

GisGen: Technical Sheet

Executive Summary

GisGen represents a significant leap forward in the harmonization of virus sequence metadata between two of the most pivotal databases in genomics: GenBank and GISAID. With its advanced algorithms and user-centric design, GisGen, the latest version of this tool, provides an unparalleled solution for researchers looking to reconcile and enrich metadata for high-resolution epidemiological research. This report delves into the functionalities, importance, and operational framework of GisGen, offering users a comprehensive understanding of its capabilities and enhancements over previous versions.

Introduction to GisGen

GisGen is a sophisticated computational tool designed to address the critical challenge of harmonizing virus sequence metadata across the GenBank and GISAID databases. Utilizing fuzzy matching algorithms and MD5 checksum validation, GisGen ensures the accurate linkage of genomic sequences with their correct metadata counterparts, thereby facilitating enhanced epidemiological research and public health surveillance.

Importance of GisGen

The harmonization of metadata is pivotal in genomic epidemiology, where discrepancies in data across databases can significantly hinder accurate models of pathogen evolution and spread. GisGen emerges as a vital resource, providing a flexible, automated solution for the integration and enrichment of virus sequence metadata. This harmonization is crucial for developing comprehensive, accurate epidemiological models and supporting effective public health interventions.

Key Features and Enhancements

User Interface

GisGen boasts an improved Python-based graphical user interface (GUI), making it more intuitive and accessible for researchers. This enhancement simplifies the navigation and operation of the tool, facilitating a smoother workflow for users.

Operational Workflow

The operational workflow of GisGen has been streamlined into a user-friendly, six-step process, from data input to the export of reconciled metadata and sequences. This process is designed to be both efficient and comprehensive, accommodating various file formats and large data volumes. The tool supports real-time updates, providing users with timely feedback throughout the reconciliation process.

Performance and Capacity

To accommodate the tool's advanced functionalities, significant optimizations have been made to improve performance. Users are advised to process no more than 1,000 samples at a time and to operate GisGen on a machine with at least 20GB of RAM and a CPU performance comparable to an Intel Core i7-11700K. These recommendations are made to ensure optimal efficiency and reliability of the tool.

Data Processing and Downloading

A critical step in the GisGen workflow involves the downloading of multi-sample FASTA files from GISAID. Users are guided through this process, which is now explicitly detailed, ensuring clarity and ease in acquiring necessary sequence data.

Focus on Omicron Variants

Reflecting the evolving landscape of public health priorities, GisGen includes specific functionalities for filtering Omicron variants, highlighting the tool's adaptability to current research needs.

Output Data

GisGen generates comprehensive outputs, including detailed analyses of virus characteristics, collection data, host information, and fuzzy matching results.

These outputs are crucial for robust genomic surveillance and epidemiological analysis.

Conclusion

GisGen stands as a testament to the advancements in genomic data analysis tools, offering an essential resource for researchers striving to understand pathogen evolution and inform public health strategies. Its development reflects a commitment to innovation, user experience, and the provision of accurate, actionable genomic data. With GisGen, the scientific community has a powerful tool at its disposal, capable of significantly enhancing the efficiency and accuracy of genomic epidemiology efforts.

Availability

GisGen is freely available on GitHub, ensuring its accessibility to the wider research community. For detailed information on implementation and source code, users are encouraged to visit [GisGen on GitHub](#).

For further assistance and support, users can contact Dr. Matthew Scotch via email at matthew.scotch@asu.edu or Amir Elyaderani via email at aelyader@asu.edu. Supplementary data and tutorials are available online to facilitate the use and understanding of GisGen.