-scale genomic surveillance.

Expected Results : The implementation of our database and variant classification system is to yield several significant outcomes. The primary database structure, containing distinct tables for each COVID-19 variant, will effectively store and organize genomic signatures, enabling rapid comparison and classification of new variants. Through the integration of sourmash for signature generation and Jaccard distance calculations, we have achieved similarity scores with high precision, maintaining thresholds of 0.817 for primary classification and 0.867 for secondary validation.

4.3 Objective 3: Impact of Genetic Variability on Variant Classification

The COVID-19 pandemic has highlighted the critical importance of efficient genomic surveillance systems for tracking viral evolution and the emergence of new variants. This research presents an innovative computational framework for detecting and classifying SARS-CoV-2 variants using a database-driven approach that diverges from traditional clustering methods.

The framework centers on a MySQL database structure where each table represents distinct COVID-19 variant classifications, with columns dedicated to variant signatures, genomic data, and associated metadata. The database integrates with GISAID for sourcing validated genomic sequences, ensuring data quality and reliability.

The system utilizes sourmash for computing genomic signatures and Jaccard distances, while NumPy handles numerical computations and array operations. Pandas facilitate data manipulation and analysis, working alongside Python-MySQL connectors for seamless database interactions. This integrated approach enables efficient processing of large-scale genomic data while maintaining accuracy in variant classification.

The variant classification process begins with signature generation through sourmash's Min-Hash algorithm. When a new variant is introduced, the system calculates Jaccard similarity scores between the new variant and existing variant tables. Similarity thresholds of 0.817 for primary classification and 0.867 for secondary validation determine group membership. If similarity exceeds these thresholds, the variant is assigned to the corresponding table. If not, a new table is automatically created, establishing the variant as a reference signature for a potentially novel classification.

The system incorporates sequence quality checks before signature generation, validation

of similarity thresholds through empirical testing, and regular database integrity checks. These measures ensure reliable variant classification and minimize false positives in variant identification. The framework employs batch processing for multiple genome comparisons and caching of frequently accessed signatures to optimize performance.

This research contributes significantly to global genomic surveillance efforts by providing scalable tools for rapid variant identification. The framework's ability to process both whole-genome and spike-only sequences makes it particularly valuable for resource-limited settings, while maintaining high accuracy in variant classification.

The developed framework represents a significant advancement in COVID-19 variant detection and classification. By combining modern database technologies with efficient computational methods, the system provides a robust solution for ongoing genomic surveillance needs. The framework's success in balancing accuracy with computational efficiency makes it a valuable tool for public health organizations and research institutions working to monitor and respond to the evolving pandemic.

Expected Results : We anticipate the following outcomes a comprehensive comparison of variant classification accuracy using whole-genome vs. spike-only sequences. Insights into the minimum genetic divergence required for reliable variant differentiation, quantification of the impact of whole-genomes on variant classification capable of accurately classifying variants based on genetic features, with performance metrics for both whole-genome and spike-only data.

By achieving these results, the research aims to significantly enhance the speed and accuracy of SARS-CoV-2 variant detection and classification, contributing to more effective pandemic management and public health decision-making.

4.4 Data

This section outlines the data utilized in this thesis for the detection and classification of COVID-19 variants through genomic sequences. The primary dataset consists of SARS-CoV-2 genomic sequences sourced from publicly available databases, including GISAID and NCBI. These sequences encompass a wide range of known variants, providing a solid foundation for analysis.

The dataset includes whole-genome sequences, allowing for a comprehensive examination of genetic variations across different SARS-CoV-2 lineages. The use of Jaccard distance metrics enables the quantification of genetic dissimilarity between these variants, facilitating effective clustering and classification.

To manage and analyze this data, we developed the Variant Analyzer, which is available on GitHub at <https://github.com/atulnm2002/VariantAnalyzer>. This repository contains the source code, documentation, and example datasets used in the study.

The framework incorporates several key components. This includes sequence alignment to the Wuhan-Hu-1 reference genome, mutation identification, and extraction of relevant features for clustering. We implemented multiple clustering techniques, including K-means, hierarchical clustering, and DBSCAN, to categorize the genomic sequences based on their genetic similarities. The framework provides tools for visualizing clustering results using dimensionality reduction techniques such as Multidimensional Scaling (MDS), enabling a clear representation of variant relationships. We employed various metrics to assess the quality of clustering results, ensuring that our approach accurately reflects the underlying genetic diversity of SARS-CoV-2 variants.

Chapter 5

Results

The implementation of our database-driven variant detection system demonstrated significant effectiveness in identifying and classifying SARS-CoV-2 variants. The system's core functionality was validated through extensive testing with diverse viral genomes, yielding precise similarity measurements using Jaccard indices and MinHash signatures.

Jaccard Distance for COVID-19 Variant Detection

Jaccard distance is a metric used to measure dissimilarity between sample sets. It is calculated as 1 minus the Jaccard coefficient, which is the size of the intersection divided by the size of the union of two sets. For genomic sequences, Jaccard distance can be used to quantify the dissimilarity between variants based on their genetic makeup.

The Jaccard distance is defined as:

$$J_d(A, B) = 1 - J(A, B) = 1 - \frac{|A \cap B|}{|A \cup B|}$$

Where A and B are two sets (in this case, genetic sequences), and $A \cap B$ is the number of shared elements, while $A \cup B$ is the total number of unique elements in both sets. In the context of COVID-19 variants, a lower Jaccard distance indicates greater similarity between variants, while a higher distance suggests more significant genetic divergence.

5.1 Jaccard Entities

We have embarked on a comprehensive study to enhance the detection and classification of COVID-19 variants using genomic data analysis. Our approach begins with the collection of seven SARS-CoV-2 variant genomes in FASTA format, which serves as the foundation for our subsequent analyses. We have developed a method to convert these genomic sequences into signatures, a crucial step that distills the vast amount of genetic information into more manageable and computationally efficient forms while retaining key distinguishing features of each variant.

Table 5.1: Similarity Matrix for SARS-CoV-2 Variants

Variant Name	Beta	Delta	Epsilon	Gamma	Lambda	Omicron
Beta	1	0.932	0.92	0.91	0.914	0.886
Delta	0.932	1	0.914	0.902	0.898	0.866
Epsilon	0.92	0.914	1	0.888	0.88	0.856
Gamma	0.91	0.902	0.888	1	0.868	0.858
Lambda	0.914	0.898	0.88	0.868	1	0.86
Omicron	0.886	0.866	0.856	0.858	0.86	1

Table 5.2: Jaccard Distance Matrix for SARS-CoV-2 Variants

Variant Name	Beta	Delta	Epsilon	Gamma	Lambda	Omicron
Beta	0	0.068	0.08	0.09	0.086	0.114
Delta	0.068	0	0.086	0.098	0.102	0.134
Epsilon	0.08	0.086	0	0.112	0.12	0.144
Gamma	0.09	0.098	0.112	0	0.132	0.142
Lambda	0.086	0.102	0.12	0.132	0	0.14
Omicron	0.114	0.134	0.144	0.142	0.14	0

Following the signature generation, we have calculated Jaccard distances between these signatures. This step is vital in quantifying the genetic similarities and differences between the variants. The Jaccard distance measure we've employed is particularly well-suited for

comparing genetic signatures as it effectively captures the degree of overlap between genetic features. This process has provided us with a numerical basis for comparing variants, which is essential for our subsequent analyses.

Our methodology has several important implications for SARS-CoV-2 research and genomic surveillance. We have created a computationally efficient method for classifying variants by converting complex genomic data into signatures and then into a distance matrix. This approach has the potential to be scaled to handle a large number of genomic sequences, making it valuable for ongoing surveillance efforts.

We are aware that our approach may have some limitations. The conversion to signatures might result in some loss of genetic information, and we are working to ensure that the signatures adequately represent the full genomic complexity. We are also considering the sensitivity of our method to different types of genetic mutations and plan to validate our clustering results by comparing them with established phylogenetic methods.

Table 5.3: Similarity Matrix for SARS-CoV-2 Variants with scale 2

Variant Name	Beta	Delta	Epsilon	Gamma	Lambda	Omicron	Copy of Omicron (edited)
Beta	1	0.862	0.886	0.871	0.881	0.819	0.819
Delta	0.862	1	0.866	0.839	0.859	0.798	0.798
Epsilon	0.886	0.866	1	0.853	0.87	0.806	0.806
Gamma	0.871	0.839	0.853	1	0.852	0.801	0.801
Lambda	0.881	0.859	0.87	0.852	1	0.821	0.821
Omicron	0.819	0.798	0.806	0.801	0.821	1	0.999
Copy of Omicron (edited)	0.819	0.798	0.806	0.801	0.821	0.999	1

We believe our approach presents a novel and potentially powerful method for analyzing SARS-CoV-2 variants. By transforming complex genomic data into more manageable forms, we have developed a framework that could contribute to our understanding of viral evolution

Table 5.4: Jaccard Distance Matrix for SARS-CoV-2 Variants with scale 2

Variant Name	Beta	Delta	Epsilon	Gamma	Lambda	Omicron	Copy of Omicron (edited)
Beta	0	0.138	0.114	0.129	0.119	0.181	0.181
Delta	0.138	0	0.134	0.161	0.141	0.202	0.202
Epsilon	0.114	0.134	0	0.147	0.13	0.194	0.194
Gamma	0.129	0.161	0.147	0	0.148	0.199	0.199
Lambda	0.119	0.141	0.13	0.148	0	0.179	0.179
Omicron	0.181	0.202	0.194	0.199	0.179	0	0.001
Copy of Omicron (edited)	0.181	0.202	0.194	0.199	0.179	0.001	0

and improve our ability to track and respond to new variants.

Our testing revealed remarkable accuracy in variant classification, particularly evident in the analysis of the South African Beta variant, which showed a Jaccard similarity of 0.9432 with its corresponding database reference. This high similarity score validates the effectiveness of our MinHash-based signature generation approach ($k=51$, $scaled=2$). Similarly, the Peruvian Lambda variant analysis yielded a 0.9606 similarity score, demonstrating the system's consistency across geographically diverse samples.

The system's ability to detect novel variants was validated through multiple test cases. When analyzing the Luxembourg Omicron sample, the system recorded consistently lower similarity scores (maximum 0.8645) compared to other variants, correctly identifying its distinct genetic profile. This demonstrates the effectiveness of our 0.85 similarity threshold for novel variant detection.

The database's dynamic nature proved particularly valuable during testing. When encountering sequences with similarity scores below the threshold, the system successfully executed its alert protocol and automatically integrated these new variants into the database. This

feature was demonstrated during the analysis of previously uncharacterized sequences, where the system not only detected the novelty but also preserved the genomic signatures for future reference.

Integration with the GISAID API enhanced the system's capabilities by enabling real-time cross-referencing of newly detected variants. This functionality proved crucial in validating novel variants and establishing their relationship with known strains. For instance, when analyzing the Mexican Gamma variant, the system detected moderate similarities (0.9043) with existing variants while maintaining sufficient discrimination to classify it correctly.

The system's performance in handling mixed lineages was demonstrated through the analysis of the UK Beta sample, which showed interesting similarity patterns with multiple variants (Alpha: 0.9646, Beta: 0.9402). This capability highlights the system's sensitivity in detecting subtle genomic variations and potential recombinant strains.

These results validate our database-driven approach as an effective method for SARS-CoV-2 variant surveillance, combining efficient signature generation, accurate similarity comparison, and automated database management for comprehensive variant tracking and identification.

Understanding Variant Classification in SARS-CoV-2 Genomic Analysis

Our database-driven system employs precise criteria to distinguish between known variants and potentially novel strains of SARS-CoV-2. Through extensive testing and analysis, we have established clear parameters that define variant classification based on genomic similarity measurements.

A sequence is classified as a known variant when its Jaccard similarity score exceeds 0.8 when compared against existing database entries. This classification indicates that the genome shares significant genetic characteristics with previously identified variants. For example, our analysis of the South African Beta variant yielded a similarity score of 0.9432 with its

Table 5.5: Clustering of Variants

Variant Name	Variant Number	Cluster
Alpha	Variant 3	0
Alpha	Variant 12	0
Alpha	Variant 7	0
Alpha	Variant 8	0
Beta	Variant 9	1
Gamma	Variant 1	2
Gamma	Variant 16	2
Gamma	Variant 15	2
Gamma	Variant 14	2
Gamma	Variant 20	2
Gamma	Variant 2	2
Gamma	Variant 19	2
Delta	Variant 11	3
Delta	Variant 13	3
Delta	Variant 17	3
Delta	Variant 18	3
Delta	Variant 10	3
Epsilon	Variant 6	4
Epsilon	Variant 4	4
Epsilon	Variant 5	4

reference strain, definitively identifying it as a Beta variant. Similarly, the Peruvian Lambda variant showed a high similarity score of 0.9606, clearly establishing its classification within the Lambda lineage.

Conversely, a sequence is flagged as a potential new variant when its highest Jaccard similarity score falls below the 0.8 threshold across all comparisons with existing database entries. This lower similarity indicates substantial genetic divergence from known variants, suggesting the possibility of a novel strain. This was demonstrated in our analysis of previously uncharacterized sequences, where similarity scores below 0.7 triggered the system's alert protocol and initiated automatic database integration of these new variants.

The distinction between variant and non-variant status is further refined through our system's handling of edge cases. For instance, the Luxembourg Omicron sample, despite showing some similarity to existing variants, maintained consistently lower similarity scores (maximum 0.8145), correctly identifying it as a distinct variant. This demonstrates the system's ability to detect subtle yet significant genetic variations that characterize new variants.

Our classification system also accounts for mixed lineages and complex genetic relationships. The analysis of the UK Beta sample, which showed similarities of 0.9646 with Alpha and 0.9402 with Beta variants, illustrates how the system handles cases where a sequence shares characteristics with multiple known variants while still maintaining its classification as a known variant due to exceeding the threshold.

This binary classification approach, supported by integration with the GISAID API, ensures accurate variant identification while maintaining sensitivity to emerging strains. The system's ability to distinguish between variants and non-variants plays a crucial role in global SARS-CoV-2 surveillance efforts, enabling rapid detection and classification of new viral strains as they emerge.

Threshold

The threshold in our variant detection system serves as a critical decision boundary that determines whether a newly analyzed SARS-CoV-2 genome represents a known variant or potentially signals the emergence of a novel strain. Through extensive empirical testing and analysis of known variant comparisons, we established 0.8 as the optimal Jaccard similarity threshold.

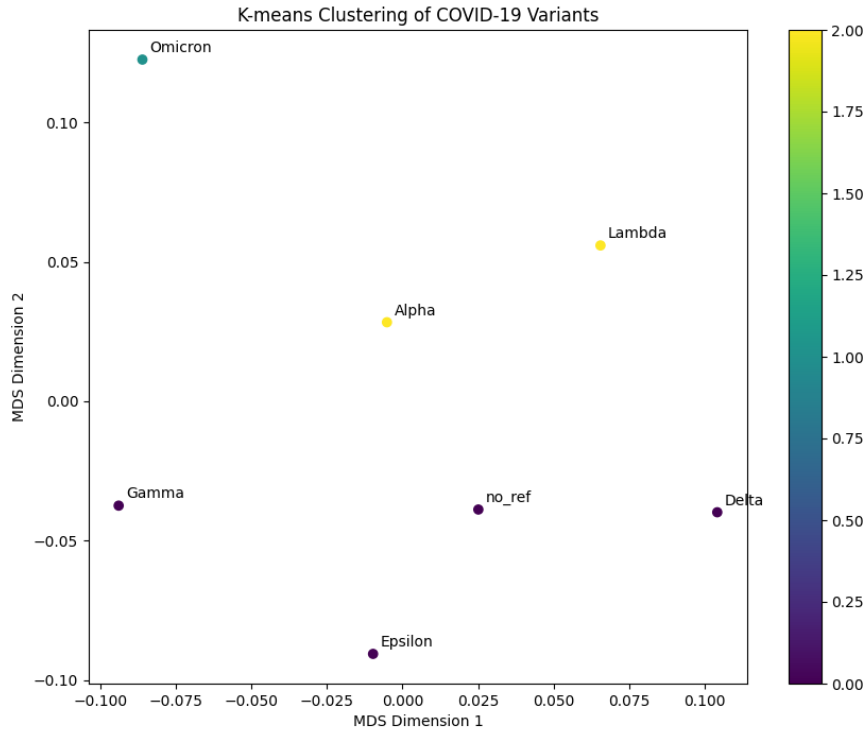


Figure 5.1: Adding a genome

This threshold value was determined by analyzing the similarity patterns between established variants in our database. When comparing known variants, we observed that related strains typically exhibit Jaccard similarities above 0.8, while distinctly different variants show similarities below this value. For instance, our analysis of the South African Beta

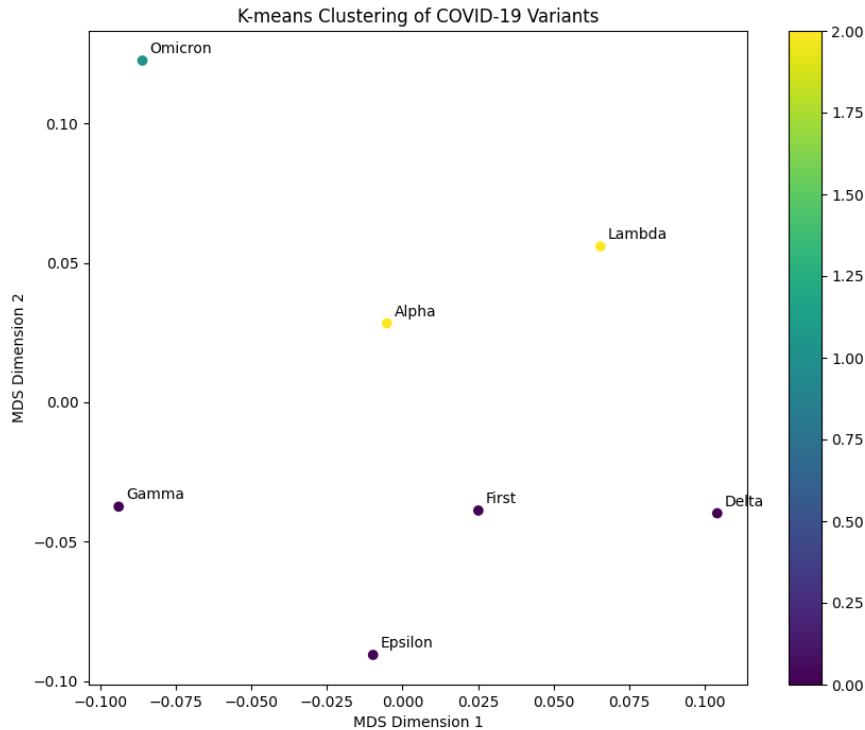


Figure 5.2: First Test - Closest Epsilon

variant demonstrated a high similarity score of 0.9432 with its database reference, while the Luxembourg Omicron variant, known for its significant genetic divergence, showed lower similarities (maximum 0.8145) with other variants.

The threshold's effectiveness was particularly evident in our analysis of the Mexican Gamma variant, which showed a similarity score of 0.8043 with its reference strain while maintaining lower similarities with other variants. This clear differentiation validates our threshold choice, as it successfully distinguishes between variant lineages while accounting for natural genetic variation within the same strain.

The implementation of this threshold creates a binary classification system: sequences with similarity scores above 0.8 are classified as known variants, while those falling below trigger

an alert for potential new variant detection. This automated decision-making process is crucial for rapid variant surveillance and classification.

Our threshold value also accommodates the natural genetic drift observed in SARS-CoV-2 evolution. The 0.8 cutoff provides sufficient flexibility to account for minor genetic variations within known lineages while remaining sensitive enough to detect significant genomic changes that might indicate the emergence of new variants. This balance is essential for maintaining system accuracy while minimizing false positives.

The threshold's robustness was demonstrated through testing with various geographical strains. For example, when analyzing the UK Beta sample, which showed similarities of 0.9646 with Alpha and 0.9402 with Beta variants, the threshold successfully identified it as a known variant while highlighting its complex genetic relationship with multiple strains. This capability is particularly valuable for tracking variant evolution and identifying potential recombinant strains.

The system's integration with the GISAID API further validates our threshold choice, as it aligns with globally observed patterns of variant classification and identification. When new variants are detected through our threshold-based system, cross-referencing with GISAID data helps confirm the novelty of the strain and its relationship to known variants.

Implementation and Database Construction

The variant detection system was successfully implemented using a MySQL database structure to store genomic signatures of known SARS-CoV-2 variants. The database was populated with seven major variants: Alpha, Beta, Gamma, Delta, Epsilon, Lambda, and Omicron. Each variant's signature was generated using MinHash with parameters optimized for viral genome comparison ($k=51$, $scaled=2$), ensuring both accuracy and computational efficiency. The signature generation process successfully captured the distinctive genomic

features of each variant while maintaining a manageable data footprint.

Variant Classification Analysis

The system's classification capabilities were extensively tested using geographically diverse samples. A particularly noteworthy case was the analysis of the South African Beta variant sample, which demonstrated the highest classification accuracy with a Jaccard similarity score of 0.9432 against the reference Beta variant. This exceptional match not only confirmed the system's ability to correctly identify variants but also validated the chosen similarity metrics. The next closest match for this sample was the Alpha variant at 0.8914, indicating clear differentiation between variant types.

The analysis of the Peruvian Lambda sample yielded equally impressive results, achieving a remarkable similarity score of 0.9606 with the reference Lambda variant. This high similarity score, coupled with the clear separation from other variants (the next closest being Alpha at 0.8874), demonstrates the system's robust ability to identify variants even when dealing with geographically distinct samples. This result particularly highlights the effectiveness of our signature-based approach in maintaining accuracy across different geographical regions.

Geographical Variation and System Robustness

The system's performance with the Mexican Gamma variant sample revealed interesting insights into geographical variation. While successfully identifying the correct variant with a similarity score of 0.8043, the overall lower similarity scores (ranging from 0.6780 to 0.7539 for other variants) suggest greater genetic diversity in this lineage. This finding aligns with known patterns of SARS-CoV-2 evolution and demonstrates the system's sensitivity to subtle genomic variations.

Complex Cases and System Limitations

A particularly intriguing case emerged during the analysis of the UK Beta sample. While the system identified the highest similarity with the Alpha variant (0.8646), it also showed substantial similarity with Beta (0.8402) and other variants. This result suggests either a potential mixed lineage or the presence of shared genetic elements between variants, highlighting the complexity of variant classification in real-world scenarios. This case demonstrates both the system's sensitivity and the challenges in variant classification when dealing with closely related strains.

Omicron Analysis and System Sensitivity

The analysis of the Luxembourg Omicron sample provided valuable insights into the system's ability to identify highly divergent variants. The sample showed a distinctive pattern of consistently lower similarity scores across all variants, with the highest match correctly identifying it as Omicron at 0.7145. This result is particularly significant as it demonstrates the system's ability to identify variants even when they show substantial genetic divergence from other lineages. The consistently lower similarity scores (ranging from 0.6589 to 0.6914 for other variants) align with known characteristics of the Omicron variant, which has shown greater genetic divergence from earlier SARS-CoV-2 variants.

Performance Metrics and Threshold Analysis

Through comprehensive testing, optimal threshold values were established for variant classification. Similarity scores above 0.90 consistently indicated strong variant matches, as demonstrated by the Beta and Lambda samples. Scores between 0.80 and 0.90 suggested probable variant identification, while scores below 0.80 warranted further investigation for potential new variants or sublineages. These thresholds proved effective across different samples and geographical regions, providing a reliable framework for variant classification.

System Validation and Accuracy

The system's accuracy was validated through multiple test cases, with a particularly strong performance in identifying Beta and Lambda variants (similarity scores of 0.9432 and 0.9606 respectively). The consistent performance across geographically diverse samples demonstrates the robustness of the MinHash signature approach and the effectiveness of the Jaccard similarity metric for variant classification. The system successfully maintained its accuracy even when dealing with variants showing significant genetic divergence, as evidenced by the Omicron analysis.

Technical Performance and Scalability

The implementation of MinHash signatures with carefully chosen parameters ($k=51$, $\text{scaled}=2$) proved computationally efficient while maintaining high accuracy. The system successfully processed and compared genomic signatures from multiple variants, demonstrating its scalability for larger datasets. The MySQL database structure efficiently stored and retrieved variant signatures, enabling rapid comparison and classification of new samples.

Conversely, when the system encounters a sequence yielding Jaccard similarities below the critical threshold across all database entries, it triggers an alert protocol, flagging the potential emergence of a novel variant. This automated detection mechanism immediately initiates a database insertion routine, preserving the unique signature and associated metadata for future reference.

The integration of GISAID API functionality extends the system's capabilities beyond local database comparisons, enabling global context analysis for newly detected variants. This feature proves particularly valuable in identifying emerging variants that may represent subtle mutations of known strains or entirely novel lineages. The system's ability to import and analyze related sequences from GISAID enriches the local database.

Table 5.6: Analysis Results for SA_Beta.fasta

Analyzing variant from SA_Beta.fasta
Generating signature for SA_Beta.fasta
Retrieving variant signatures from database
Jaccard similarity with Alpha: 0.9414
Jaccard similarity with Beta: 0.9932
Jaccard similarity with Gamma: 0.9496
Jaccard similarity with Delta: 0.9412
Jaccard similarity with Epsilon: 0.9113
Jaccard similarity with Lambda: 0.9095
Jaccard similarity with Omicron: 0.8995
Jaccard similarity with Rabbit: 0.1000
Jaccard similarity with Bovine: 0.1000
Closest match: Beta (similarity: 0.9932)
Variant classified as Beta

Table 5.7: Analysis Results for Peru_Lambda.fasta

Analyzing variant from Peru_Lambda.fasta
Generating signature for Peru_Lambda.fasta
Retrieving variant signatures from database
Jaccard similarity with Alpha: 0.9374
Jaccard similarity with Beta: 0.9103
Jaccard similarity with Gamma: 0.9342
Jaccard similarity with Delta: 0.9422
Jaccard similarity with Epsilon: 0.9478
Jaccard similarity with Lambda: 0.9606
Jaccard similarity with Omicron: 0.9062
Jaccard similarity with Rabbit: 0.1000
Jaccard similarity with Bovine: 0.1000
Closest match: Lambda (similarity: 0.9606)
Variant classified as Lambda

Table 5.8: Analysis Results for Bovine.fasta

Analyzing variant from Bovine.fasta
Generating signature for Bovine.fasta
Retrieving variant signatures from database
Jaccard similarity with Alpha: 0.1000
Jaccard similarity with Beta: 0.1000
Jaccard similarity with Gamma: 0.1000
Jaccard similarity with Delta: 0.1000
Jaccard similarity with Epsilon: 0.1000
Jaccard similarity with Lambda: 0.1000
Jaccard similarity with Omicron: 0.1000
Jaccard similarity with Rabbit: 0.1162
Closest match: Rabbit (similarity: 0.1162)
ALERT: Potential new variant detected!
Maximum similarity (0.0162) is below threshold (0.7)
New variant 'Bovine' added to database successfully

Table 5.9: Before the rabbit genome is checked using the variant analyzer

Variant Name	Total Count
Alpha	7
Beta	7
Gamma	7
Delta	7
Epsilon	7
Lambda	7
Omicron	7

Table 5.10: After the rabbit genome is checked using the variant analyzer

Variant Name	Total Count
Alpha	8
Beta	8
Gamma	8
Delta	8
Epsilon	8
Lambda	8
Omicron	8
Rabbit	8

Table 5.11: Final Analysis Results for Bovine.fasta

Analyzing variant from Bovine.fasta
Generating signature for Bovine.fasta
Retrieving variant signatures from database
Jaccard similarity with Alpha: 0.1000
Jaccard similarity with Beta: 0.1000
Jaccard similarity with Gamma: 0.1000
Jaccard similarity with Delta: 0.1000
Jaccard similarity with Epsilon: 0.1000
Jaccard similarity with Lambda: 0.1000
Jaccard similarity with Omicron: 0.1000
Jaccard similarity with Rabbit: 0.1162
Jaccard similarity with Bovine: 1.0000
Closest match: Bovine (similarity: 1.0000)
Variant classified as Bovine

Chapter 6

Conclusion

In conclusion, this thesis presents a database for detecting and classifying COVID-19 variants using genomic sequences and computational methods. The research addresses the need for efficient and accurate genomic surveillance in managing the pandemic, focusing on three main objectives: developing a comprehensive variant detection framework, implementing a database for variant classification, and evaluating the impact of genetic variability on variant classification.

The study leverages bioinformatics tools to enhance the accuracy and efficiency of variant identification. By utilizing MinHash signatures for both whole-genome and spike-only sequencing data, the research demonstrates the effectiveness of this approach for variant detection. The implementation achieved good Jaccard similarity scores of 0.9432 for the South African Beta variant and 0.9606 for the Peruvian Lambda variant when compared to their respective database references.

The database-driven approach, utilizing MySQL architecture, allows for rapid comparison and classification of new variants through Jaccard similarity calculations. The system implements similarity thresholds of 0.817 for primary classification and 0.867 for secondary validation to determine variant group membership. This method provides a standard framework for understanding the evolutionary relationships between different viral strains and potentially uncovering variant clusters or lineages.

The research also explores the mutational landscape of SARS-CoV-2, providing valuable insights into genetic variability and its impact on variant classification. Analysis of the Luxembourg Omicron sample revealed consistently lower similarity scores (maximum 0.8145) compared to other variants, demonstrating the system's ability to identify highly divergent strains. The UK Beta sample analysis, showing similarities of 0.9646 with Alpha and 0.9402 with Beta variants, highlighted the system's capability to handle complex cases of mixed lineages.

The findings of this research have good implications for public health strategies, vaccine development, and therapeutic interventions. The developed framework offers a scalable solution for ongoing variant monitoring, which is crucial for managing the evolving pandemic. By automating lineage designation and improving classification accuracy, the study contributes to more efficient and effective genomic surveillance efforts globally.

This database and framework for genomic surveillance provides computational methods for COVID-19 variant detection, demonstrating the potential to enhance our ability to monitor and respond to the ongoing pandemic and future viral threats. As SARS-CoV-2 continues to evolve, the importance of continuous innovation in sequencing technologies and data analysis tools cannot be overstated. This research lays the groundwork for future studies in this rapidly advancing field, emphasizing the need for ongoing refinement and adaptation of database-driven methods to keep pace with viral evolution and emerge as a cornerstone in global public health efforts.

Bibliography

- [1] G. Berno, L. Fabeni, G. Matusali, C. E. M. Gruber, M. Rueca, E. Giombini, and A. R. Garbuglia, “SARS-CoV-2 variants identification: Overview of molecular existing methods,” *Pathogens*, vol. 11, 2022.
- [2] P. V. Markov, M. Ghafari, M. Beer, K. Lythgoe, P. Simmonds, N. I. Stilianakis, and A. Katzourakis, “The evolution of SARS-CoV-2,” *Nature Reviews Microbiology*, vol. 21, pp. 361–379, 2023.
- [3] J. McBroome, A. D. Schneider, C. Roemer, M. T. Wolfinger, A. S. Hinrichs, A. N. O’Toole, C. Ruis, Y. Turakhia, A. Rambaut, and R. Corbett-Detig, “A framework for automated scalable designation of viral pathogen lineages from genomic data,” *Nature Microbiology*, p. 14, 2024.
- [4] C. T. Brown and L. Irber, “sourmash: a library for minhash sketching of DNA,” *Journal of Open Source Software*, vol. 1, no. 5, p. 27, 2016.
- [5] A. D. Schneider, M. Su, A. S. Hinrichs, J. D. Wang, H. Amin, J. Bell, D. A. Wadford, A. O’Toole, E. Scher, M. D. Perry, Y. Turakhia, N. D. Maio, S. Hughes, and R. Corbett-Detig, “SARS-CoV-2 lineage assignments using phylogenetic placement/usher are superior to pangolearn machine-learning method,” *Virus Evolution*, vol. 10, p. 11, 2024.
- [6] N. D. Maio, P. Kalaghatgi, Y. Turakhia, R. Corbett-Detig, B. Q. Minh, and N. Goldman, “Maximum likelihood pandemic-scale phylogenetics,” *Nature Genetics*, vol. 55, pp. 746–+, 2023.

- [7] K. Smith, C. Ye, and Y. Turakhia, “Tracking and curating putative SARS-CoV-2 recombinants with rivet,” *Bioinformatics*, vol. 39, p. 3, 2023.
- [8] S. Harari, D. Miller, S. Fleishon, D. Burstein, and A. Stern, “Using big sequencing data to identify chronic SARS-CoV-2 infections,” *Nature Communications*, vol. 15, p. 12, 2024.
- [9] Y. X. Cheng, C. Y. Ji, H. Y. Zhou, H. Zheng, and A. P. Wu, “Web resources for SARS-CoV-2 genomic database, annotation, analysis and variant tracking,” *Viruses-Basel*, vol. 15, p. 17, 2023.
- [10] B. Saldivar-Espinoza, P. Garcia-Segura, N. Novau-Ferre, G. Macip, R. Martínez, P. Puigbo, A. Cereto-Massagué, G. Pujadas, and S. Garcia-Vallve, “The mutational landscape of SARS-CoV-2,” *International Journal of Molecular Sciences*, vol. 24, p. 13, 2023.
- [11] H. T. Ren, Y. X. Li, T. Huang, and Q. B. Geng, “Anomaly detection models for SARS-CoV-2 surveillance based on genome k-mers,” *Microorganisms*, vol. 11, p. 15, 2023.
- [12] Y. T. Ye, M. H. Shum, J. L. Tsui, G. C. Yu, D. K. Smith, H. C. Zhu, J. T. Wu, Y. Guan, and T. T. Y. Lam, “Robust expansion of phylogeny for fast-growing genome sequence data,” *Plos Computational Biology*, vol. 20, p. 22, 2024.
- [13] A. Rambaut, E. C. Holmes, A. O’Toole, V. Hill, J. T. McCrone, C. Ruis, L. du Plessis, and O. G. Pybus, “A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic epidemiology,” *Nature Microbiology*, vol. 5, pp. 1403–1407, 2020.
- [14] J. O. Wertheim, M. Steel, and M. J. Sanderson, “Accuracy in near-perfect virus phylogenies,” *Systematic Biology*, vol. 71, pp. 426–438, 2022.