# RCoV19: A One-stop Hub for SARS-CoV-2 Genome Data Integration, Variants Monitoring, and Risk Pre-warning

Cuiping Li[1,#], Lina Ma[1,2,3,#], Dong Zou[1,2,#], Rongqin Zhang[1,3,4,#], Xue Bai[1], Lun Li[1], Gangao Wu[1,2,3], Tianhao Huang[1,2,3], Wei Zhao[1,2,3], Enhui Jin[1,2,3], Yiming Bao[1,2,3,*], Shuhui Song[1,2,3,4,*]

[1] *National Genomics Data Center, Beijing Institute of Genomics, Chinese Academy of Sciences and China National Center for Bioinformation, Beijing 100101, China*

[2] *CAS Key Laboratory of Genome Sciences and Information, Beijing Institute of Genomics, Chinese Academy of Sciences and China National Center for Bioinformation, Beijing 100101, China*

[3] *University of Chinese Academy of Sciences, Beijing 100049, China*

[4] *Sino-Danish College, University of Chinese Academy of Sciences, Beijing 100049, China*

[#] Equal contribution.

[*] Corresponding authors.

E-mail: songshh@big.ac.cn (Song S), baoym@big.ac.cn (Bao Y)

**Running title:** *Li C et al / Platform of SARS-CoV-2 Variants Monitoring and Pre-warning*

Total word counts: 5023

Total References: 30

Total figures: 7

## Abstract

The Resource for Coronavirus 2019 (RCoV19, https://ngdc.cncb.ac.cn/ncov/) is an open-access information resource dedicated to providing valuable data on the genomes, mutations, and variants of the severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2). In this updated implementation of RCoV19, we have made significant improvements and advancements over the previous version. Firstly, we have implemented a highly refined genome data curation model. This model now features an automated integration pipeline and optimized curation rules, enabling efficient daily updates of data in RCoV19. Secondly, we have developed a global and regional lineage evolution monitoring platform, alongside an outbreak risk pre-warning system. These additions provide a comprehensive understanding of SARS-CoV-2 evolution and transmission patterns, enabling better preparedness and response strategies. Thirdly, we have developed a powerful interactive mutation spectrum comparison module. This module allows users to compare and analyze mutation patterns, assisting in the detection of potential new lineages. Furthermore, we have incorporated a comprehensive knowledgebase on mutation effects. This knowledgebase serves as a valuable resource for retrieving information on the functional implications of specific mutations. In summary, RCoV19 serves as a vital scientific resource, providing access to valuable data, relevant information, and technical support in the global fight against COVID-19.


**KEYWORDS:** SARS-CoV-2; Mutation; Variants; Surveillance; Pre-warning

## Introduction

SARS-CoV-2 is responsible for the COVID-19 pandemic, and continues to evolve and spread to threat public health worldwide. Genome data play a crucial role in understanding mutations (refers to an actual nucleotide or amino acid change in a

58  viral genome), functions, and supporting the design of candidate vaccines. While

59  there are various data deposition repositories available, such as EpiCoV$^{TM}$ [1],

60  GenBank [2, 3], and GenBase (https://ngdc.cncb.ac.cn/genbase/), none of them

61  encompass all worldwide genome data, and redundancies exist among these

62  repositories. Therefore, the need for a comprehensive SARS-CoV-2 database arises to

63  integrate genome data, monitor evolution, and provide pre-warnings for high-risk

64  variants. Such a database is essential to comprehend the ongoing pandemic and

65  facilitate timely adjustments to public health interventions.

66      With millions of genome sequences now available, several platforms have been

67  developed to track SARS-CoV-2 mutations. These platforms, including COVID-19

68  CG [4], Outbreak [5], and CoV-Spectrum [6], enable tracking of mutations by

69  location, date of interest, and known variants globally. VarEPS [7] assesses the risk

70  level of mutations and variants based on their transmissibility and affinity to

71  neutralizing antibodies. Additionally, databases like CoV-RDB [8, 9] and

72  COG-UK-ME [10] have compiled mutations associated with reduced susceptibility to

73  various factors, such as clinical stage SARS-CoV-2 Spike monoclonal antibody

74  (mAb), RNA-dependent RNA polymerase (RdRP) inhibitor, 3C-like protease

75  (3CLpro) inhibitor, or mutations on T cell epitope. However, despite these significant

76  efforts, there are limitations in terms of efficiency and comprehensiveness. Most of

77  these platforms and databases only focus on specific aspects of SARS-CoV-2

78  monitoring or prevention.

79      Furthermore, numerous important mutations affecting transmissibility, infectivity,

80  or expression are scattered throughout published literature. Consequently, there is an

81  urgent need to build an integrated and comprehensive system that encompasses

82  "data-information-knowledge-application". This system should provide real-time

83  services for sequence monitoring, evolution tracking, and pre-warning of high-risk

84  variants.

85      RCoV19, previously known as 2019-nCoVR [11, 12], is an open-access

86  information resource for SARS-CoV-2. It has been available online and has already

87    provided data services to over 3.2 million visitors from 182 countries/regions

88    worldwide, with more than 14 billion data downloads in total. In this updated release

89    of RCoV19, significant improvements have been made in data curation, integration,

90    sequence growth and lineage evolution surveillance, and mutation comparisons of

91    sequences and lineages. Additionally, a weekly report on potentially high-risk

92    haplotypes (a distinct virus genome sequence) and variants (a viral genome that may

93    contain one or more mutations, which may affect virus's properties) is provided,

94    combining genetic mutation effects and haplotype network features [13, 14].

95    Furthermore, RCoV19 curates an integrated knowledge of mutation effects from

96    literatures and databases, offering critical insights into virus evolution, immune

97    escape, and medical countermeasures. Ultimately, RCoV19 establishes a one-stop hub

98    for SARS-CoV-2 genome data integration and variant monitoring, as illustrated in

99    **Figure 1**.

## Database content and features

### Efficient integration and retrieval of worldwide SARS-CoV-2 genome data

102    RCoV19 is an extensive data resource for SARS-CoV-2 that collects genome data

103    from multiple repositories, performs de-redundancy processing, and assesses

104    sequence quality to ensure a comprehensive and curated collection of worldwide

105    genomes (**Figure 2**). The resource incorporates data from repositories such as

106    EpiCoV$^{TM}$ [1], GenBank [2, 3], CNGBdb [15], and Novel Coronavirus Service

107    System of NMDC [16], and has included data from GenBase since the beginning of

108    2023. To eliminate redundancies, RCoV19 identifies identical genomes across

109    different sources and cross-references related accession IDs. Notably, 91.3%

110    GenBank sequences overlap with EpiCoVTM sequences, while 56.7% of EpiCoVTM

111    sequences are unique (**Figure S1**). It determines completeness of the protein-coding

112    region, evaluates sequences in five aspects (Ns, degenerate bases, gaps, mutations,

113    and mutation density) and defines high-quality sequences based on Ns and degenerate

114    bases. These processes enable RCoV19 to provide a comprehensive and reliable list

115    of SARS-CoV-2 genomes for global monitoring and pre-warning purposes.

116    In the new version, the SARS-CoV-2 genome data curation model has been

117    significantly enhanced with an automated integration pipeline and optimized curation

118    rules (**Figure 2**), ensuring efficient daily updates in RCoV19. The automated pipeline,

119    activated by a timer every day, collects genome data from various repositories through

120    the Chrome Browser on Linux, standardizes genome metadata, and performs

121    de-redundancy processing. This automated approach improves efficiency compared to

122    semi-automated methods and enables regular and constant updates. Curation rules

123    have also been optimized to achieve more accurate de-redundancy, by comparing

124    genome sequences (with removal of Ns and uniform letter case) in addition to key

125    metadata (virus name, sampling date and location). Furthermore, the curation rule for

126    assessing abnormally high mutations has been improved. The expected number of

127    mutations for each sequence is now calculated based on its sampling date and

128    empirical mutation rate [17], providing a more realistic assessment. If the observed

129    number of mutations exceeds the expected number, the genome sequence is

130    highlighted with a red dot, indicating the need for further investigation into

131    sequencing quality issues.

132    With the automated integration pipeline and optimized curation rules, RCoV19

133    accommodated a total of 16,119,080 non-redundant genome sequences from 193

134    countries/regions as of June 10, 2023. A comprehensive and up-to-date list of all

135    released SARS-CoV-2 genome metadata can be freely accessed and downloaded by

136    users at https://bigd.big.ac.cn/ncov/release_genome. The majority of these genomes

137    are contributed by countries such as the United States (31.6%), United Kingdom

138    (19.3%), Germany (5.9%), France (4.4%), Denmark (4.0%), Japan (3.8%), and

139    Canada (3.4%). Among the released human-derived genome sequences (16,103,219),

140    87.7% are complete, and 47.0% are both complete and high-quality. Additionally,

141    RCoV19 offers the service of collapsing identical sequences, resulting in a total of

142    5,832,804 unique sequences (1:1.3) among the complete and high-quality

143    human-derived genome sequences, and 13,762,271  unique sequences (1:1.2) among

144    all released genomes, highlighting the rapid evolution and high diversity of

145    SARS-CoV-2 genomes.

146          To facilitate fast and customized retrieval of SARS-CoV-2 genomes from this

147    vast collection, RCoV19 has developed an advanced search module at

148    https://ngdc.cncb.ac.cn/ncov/genome/search. Users can query by accession ID, Pango

149    lineage, WHO variant label, country/region, host, nucleotide completeness, quality

150    assessment, database resource, sampling date, and sequence length range. The search

151    results are complemented by statistics displayed on the right side of the search page,

152    showcasing distributions in nucleotide completeness, sequence quality, data source,

153    WHO variant label, lineage, country/region, and host. Furthermore, all filtered results

154    can be easily downloaded to support downstream analysis.

155    **Timely monitoring of sequence growth and lineage evolution**

156    With the rapid accumulation of SARS-CoV-2 genome sequences, the emergence of

157    new lineages in specific regions or the whole world has become increasingly

158    prevalent. To enhance our understanding of SARS-CoV-2 evolution and transmission

159    characteristics, we have developed specific modules for monitoring global and

160    regional sequence growth and lineage evolution.

161          Sequence growth serves as an indicator of a country's monitoring capability and

162    level. By examining the cumulative curve of genome sequence growth based on

163    release dates, we can identify three distinct periods: slow growth (January 2020 to

164    March 2021), fast growth (April 2021 to April 2022), and relatively slow growth

165    (May 2022 to present) (**Figure 3A**). We dynamically display the sequence numbers

166    for the top ten countries each month to visualize their contributions (**Figure 3B**).

167    Moreover, we organize sequence numbers for each country/region in a tabular format

168    to provide various detailed data (**Figure 3C**). For example, as of June 10, 2023, a total

169    of 67,149 sequences have been released for China (include Taiwan, HongKong and

170    Maco), with an average release rate of dozens of sequences per day in May 2023.

171          As SARS-CoV-2 spreads, mutations constantly occur and accumulate, leading to

172    the emergence of new lineages and variants. To monitor mutation rates, we calculate

173    the mutation frequency (mutation numbers / genome length) for each genome and plot

174    the daily median mutation frequency as a curve (**Figure 4A**). By observing the slope

175    of curve growth, it is facilitated to timely monitor signals indicating accelerated

176    mutation. For instance, the median mutation frequency rapidly increased to 2.1‰ in

177    mid-December 2021 due to the rapid spread of Omicron variant and reached 3.28‰ in

178    March 2023 due to the spread of XBB.1.5 variant. As sequences with similar mutation

179    spectra are always classified into a Pango lineage [18] or named as a WHO-defined

180    variant

181    (https://www.who.int/news/item/31-05-2021-who-announces-simple-easy-to-say-label

182    s-for-sars-cov-2-variants-of-interest-and-concern), we display the weekly sequence

183    proportion for each lineage or variant. To highlight the main lineages or variants that

184    are currently or previously popular, we interactively display only the top three Pango

185    lineages or WHO-defined variants (**Figure 4B**). Additionally, the sequence proportion

186    for each lineage is represented in a heat map (**Figure 4C**), providing informative

187    insights into lineage trends over time. Taking China as an example, it experienced a

188    wave of COVID-19 infections from late 2022 to early 2023. We developed an

189    interactive map panel to dynamically display the sequence proportions of different

190    lineages and monitor the prevalence and transmission of SARS-CoV-2 at the

191    provincial level (**Figure 4D**).

192    **Pre-warning of potential high-risk haplotypes and lineages**

193    Early and accurate detection of potential high-risk SARS-CoV-2 haplotypes or

194    lineages is a shared challenge for the scientific community in combating the virus.

195    Leveraging the vast amount of genome sequences, we have developed a machine

196    learning model called HiRiskPredictor[13] to predict potential high-risk haplotypes

197    and update these predictions weekly in RCoV19. For each haplotype, a risk score

198    ranging from 0 to 1 is calculated based on the available sequences at that time.

199    Haplotypes with higher risk scores (> 0.5) are identified as potential high-risk

200    haplotypes. A tabular table (**Figure 5A**) organizes the risk score, associated lineage,

201    and transmission-related values (e.g., geographic entropy, betweenness, etc.) for each

202    high-risk haplotype. Users can quickly search for specific haplotypes or lineages

203    using different keywords, or sort the table by 'Risk score' to identify haplotypes with

204    the highest risk scores. Additionally, a boxplot displays the top lineages (20 at most),

205    ranked in descending order based on the median risk scores of all associated

206    haplotypes. **Figure 5B** illustrates the prediction of 12 potential high-risk lineages as

207    of May 31, 2023, with BN.1.2.3, XBB.1.5.24, XBB.1.9.1, XBB.1.16.1, and

208    XBB.1.9.2 identified as the top five lineages. Importantly, the weekly predicted risk

209    scores for all lineages are recorded, allowing users to track historical predictions to

210    detect new warning lineages and understand their development trends (**Figure 5C**).

211    Furthermore, the lineage prevalence, represented by the sequence proportion, is

212    plotted to visualize global changes in epidemic variants (**Figure 5D**). For example,

213    the dominant lineage XBB.1.5 accounted for 20% of all Omicron lineages but is

214    gradually diminishing and being replaced by XBB.1.9.1.

215         These tools and visualizations provided by RCoV19 empower users to identify

216    potential high-risk haplotypes and track the prevalence and evolution of lineages,

217    contributing to early warning systems and informed decision-making in the fight

218    against SARS-CoV-2.

219    **Mutation spectrum comparison between selected lineages or sequences**

220    To facilitate the analysis of mutation spectra and comparisons between different

221    lineages and sequences of SARS-CoV-2, we have developed two interactive modules

222    within RCoV19. These modules allow users to explore mutation distributions and

223    construct mutation maps on lineage level or sequence level.

224         In the inter-lineage or variants comparison module, users can examine the

225    mutation patterns across WHO defined variants (e.g. Delta and Omicron) or Pango

226    lineages (e.g. B.1.177 and XBB.1.5) and analyze mutations by genes or mutation

227    frequency. For example, considering the top three prevalent lineages in the tenth week

228 of 2023 (XBB.1, XBB.1.5, BQ.1) and previous VOCs (Alpha, Beta, Gamma, Delta,

229 Omicron), it is evident that these lineages exhibit more mutations in the *S* gene

230 (**Figure 6A**). Moreover, several novel mutations with high frequencies, such as S371F,

231 T376A, and S477N (frequency > 0.89), have emerged in XBB.1 and XBB.1.5.

232 Additionally, well-known mutations like D614G, known to enhance SARS-CoV-2

233 infectivity in human lung cells, and N501Y, associated with reduced vaccine

234 protection in Delta, may explain the prevalence of ongoing XBB variants [19]. In

235 addition to the extensively studied S protein, N protein mutations like R203K and

236 G204R, implicated in increased transmission [20], are commonly observed in the top

237 three ongoing lineages (**Figure 6B**). Notably, the N protein mutation P13L, which

238 occurs at a high frequency of 90% in the top three ongoing variants, can significantly

239 impair the CD8 T cell epitope (QRNAPRITF), leading to a loss of T cell recognition

240 [21-23]. Similarly, amino acid deletions from position 31 to 33 in the N protein, with

241 a high frequency of 90% among ongoing lineages, may contribute to improved

242 replication efficacy or breakthrough infections, warranting further investigation in the

243 future.

244    In the multiple sequence comparison module, users can sensitively detect

245 potential new lineages by comparing newly released sequences with the representative

246 sequences of the latest lineages in our database. By inputting accession IDs and

247 selecting the lineages of interest, the module displays a mutation matrix for

248 comparison, which can be further refined interactively by genes or differential

249 mutation sites. Additionally, the mutation matrix can be color-coded based on lineage,

250 sampling date, or location. This module is particularly useful in narrowing down the

251 breakpoint range in recombinant variants without the need for intensive sequence

252 similarity calculations. For example, when analyzing the XBB recombinant lineage,

253 comparing it with its parental sequences (BJ.1: EPI_ISL_14891585; BM.1.1.1:

254 EPI_ISL_14733830) reveals that the breakpoint likely lies between V445 and N460 in

255 *S* gene since XBB harbors V445 from BJ.1 and N460 from BM.1.1.1 (**Figure 6C**).

256 Overall, this module complements existing platforms [24, 25] and aids in assessing

257    the validity of newly assigned lineages.

258    These interactive modules within RCoV19 empower users to explore and

259    compare mutation spectra across different lineages and variants, providing valuable

260    insights into the evolution and characteristics of SARS-CoV-2 lineages.

**261    Investigation of the mutation effects on transmissibility and immune escape**

262    A number of mutations have been confirmed to affect viral characteristics, including

263    pathogenicity, infectivity, transmissibility, and antigenicity [8-10, 26, 27]. However,

264    these knowledges are scattered across publications and always focuses on one aspect

265    of a mutation or a variant. To facilitate the effective retrieval of mutation function, we

266    have constructed an integrated knowledgebase by curating information from

267    literatures and databases. Specifically, mutation knowledges are recorded and

268    organized according to their impacts on infectivity/transmissibility, and effectiveness

269    to antibodies, drug, and T cell epitopes.

270    Mutation-related information is collected and categorized based on their specific

271    impacts. For each mutation, we have gathered details on its effects, including a

272    comprehensive description, experimental methods used for characterization, and

273    corresponding PubMed IDs (PMIDs) for reference. In the case of T cell epitope

274    mutations, information on epitopes, HLA restriction, and corresponding T cell types

275    has also been integrated. Overall, we have collected and summarized a total of 2696

276    single mutations, as well as other mutations such as SNPs and Indels, along with 19

277    combined mutations. Among these mutations, 76 affect infectivity/transmissibility,

278    131 are associated with drug resistance, 734 are related to antibody resistance, and

279    1817 mutations are located in T cell epitopes (**Figure 7A**). When considering the

280    distribution of mutations across genes and open reading frames (ORFs), there is an

281    uneven distribution. Specifically, in the S protein, 73 mutations (4%) have been

282    reported to affect infectivity/transmissibility, while 733 mutations (58%) are

283    associated with antibody resistance. This is understandable as the receptor-binding

284    domain (RBD) of the S protein is responsible for virus binding to the ACE2 receptor

10

285 and is a target for neutralizing antibodies. In the *ORF1ab* gene, 127 mutations (49%)

286 are related to drug resistance, which may be attributed to *ORF1ab* being the target of

287 most small molecule inhibitors.

288 Mutations located in T cell epitopes are of particular concern as they are

289 dispersed across different proteins, posing challenges for the immune system to

290 recognize and mount an effective response against various variants. Moreover,

291 mutations in CD4 and CD8 T cell epitopes have the potential to disrupt HLA-peptide

292 binding, leading to immune escape. The diverse epitopes found in different proteins

293 exhibit distinct mutation patterns, which need to be carefully considered during the

294 design of epitope-based vaccines (**Figure 7B**).

295 By providing a comprehensive and organized knowledgebase, researchers and

296 users can easily access and retrieve information regarding the functional impacts of

297 specific mutations. This integrated resource

298 (https://ngdc.cncb.ac.cn/ncov/knowledge/mutation) enhances our understanding of the

299 effects of mutations on viral characteristics and assists in the development of effective

300 countermeasures against SARS-CoV-2 variants.

## Discussion

302 RCoV19 has been continuously updated and developed to support precise prevention

303 of COVID-19. As an integrated repository for SARS-CoV-2 genome data, we have

304 addressed various challenges by implementing a one-stop curation pipeline. This

305 pipeline resolves issues such as sequence redundancy across different repositories,

306 cross-linking between resources, and sequence quality evaluation. However, due to

307 the lack of comprehensive clinical phenotype data, conducting in-depth association

308 studies between massive genomic data and clinical outcomes, as well as unraveling

309 the clinical significance of mutations, remains challenging. To enhance our

310 understanding of disease spread and pathogenesis, we urge the collection and

311 integration of clinical phenotype data of infected individuals to create a more

312 comprehensive platform.

313    Timely monitoring and precise pre-warning based on genomic data are crucial

314    for epidemic prevention. While there are platforms [4-6] available for spatiotemporal

315    surveillance of mutations and variant evolution, there is a deficiency in platforms

316    specifically focused on pre-warning of high-risk variants. In recent years, various

317    machine learning-based prediction models have been proposed, such as PyR0 and

318    VarEPS. PyR0, a hierarchical Bayesian multinomial logistic regression model, can

319    identify mutations that are likely to increase SARS-CoV-2 fitness [28], while VarEPS

320    evaluates the risk level of mutations and variants based on their transmissibility and

321    affinity to neutralizing antibodies using a random forest model [7]. In RCoV19, we

322    have developed a LightGBM model called HiRiskPredictor [13], which calculates a

323    comprehensive risk score and predicts potential high-risk haplotypes on a weekly

324    basis. In the future, we aim to provide multidimensional pre-warning by combining

325    the strengths of different AI models and features.

326    Genetic mutation spectra play a critical role in determining the virological

327    characteristics of different virus strains. Sequence comparison remains the primary

328    approach for identifying differences in mutation spectra. Although the Pango dynamic

329    phylogeny-informed nomenclature system has made significant contributions to

330    tracking genetic diversity and classifying SARS-CoV-2 lineages, there is often a time

331    gap before sporadic variants occurring in specific regions are designated as new

332    lineages. To stay updated on SARS-CoV-2 mutations more sensitively and identify

333    novel lineages earlier, RCoV19 now supports the comparison of newly released

334    sequences with representative sequences of the latest lineages. This feature

335    complements existing public platforms [24, 25] and assists in verifying the assigned

336    lineages of newly released sequences.

337    Numerous mutations have been identified that can increase the severity of

338    infections, enhance transmissibility, and enable evasion of natural and

339    vaccine-induced immunity [29]. Through comprehensive literature curation, we have

340    consolidated a wealth of knowledge regarding the effects of mutations on viral

341    infectivity, resistance to antibodies and therapeutic drugs, and alterations to T cell

342  epitopes. However, further investigation is needed on mutations that impact disease

343  severity. For example, the mutation S194L in the nucleocapsid (N) protein has a

344  notably high frequency among individuals with severe clinical manifestations [30],

345  suggesting its potential contribution to disease progression. Additionally, most of the

346  knowledge on mutation effects is curated from published literature or databases.

347  Future improvements could focus on structural bioinformatics-based prediction of

348  mutation effects, which would enhance our understanding of future pandemics and aid

349  in the development of preventive measures and treatment strategies. In conclusion,

350  knowledge of mutation effects is essential for effective public health interventions, the

351  development of therapeutics, and the creation of pre-warning models.

## Methods

352

### Pre-warning of potential high-risk haplotypes

353

354  All the complete and high-quality SARS-CoV-2 sequences and metadata in RCoV19

355  were used to predict potential high-risk haplotypes weekly. First, we calculated the

356  population mutation frequency (PMF) for each mutated site within every month. Then,

357  those non-UTR mutations with PMF > 0.005 were selected for haplotype network

358  construction by McAN with default parameters. Next, the result of the haplotype

359  network was loaded into HiRiskPredictor with a pre-trained machine learning

360  algorithm (LightGBM) to perform the forewarning analysis process. The

361  HiRiskPredictor automatically extracts features, such as out degree, geographic

362  information entropy, betweenness, etc., for each haplotype in the network. And

363  HiRiskPredictor infers a risk score indicating the likelihood of a haplotype being

364  positive or classified as high-risk according to those features via the pretrained model.

365  If the predicted risk score of a haplotype is greater than 0.5, it is defined as a high-risk

366  haplotype.

### Mutation spectrum comparison between selected lineages or sequences

367

368  Only complete and high-quality genome sequences that have been previously

369 evaluated were employed for the following sequence comparison. To achieve this, the

370 genome sequences were aligned using MUSCLE (version 3.8.31) [15034147] and

371 compared against the initial SARS-CoV-2 genome release (GenBank: MN908947.3).

372 The identification of sequence variations was accomplished using a custom Perl

373 program. At lineage level, mutations among all complete and high-quality sequences

374 of selected variants are displayed in heatmap with customized population mutation

375 frequency. At sequence level, newly emerged sequences with fixed mutations in

376 specific lineage are chosen as representative sequences. After compared with

377 reference genome, mutations among different sequences are displayed in heatmap.

378 Instead of representative sequences, this module also supports to conduct sequence

379 comparison according to Input Sequence Accession within the database.

380 **Investigation of the mutation effects on transmissibility and immune escape**

381 Through a comprehensive literature curation, we have collected a curated list of

382 epitopes that have been experimentally validated. These experiments involved

383 interferon-γ (IFN-γ) enzyme-linked immunospot (ELISpot) assays, complex class I

384 (pMHCI) tetramer staining, and peptide-stimulated activation-induced marker (AIM)

385 assays .etc. Subsequently, we employed an in-house program to integrate all available

386 mutation data across the genome with those effective epitopes and filter mutations

387 with sequences account lower than 2000. Following this, we have conduct a more

388 precise literature curation to search for mutation effect occurring on epitopes to

389 illustrate their functions in T cell recognitions.

## Data availability

391 SARS-CoV-2 genomes, mutations in vcf/tab format and their annotations are publicly

392 available at https://ngdc.cncb.ac.cn/ncov/.

## CRediT author statement

394 **Cuiping Li**: Methodology, Formal analysis, and Writing - Original Draft. **Lina Ma**:

395 Data curation, Methodology, and Writing - Original Draft. **Dong Zou:** Software.

396  **Rongqin Zhang**: Data curation, Methodology, and Writing - Original Draft. **Xue Bai**:

397  Data curation, Writing. **Lun Li**: Methodology. **Gangao Wu**, **Tianhao Wu**, **Wei Zhao**

398  and **Enhui Jin**: Data curation. **Yiming Bao**: Conceptualization, Supervision, and

399  Writing - Review & editing. **Shuhui Song**: Conceptualization, Methodology, and

400  Writing - Review & editing. All authors have read and approved the final manuscript.

## Competing interests

402  The authors have declared no competing interests.

## Acknowledgments

## Authors' ORCID IDs

418  0000-0002-7144-7745 (Cuiping Li)

419  0000-0001-6390-6289 (Lina Ma)

420  0000-0002-7169-4965 (Dong Zou)

421  0009-0000-1570-3292 (Rongqin Zhang)

422  0000-0002-0085-5944 (Xue Bai)

423   0000-0003-3242-031X (Lun Li)

424   0000-0002-3036-5997 (Gangao Wu)

425   0009-0009-9017-7267 (Tianhao Huang)

426   0009-0009-2478-128X (Wei Zhao)

427   0009-0000-9916-9508 (Enhui Jin)

428   0000-0002-9922-9723 (Yimin Bao)

429   0000-0003-2409-8770 (Shuhui Song)

## References

431   [1] Shu Y, McCauley J. GISAID: Global initiative on sharing all influenza data - from

432   vision to reality. Euro Surveill 2017;22:30494.

433   [2] Benson DA, Cavanaugh M, Clark K, Karsch-Mizrachi I, Lipman DJ, Ostell J, et al.

434   GenBank. Nucleic Acids Res 2013;41:D36–42.

435   [3] Brister JR, Ako-Adjei D, Bao Y, Blinkova O. NCBI viral genomes resource.

436   Nucleic Acids Res 2015;43:D571–7.

437   [4] Chen AT, Altschuler K, Zhan SH, Chan YA, Deverman BE. COVID-19 CG enables

438   SARS-CoV-2 mutation and lineage tracking by locations and dates of interest. elife

439   2021;10:e63409.

440   [5] Gangavarapu K, Latif AA, Mullen JL, Alkuzweny M, Hufbauer E, Tsueng G, et al.

441   Outbreak.info genomic reports: scalable and dynamic surveillance of SARS-CoV-2

442   variants and mutations. Res Sq 2022:rs.3.rs-1723829.

443   [6] Chen C, Nadeau S, Yared M, Voinov P, Xie N, Roemer C, et al. CoV-Spectrum:

444   analysis of globally shared SARS-CoV-2 data to identify and characterize new variants.

445   Bioinformatics 2022;38:1735–7.

446   [7] Sun Q, Shu C, Shi W, Luo Y, Fan G, Nie J, et al. VarEPS: an evaluation and

447   prewarning system of known and virtual variations of SARS-CoV-2 genomes. Nucleic

448   Acids Res 2022;50:D888–97.

449   [8] Tzou PL, Tao K, Pond SLK, Shafer RW. Coronavirus Resistance Database

450   (CoV-RDB): SARS-CoV-2 susceptibility to monoclonal antibodies, convalescent

451   plasma, and plasma from vaccinated persons. PLoS One 2022;17:e0261045.

452   [9] Tzou PL, Tao K, Sahoo MK, Kosakovsky Pond SL, Pinsky BA, Shafer RW. Sierra

453   SARS-CoV-2 sequence and antiviral resistance analysis program. J Clin Virol

454   2022;157:105323.

455 [10] Wright DW, Harvey WT, Hughes J, Cox M, Peacock TP, Colquhoun R, et al.
456 Tracking SARS-CoV-2 mutations and variants through the COG-UK-Mutation
457 Explorer. Virus Evol 2022;8:veac023.

458 [11] Song S, Ma L, Zou D, Tian D, Li C, Zhu J, et al. The global landscape of
459 SARS-CoV-2 genomes, variants, and haplotypes in 2019nCoVR. Genomics
460 Proteomics Bioinformatics 2020;18:749−59.

461 [12] Zhao WM, Song SH, Chen ML, Zou D, Ma LN, Ma YK, et al. The 2019 novel
462 coronavirus resource. Yi Chuan 2020;42:212−21(in Chinese with an English abstract).

463 [13] Li L, Li C, Li N, Zou D, Zhao W, Xue Y, et al. Machine learning detection of
464 SARS-CoV-2 high-risk variants. biorxiv 2023;
465 https://doi.org/10.1101/2023.04.19.537460.

466 [14] Li L, Xu B, Tian D, Wang A, Zhu J, Li C, et al. McAN: a novel computational
467 algorithm and platform for constructing and visualizing haplotype networks. Brief
468 Bioinform 2023;24:bbad174.

469 [15] Chen FZ, You LJ, Yang F, Wang LN, Guo XQ, Gao F, et al. CNGBdb: China
470 National GeneBank DataBase. Yi Chuan 2020;42:799−809.

471 [16] Shi W, Qi H, Sun Q, Fan G, Liu S, Wang J, et al. gcMeta: a global catalogue of
472 metagenomics platform to support the archiving, standardization and analysis of
473 microbiome data. Nucleic Acids Res 2019;47:D637−48.

474 [17] Liu Q, Zhao S, Shi CM, Song S, Zhu S, Su Y, et al. Population genetics of
475 SARS-CoV-2: disentangling effects of sampling bias and infection clusters. Genomics
476 Proteomics Bioinformatics 2020;18:640−7.

477 [18] Rambaut A, Holmes EC, O'Toole A, Hill V, McCrone JT, Ruis C, et al. A dynamic
478 nomenclature proposal for SARS-CoV-2 lineages to assist genomic epidemiology. Nat
479 Microbiol 2020;5:1403−7.

480 [19] Fibke CD, Joffres Y, Tyson JR, Colijn C, Janjua NZ, Fjell C, et al. Spike mutation
481 profiles associated with SARS-CoV-2 breakthrough infections in Delta emerging and
482 predominant time periods in British Columbia, Canada. Front Public Health
483 2022;10:915363.

484 [20] Wu H, Xing N, Meng K, Fu B, Xue W, Dong P, et al. Nucleocapsid mutations
485 R203K/G204R increase the infectivity, fitness, and virulence of SARS-CoV-2. Cell
486 Host Microbe 2021;29:1788−801.

487 [21] Peng Y, Mentzer AJ, Liu G, Yao X, Yin Z, Dong D, et al. Broad and strong
488 memory CD4(+) and CD8(+) T cells induced by SARS-CoV-2 in UK convalescent
489 individuals following COVID-19. Nat Immunol 2020;21:1336–45.

490 [22] Nelde A, Bilich T, Heitmann JS, Maringer Y, Salih HR, Roerden M, et al.
491 SARS-CoV-2-derived peptides define heterologous and COVID-19-induced T cell
492 recognition. Nat Immunol 2021;22:74–85.

493 [23] de Silva TI, Liu G, Lindsey BB, Dong D, Moore SC, Hsu NS, et al. The impact of
494 viral mutations on recognition by SARS-CoV-2 specific T cells. iScience
495 2021;24:103353.

496 [24] Hadfield J, Megill C, Bell SM, Huddleston J, Potter B, Callender C, et al.
497 Nextstrain: real-time tracking of pathogen evolution. Bioinformatics 2018;34:4121-3.

498 [25] O'Toole A, Hill V, Pybus OG, Watts A, Bogoch, II, Khan K, et al. Tracking the
499 international spread of SARS-CoV-2 lineages B.1.1.7 and B.1.351/501Y-V2 with
500 grinch. Wellcome Open Res 2021;6:121.

501 [26] Peng Q, Zhou R, Liu N, Wang H, Xu H, Zhao M, et al. Naturally occurring spike
502 mutations influence the infectivity and immunogenicity of SARS-CoV-2. Cell Mol
503 Immunol 2022;19:1302–10.

504 [27] Tian D, Sun Y, Zhou J, Ye Q. The global epidemic of SARS-CoV-2 variants and
505 their mutational immune escape. J Med Virol 2022;94:847–57.

506 [28] Obermeyer F, Jankowiak M, Barkas N, Schaffner SF, Pyle JD, Yurkovetskiy L, et
507 al. Analysis of 6.4 million SARS-CoV-2 genomes identifies mutations associated with
508 fitness. Science 2022;376:1327–32.

509 [29] Thakur S, Sasi S, Pillai SG, Nag A, Shukla D, Singhal R, et al. SARS-CoV-2
510 mutations and their impact on diagnostics, therapeutics and vaccines. Front Med
511 (Lausanne) 2022;9:815389.

512 [30] Dao TL, Hoang VT, Colson P, Lagier JC, Million M, Raoult D, et al.
513 SARS-CoV-2 infectivity and severity of COVID-19 according to SARS-CoV-2
514 variants: current evidence. J Clin Med 2021;10:2635.

515

## Figure legends

517 **Figure 1   Logical architecture diagram of RCoV19 database**

18

518 **Figure 2    Framework of genome data curation model for SARS-CoV-2**

519 RCoV19 integrates genome data from different repositories and provides

520 valued-added curations. It collects metadata and genome sequences from different

521 resources, standardizes metadata, and performs de-redundancy processing based on

522 metadata and sequence comparison. These steps have been chained together as one

523 workflow, which is activated automatically every day and sends the integration

524 statistics to mobile phone client at the end. After integration, RCoV19 performs a

525 series of assessments; it determines completeness of the protein-coding region,

526 assesses sequence quality in five aspects, and defines high-quality sequences. We

527 consider a sequence to be of high quality if it could pass quality control for both Ns ($<=$

528 15) and degenerate bases ($<=$ 50). Otherwise, it is of low quality.

529 **Figure 3    The monitoring platform of SARS-CoV-2 sequence growth globally**

530 **and regionally**

531 **A**. The dynamic growth curve of globally and China released genome sequences, and

532 globally released complete genome sequences. **B**. A bar chart shows the top ten

533 countries with the most public released sequences as of June 3, 2023. **C**. A tabular

534 table shows the statistic of sequences in country/region.

535 **Figure 4    Monitoring of SARS-CoV-2 lineage evolution globally and regionally**

536 **A**. Number of released sequences and the average mutation frequency along sequence

537 sampling date. The mutation frequency is calculated by dividing the total mutation of

538 each sequence by the genome length. **B**. The stacking diagram shows the proportion

539 of top three prevalent Pango lineages or WHO define abbrevaitions used variants per

540 week. **C**. Heatmap of the frequency of the cumulative sequences for selected lineage

541 in China. **D**. Geographical distribution of the sequences number in China, and the pia

542 chart shows the lineage proportions in provinces from May 1st to June 3 in 2023.

543 **Figure 5    Pre-warning of potential high-risk haplotypes and lineages**

544 **A**. A screenshot of the tabular table for all haplotypes with values of haplotype

545 network features and its risk score. **B**. Boxplot of predicated risk score for all

546 haplotypes of the top twenty lineages. As of May 31, 2023, 12 lineages have been

547  predicted as potential high-risk lineages. **C**. Distribution of the historical risk scores

548  for user selected lineages. **D**. Genomic prevalence of lineages based on sequence

549  collection date.

550  **Figure 6    Mutation spectrum comparison among selected lineages and**

551  **sequences**

552  **A**. Lineage mutation comparison on *S* gene among top 3 prevalence lineages in 10$^{th}$

553  week of 2023 (XBB.1, XBB.1.5, BQ.1) and previous VOC defined by WHO (Alpha,

554  Beta, Gamma, Delta, Omicron) with mutation frequency. **B**. Lineage mutation

555  comparison on *N* gene among top 3 prevalence lineages in 10$^{th}$ week of 2023 (XBB.1,

556  XBB.1.5, BQ.1) and previous VOC defined by WHO (Alpha, Beta, Gamma, Delta,

557  Omicron) with mutation frequency. **C**. Sequence mutation comparison among

558  sequences (XBB: EPI_ISL_15854782, BJ.1: EPI_ISL_14891585; BM.1.1.1:

559  EPI_ISL_14733830) presented by differential mutations (refers to those after

560  removing common mutations among sequences) in each sequence, the range between

561  mutations in red color indicating possible recombination breakpoint.

562  **Figure 7    Mutation effects on SARS-CoV-2 viral characteristics**

563  **A**. Collection of mutation effect knowledge. The horizontal axis represents the

564  number of mutation types. **B**. Mutations occurring on experimentally verified T cell

565  epitopes. The magnitude of the circles represents the number of mutations on each

566  epitope, and different colors indicate T cell epitopes on different proteins.

567  **Supplementary material**

568  **Figure S1 SARS-CoV-2 genome sequence overlaps among different sources (as of**

569  **August 16, 2023)**

Effective knowledge

Mutation comparison

RCoV19

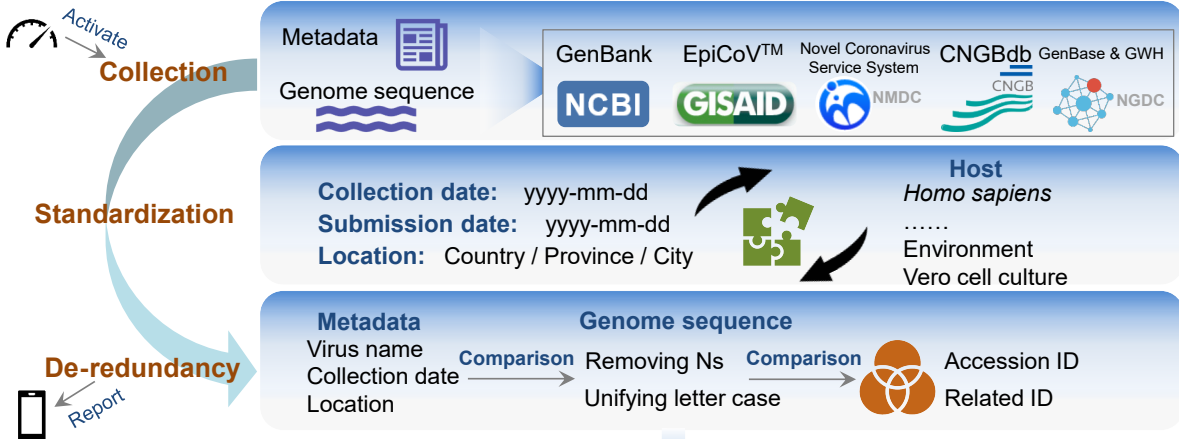https://ngdc.cncb.ac.cn/ncov

Data integration

Monitoring

Pre-warning

| Genome completeness | | N(s) | Degenerate base(s) | Gap(s) | Mutation(s) | Mutation density | Sequence quality |
|---|---|---|---|---|---|---|---|
| Complete | 🟢 | <= 15 | <= 50 | <= 2 | <= Expected+1 | < 0.25 | High |
| Partial | 🔴 | > 15 | > 50 | > 2 | > Expected+1 | >= 0.25 | Low |

**A**



**B**

No. of sequences

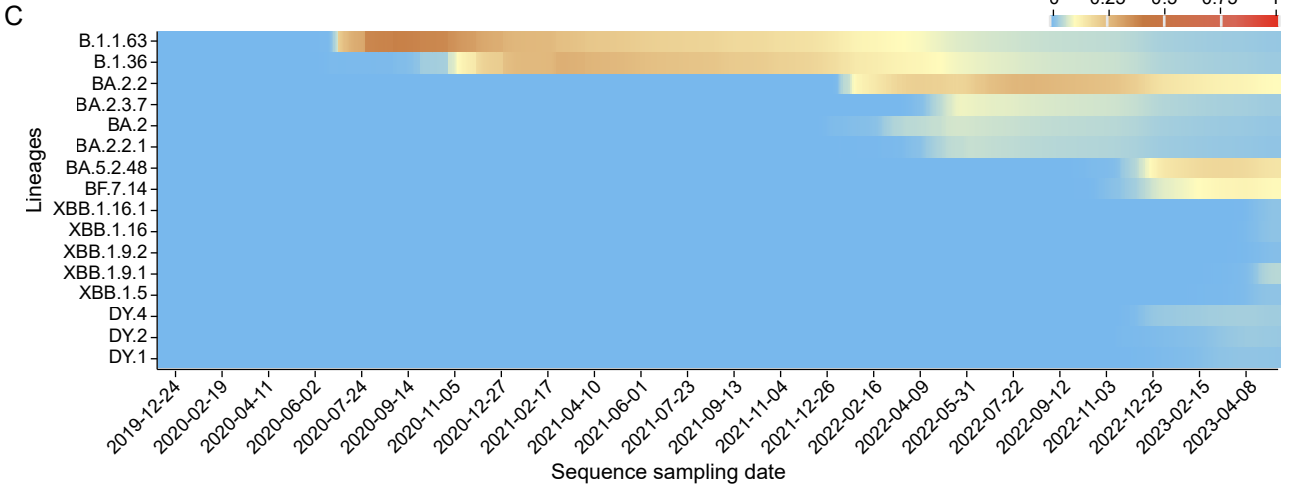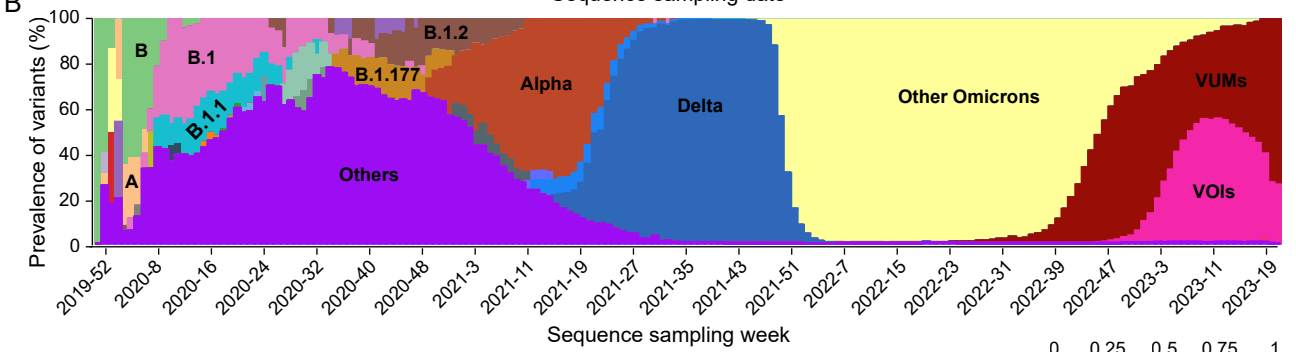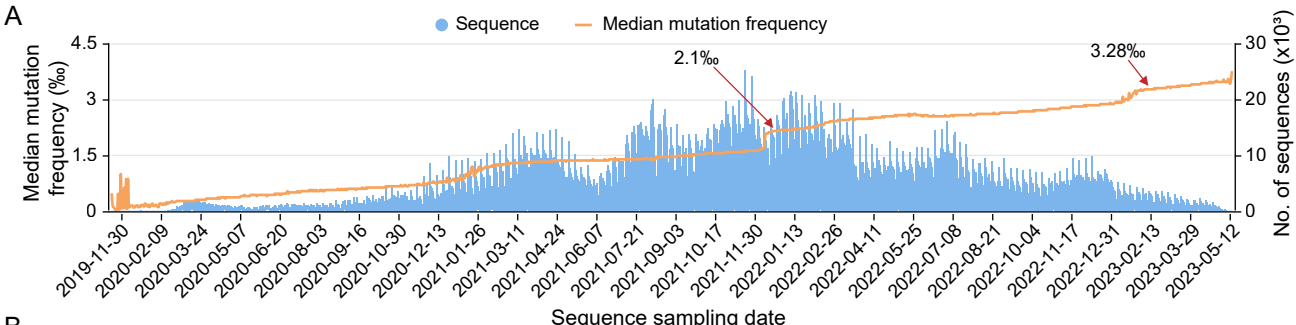| Country | No. of sequences |
|---|---|
| United States | 5,086,799 |
| United Kingdom | 3,110,610 |
| Germany | 948,969 |
| France | 704,249 |
| Denmark | 637,806 |
| Japan | 608,903 |
| Canada | 547,991 |
| India | 292,477 |
| Austria | 270,609 |
| Sweden | 255,247 |

**2023-06**

**C**

Search:

| Continent | Country/Region | Genome Sequences | Complete Genome Sequences | Human Genome Sequences | Human Complete Genome Sequences | Monitoring report |
|---|---|---|---|---|---|---|
| North America | 🇺🇸 United States | 5087879 | 4461690 | 5082702 | 4456569 | 🖥 |
| Europe | 🇬🇧 United Kingdom | 3111073 | 2836424 | 3111064 | 2836420 | 🖥 |
| Europe | 🇩🇪 Germany | 948969 | 868074 | 948905 | 868011 | 🖥 |
| Europe | 🇫🇷 France | 704249 | 573270 | 703992 | 573017 | 🖥 |
| Europe | 🇩🇰 Denmark | 637806 | 610682 | 637335 | 610211 | 🖥 |

Showing 1 to 5 of 193 entries
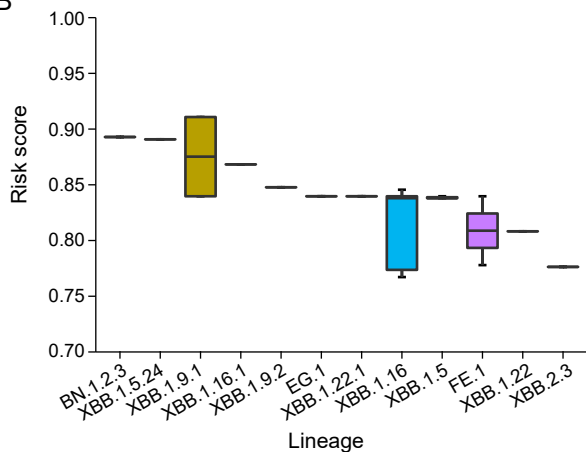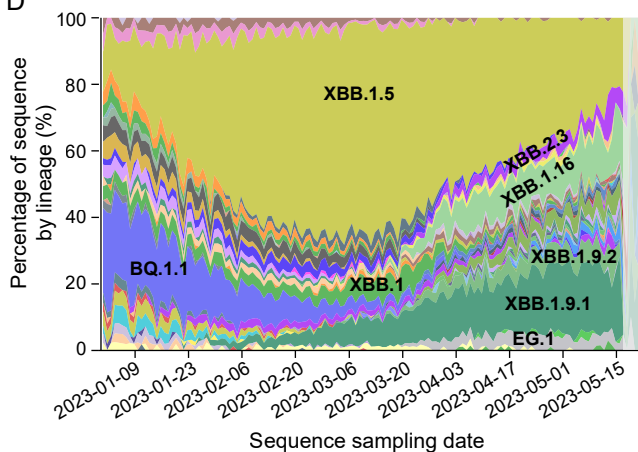
Previous 1 2 3 4 5 Next

A

| Haplotype ID | Lineage | Geographic information entropy | Betweenness | Sequences number of haplotype | Out-degree | Mutation scores | Sequential growth ratio | Connectivity of nodes | Risk score |
|---|---|---|---|---|---|---|---|---|---|
| Node_7430 | XBB.1.9.1 | 0.0980 | 324 | 50 | 26 | 70 | 0.7800 | 2 | 0.9108 |
| Node_11323 | XBB.1.9.1 | 0.1861 | 737 | 88 | 23 | 70 | 0.8182 | 2 | 0.9108 |
| Node_11729 | BN.1.2.3 | 0.5623 | 63 | 4 | 3 | 67 | 1.0000 | 1 | 0.8929 |
| Node_40026 | XBB.1.5.24 | 0.4101 | 112 | 7 | 4 | 72 | 1.0000 | 1 | 0.8908 |
| Node_1293 | XBB.1.16.1 | 0.5004 | 33 | 5 | 3 | 72 | 1.0000 | 1 | 0.8683 |
| Node_3743 | XBB.1.9.2 | 0.6931 | 57 | 2 | 2 | 70 | 1.0000 | 1 | 0.8478 |
| Node_730 | XBB.1.16 | 0.6365 | 13 | 3 | 1 | 71 | 1.0000 | 1 | 0.8456 |
| Node_1452 | XBB.1.9.1 | 0.6931 | 11 | 2 | 1 | 72 | 1.0000 | 1 | 0.8397 |
| Node_5500 | EG.1 | 0.6931 | 22 | 2 | 1 | 72 | 1.0000 | 1 | 0.8397 |
| Node_5507 | XBB.1.16 | 0.6931 | 27 | 2 | 2 | 70 | 1.0000 | 1 | 0.8397 |

A

**T cell epitopes** 1816

**Antibody resistance** 734

**Drug resistance** 131

**Infectivity transmission** 76

Effect

No. of mutations

Legend:
E  N  ORF1ab  ORF6  ORF8
M  S  ORF3a  ORF7a  ORF10

B

**Proteins**

FYVYSRVKNLNSSRV
PINLVRDLPQGFSAL
NVTWFHAIHVSGTNG
ITRFQTLLALHRSYL
YNYLYRLFRKSNLKP
CTFEYVSQPFLMDLE
RFDNPVLPF
GVYFASTEK
YEQYIKWPWYIWLGF
SPRRARSVA
VVLSFELLHAPATVC
FNGLTVLPPLLTDEM
TDEMIAQYTSALLAG
NFSQILPDPSKPSKR
LLFNKVTLA
APHGVVFL
QLIRAAEIRASANLAATK
GVSPTKLNDLCFTNV
YAWNRKRISNCVADY
VLNDILSRL
YLQPRTFLL
GTHWFVTQR
NLLLQYGSFCTQLNR
PFFSNVTWFHAIHVS
FIAGLIAIV
KLPDDFTGCV
RLNEVAKNL
DGVKHVYQLRARSVSPKL
QEEVQELYSPIFLIV
IWNLDYIINLIIKNL
SKWYIRVGARKSAPL
LLLLDRLNQLESKMS
MKDLSPRWYFYYLGTGPEAG
GTWLTYTGAIKLDDK
ASAFFGMSRIGMEVT
FPRGQGVPI
ATEGALNTPK
ASWFTALTQHGKEDL
YKHWPQIAQFAPSAS
KTFPPTEPK
QRNAPRITF
KAYNVTQAFGRRGPE
MEVTPSGTWL
AIVLQLPQGTTLPKG
DDQIGYYRRATRRIR
AADLDDFSKQLQQSM
PNFKDQVILLNKHIDAYK
KDGIIWVATEGALNT
KPRQKRTATKAYNVT
LIRQGTDYKHWPQIA
TWLTYTGAIKLDDKDPNF
IGYYRRATRRIRGGD
INVFAFPFTIYSLLL
FMRIFTIGTVTLKQG
YLYALVYFL
ALSKGVHFV
LLYDANYFL
FTSDYYQLY
VYFLQSINF
TTDPSFLGRY
ASMPTTIAK
STFNVPMEK
SAFAMMFVK
VTNNTFTLK
PTDNYITTY
CTDDNALAYY
RVESSSKLWAQCVQL
KTIQPRVEK
FLLPSLATV
ALWEIQQVV
GTDLEGNFY
FLLNKEMYL
HTTDPSFLGRY
IPRRNVATL
GAVILRGHLRIAGHHLGR
LSYYKLGASQRVAGD
TSRTLSYYKLGASQRVA
LRIAGHHLGRCDIKD
LRGHLRIAGHHLGRC
SELVIGAVIL
LGASQRVAGDSGFAA

S
M
E
ORF1ab
ORF3a
ORF10
N
ORF7a
ORF6
ORF8

No. of mutations
50  100  500  1000

E  N  ORF1ab  ORF6  ORF8
M  S  ORF3a  ORF7a  ORF10