

1 **Factors affecting SARS-CoV-2 sequence conservation:**
2 **SARSNTdb database**

3 **By**

4 **JOHN ORGERA**

5 A thesis submitted to the

6 Graduate School-Camden

7 Rutgers The State University of New Jersey

8 In partial fulfillment of the requirements

9 For the degree of Master of Science

10 Graduate Program in Computational and Integrative Biology

11 Written under the direction of

12 Dr. Andrey Grigoriev

13 And approved by

14 _____
15 Dr. Andrey Grigoriev

16 _____
17 Dr. Sunil Shende

18 _____
19 Dr. Marien Solesio

20 Camden, New Jersey

21 August 2022

22
23
24
25

26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49

ABSTRACT OF THE DISSERTATION

Factors affecting SARS-CoV-2 sequence conservation: SARSNTdb database

By John Orgera

Dissertation Director:

Dr. Andrey Grigoriev

SARSNTdb offers a curated, nucleotide-centric database for users of varying SARS-CoV-2 knowledge. Its user-friendly interface enables querying coding regions and coordinate intervals to find out the various functional and selective constraints that act upon the corresponding nucleotides and amino acids. Users can easily obtain information about viral genes and proteins, functional domains, repeats, secondary structure formation, intragenomic interactions, and mutation prevalence. Currently, many databases are focused on the phylogeny and amino acid substitutions, mainly in the spike protein. While providing mutation data, SARSNTdb takes a more nucleotide-focused approach as RNA does more than just code for proteins and many insights can be gleaned from its study. For example, RNA-targeted drug therapies for SARS-CoV-2 are currently being developed and it is essential to understand the features only visible at that level. This database enables the user to identify regions that are more prone to forming secondary structures that drugs can target. Finally, the database allows for comparing SARS-CoV-2 and SARS-CoV domains and sequences. SARSNTdb can serve the research community by being a curated repository for information that gives a jump start to analysing a mutation's effect far beyond just determining synonymous/non-synonymous substitutions in protein sequences.

50 INTRODUCTION

51 Background

52 Coronavirus disease 2019 (COVID-19) pandemic caused by the Severe Acute
53 Respiratory Syndrome Coronavirus 2 (SARS-CoV-2) virus was first identified in
54 December 2019 in Wuhan China(1). The WHO declared the outbreak a pandemic in
55 March 2020(2) and since then the scientific community has generated a massive
56 amount of data and that needs to be organized in order to understand the virus. This
57 data includes sequences, papers, medical information, and proteomics data. Over
58 time the virus has continually evolved and has branched into several variants, some of
59 them being Variants of Concern (VOC) such as the Omicron and Delta variants(3).
60 These variants are identified primarily by sequencing, further stressing the importance
61 of competent SARS-CoV-2 databases that allow scientists to interrogate mutations
62 understand the functions of the various viral proteins that the variants may affect.

63 The virus is a positive sense RNA virus with 16 non-structural proteins (NSPs) and 10
64 structural and accessory proteins(4). The 16 non-structural proteins are cleaved from
65 two open reading frames (ORF) known as ORF1a and ORF1b(5). ORF1b is a
66 polyprotein that is cleaved into NSP1 to NSP11. ORF1a is cleaved into NSP12-
67 NSP16(5). ORF1a and 1b are separated by a -1 ribosomal frameshift upstream of
68 the ORF1a stop codon. The accessory proteins are transcribed to a negative sense
69 template by the RNA-Dependent RNA Polymerase (RdRP), a complex of nsp12, 7,
70 and 8(5). This is done in the 3'-5' direction from the 3' end of the genome. When the
71 RdRP reaches the 5' start of a protein it is transcribing it encounters a transcription
72 regulatory sequence B (TRS-B). This initiates a jump to the TRS-Leader at
73 coordinate 70 on the 5' end, skipping the nucleotides between the TRS-B and TRS-L

74 (5). This attaches the leader sequence to the 5' end of the template, allowing the
75 host's ribosomes to recognize and translate the RNA.

76 The virus has several unique evolutionary constraints placed upon it. For example, its
77 relatively short genetic material caused it to evolve several unique methods to create
78 novel proteins. As previously stated, the frameshift separating ORF1a and ORF1b is
79 one of them. This frameshift is caused by a unique 7nt "slippery sequence" at the
80 RNA level(5). Another example of short RNA motifs regulating the genome is the
81 TRS-B. This 6nt long sequence is responsible for the creation of all the accessory and
82 structural proteins. Several proteins overlap at the RNA level such as ORF7a and
83 ORF7b or the Nucleocapsid (N) protein and ORF9b(4). ORF7b proteins are thought
84 to be created through leaky scanning of ORF7a(4). It is likely features at the RNA
85 level that are responsible for these efficient methods of translation. Finally, RNA
86 secondary structures are often formed by SARS-CoV-2(6). RNA secondary
87 structures are often targeted by host defences like RIG-I and are placed under
88 evolutionary constraint(7).

89 **Defining a Niche**

90 Several databases related to SARS-CoV-2 have studied and reported on the evolution
91 of the virus and identifying variants(8). For example, GISAID(9), the primary
92 repository of assembled SARS-CoV-2 genomes at the time of writing, has over 11
93 million genome sequence samples submitted and list on their website several
94 databases that use this new data to track the evolution of the virus over time.
95 However, interpreting a reported nucleotide or amino acid substitution often requires
96 sifting through pieces of information related to affected proteins of the virus. Such
97 information is scattered throughout the web in many papers and databases. If a
98 mutation is found at a certain coordinate, a thorough investigation delving into

99 multiple papers is required to understand the functional importance of that one
100 nucleotide and its surroundings. To remedy this, we created SARSNTdb, a compact
101 database of highly interlinked data records that can allow the user to rapidly navigate
102 from genome positions to functional/selective constraints on the corresponding
103 nucleotides and amino acids.

104 Genome databases typically list coordinates of coding and non-coding regions and
105 provide their annotations per such region. In contrast, SARSNTdb is nucleotide-
106 centric, it allows querying annotations for every position in the genome from the
107 perspective of potential selection factors affecting the corresponding nucleotide.

108 Public attention to SARS-CoV-2 virus has generally focused on mutations occurring
109 in its genome (and their impact on vaccine efficacy and virus spread). Most
110 frequently, SARS-CoV-2 mutations are viewed through the prism of immune system
111 evasion (3, 10). While this is relevant for the (most widely known) viral spike protein,
112 the general public and scientific community are often at a loss when other
113 substitutions are considered, especially silent ones or short insertions/deletions
114 (indels). Given the significant interest in variants of concern (VOC), strong focus on
115 selection would provide complementary functional context for respective VOC
116 mutations, beyond the trivial synonymous/non-synonymous designations. Examples
117 of such context include repeats, secondary structure formation, intragenomic
118 interactions, nucleotide and amino acid conservation, and mutation prevalence.

119 For example, it is known that repeats and their variations play a critical role in
120 production of subgenomic mRNAs in coronaviruses (5). These repeats direct the
121 RdRP to jump from one coordinate to the leader sequence at the 5' end of the
122 genome, creating a template recognizable by the host's ribosomes (5). In addition,
123 variations in the RNA declared as synonymous may alter the RNA secondary

124 structure formation. Currently, RNA secondary structure is being investigated for
125 drug targeting and variations on the structure could affect therapeutic effects of drugs
126 currently in development (11). Finally, VOCs, such as Omicron, display large
127 number of both spike (12) and non-spike substitutions or indels. Hence, it is
128 important to recognize the effect of mutations commonly passed over as non-
129 synonymous or taking place in regions not under intense scrutiny.

130 To ensure the consistent cataloguing of the nucleotide and amino acid substitutions,
131 we re-evaluated mutations across >22,000 patient samples, for which raw
132 metatranscriptome datasets of sufficient quality were available in NCBI's SRA(13).
133 We avoided taking mutations reported in GISAID since it contained already
134 assembled genomes. In many cases, such genomes contained unresolved regions or
135 segments of low genome coverage (jeopardizing mutation calling or producing
136 massive sections of missing data) and it was not possible to tell how reliable these
137 assemblies were. By calling mutations *de novo* we increased the consistency of the
138 substitution data collected.

139 **Perspectives**

140 We feel that this database fills a niche yet unfilled by other SARS-CoV-2
141 databases. Other databases often have a focus on the amino acid coordinates and non-
142 synonymous substitutions, passing over other important SNVs at the RNA level. In
143 addition, data was scattered so that upon sequencing a novel mutation, an
144 investigation delving into several papers and databases was required to understand the
145 functional importance of affected nucleotide. Now, with SARSNTDB, one can
146 simply input that coordinate into our *Genome Search* function and find what proteins
147 it is included in and the selective factors that act upon it.

148

149 This database will continue to be useful to researchers or students new to SARS-CoV-
150 2 as most of the data it provides remains relatively constant, even with new variants
151 emerging. The data in the database is subject to change of course as our
152 understanding of the virus continues to evolve. One example of this is data related to
153 SNVs, of which there is a great deal available online already. However, instead of
154 using GSAID's available SNV data we downloaded a large amount of data from
155 NCBI's website and aligned them using GROM. This was done as we found
156 excessive amounts of soft clipping with the commonly used Minimap2. GROM, on
157 the other hand, had less. Overall, we realigned and called variants on >16000 samples
158 and found ~33000 unique SNVs. These data are also available on our database.

159

160 Unique datasets also available on our website include both visualized comparisons of
161 SARS-CoV-2 and SARS-CoV, and repeats of SARS-CoV-2. Both of these data were
162 of interest to our lab and we wanted to make them available to others as well.

163

164 **MATERIAL AND METHODS**

165 **Data Collection and Processing**

166 We downloaded mutation data from the NCBI's SRA using the prefetch feature.
167 Datasets from the UK had been originally aligned using minimap2(14). We found that
168 minimap2 produced in these SARS-CoV-2 datasets excessive soft clipping, leading to
169 lower accuracy when calling variants. To solve this, we unaligned the reads using
170 Samtools (15) to convert them to fastq files. We then re-aligned the fastq files to the
171 reference sequence using BWA mem, producing what we found to be more consistent

172 alignment patterns with less soft clipping. We did the same alignment step for the
 173 samples, where fastq files were available in SRA. We then used GROM (16) to find
 174 SNVs in the data. In total we found ~33,000 unique SNVs using GROM
 175 across >22,000 samples. The output VCF files were then converted to SQL files
 176 detailing the sequencing platform, coordinates, and alternate nucleotides for each
 177 sample.

178 To identify repeats in the SARS-CoV-2 genome, we analysed the Wuhan reference
 179 sequence(1) (NC_045512.2) using UGENE(17) with the default settings. We then
 180 organized the repeats by coordinate and identified repeats that were super-repeats of
 181 one another (superstrings of shorter repeat strings) using in-house scripts.

182

Data Type	Tool Used	Source
Protein Structure Visualizations	I-TASSER – M.L. based protein structure predictor	Zhang group(18)
SHAPE reactivities of RNA	SHAPE-MaP	Yang et al (19)
SHAPE reactivities of RNA	icSHAPE	Sun et al(11)
Normalized SHAPE reactivities of RNA	SHAPE-MaP and DMS-MaPseq	Manfredonia et al(6)
Intragenome RNA interactions	SPLASH	Yang et al(19)
Repeat Detection and Coordinates	UGENE	Scripts ran in-house
SNV Data	GROM	Produced in-house

183 Table 1: A display of the types, sources, and tools used to generate data that populate the database

184

185 **Data on Protein and RNA Structure**

186 Protein structures were obtained from the Zhang group who has used I-TASSER to
187 predict protein structure for all SARS-CoV-2 proteins (18). Their predictions are
188 highly accurate for the SARS-CoV-2 proteins despite relatively few homologous
189 sequences with available protein structures.

190 To show the secondary structure of SARS-CoV-2 genomic RNA we collected several
191 datasets from groups that have measured the viral RNA accessibility at a single base
192 resolution. The first of these was taken from Manfredonia et al. (6) who has used
193 SHAPE and DMS mutational profiling to find secondary structure maps with single
194 base resolution. Yang et al.(19) has used SHAPE-MaP to find the reactivities of the
195 reference sequence as well as a delta variant sequence. Finally, Sun et al (11) has
196 used icSHAPE to map reactivities.

197 Data presented as *Intragenome Interaction Data* represent regions of pairwise RNA
198 interactions across the genome. Such regions have been detected via proximity
199 ligation sequencing was performed using SPLASH to find these regions in Vero-E6
200 infected cell (19).

201 **Gene, Protein and Functional Domain data**

202 We obtained the coordinates of viral non-coding regions, its genes and proteins, their
203 respective nucleotide and amino acid sequences from the NCBI record of the SARS-
204 CoV-2 (NC_045512.2). SARS-CoV's information was retrieved in the same way
205 from the NCBI record of the Tor reference sequence (NC_004718.3). We then
206 performed a thorough literature review (across hundreds of papers) of proteins in
207 SARS-CoV-2 and SARS-CoV to obtain their functional descriptions. Next, we

208 identified the available coordinates of functional domains in both viruses. Using
209 BLAST(20) and CLUSTAL-W(21), we further performed pairwise alignments of the
210 proteins of SARS-CoV-2 and SARS-CoV to evaluate the levels of amino acid identity
211 of the homologous functional domains. We manually curated mismatched coordinates
212 of such homologous domains between different studies, produced reconciled
213 coordinates and transferred the domain annotations, further accompanied on
214 respective pages by the publications describing them.

215 **RESULTS AND DISCUSSION**

216 Succinctly, the data in the database is retrieved by users via two main query hubs.
217 One is the *Genome Search* page and is comprised of several datasets and information
218 retrieved from literature. The other is made available in the *Mutation Search* page
219 (and *Repeat* page), presenting results of our re-analysis of >22,000 patient samples
220 obtained from NCBI's publicly available SRA SARS-CoV-2 genomes. We
221 interlinked these sections comprehensively in order to provide the user an easy way to
222 carry over the findings gained in one section to another.

223 **Web App Implementation and User Interface**

224 Users can access SARSNTdb at <https://grigoriev-lab.camden.rutgers.edu/sarsdb/>

225 The website is implemented in PHP (version 7.4.29) and the SQL server through
226 mysql (version 15.1) with MariaDB (distribution 10.3.34).

227 The interface of the database consists of several tabs. The *Search* tab has a dropdown
228 menu that brings the user to a *Genome Search*, *Mutation Search*, and a *Repeat Search*.
229 These searches are interconnected to allow the user to take information gleaned from
230 one search into another. The *Help* tab instructs the user on how to use the website by

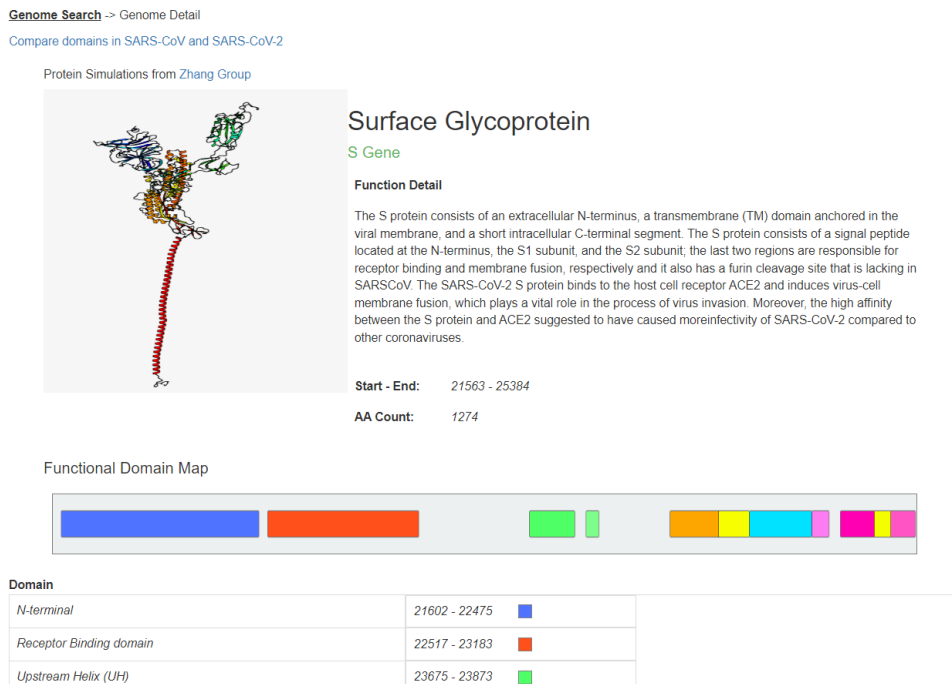
231 providing an example. The *Reference* tab brings the user to this article where they can
232 learn about the data sources and how the website was constructed.

233 **Accessing Gene, Protein and Functional Domain Details**

234 The *Genome Search* page allows the user to specify nucleotide coordinate intervals
235 and find information about functionally relevant regions of the SARS-CoV-2 virus
236 that overlap or are contained between these coordinate pairs. Such regions most often
237 correspond to genes and functional domains they encode. Also, this search reports
238 about nearby repeats and intragenomic interactions obtained using a SPLASH
239 technique (19).

240 One can also select a single ORF or Nsp from a menu to get to such genes. Their
241 protein products are described on the *Protein Detail* page [Fig.1]. In addition to
242 images of the predicted structure of the SARS-CoV-2 proteins, their functional
243 domains, smaller motifs, and certain amino acid residues with annotated functionality
244 are also displayed graphically. At the bottom of this page there are the relevant RNA
245 and Protein sequences derived from the respective NCBI reference.

246 Since SARS-CoV-2 domains are typically derived from the previous body of work on
247 SARS-CoV, we devoted a special page for each protein in both viruses for comparing
248 domains. This page is linked from the *Protein* page and contains a table detailing the
249 similarities of the two viruses and an alignment of both protein sequences created
250 using CLUSTALW(21) and BLAST(20). The coordinates in the table are derived
251 from primary literature and review papers (that can be accessed by clicking the
252 hyperlinks on the coordinates) and sometimes they differ, despite being reasonably
253 well aligned. The residue identities and positives were derived by performing a
254 BLAST alignment of the genome sequences.



262 Fig.1: The Genome Detail page of the S protein

263 Visualization of Mutation and RNA Structure Details

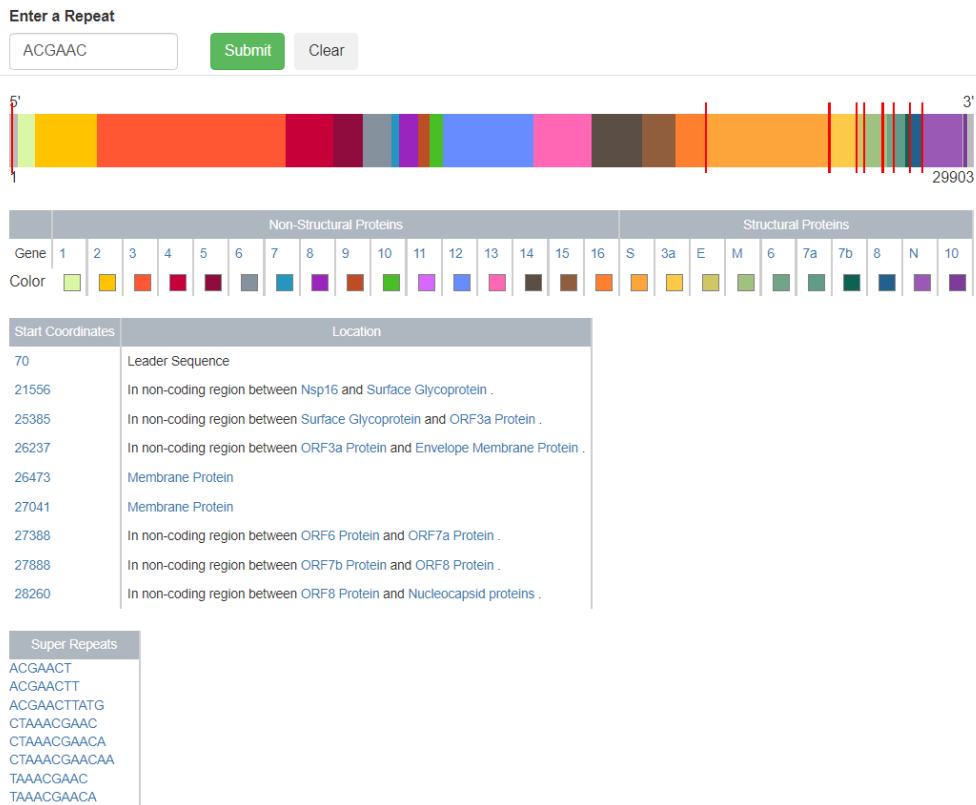
264 The *Mutation Search* page allows the user to search for mutations in a nucleotide
 265 range or within a gene. The search results are bar graphs depicting the number and
 266 type of substitutions in the range. The bars are also subdivided by sequencing
 267 platforms. If a more granular view of the mutations is needed the user can click on
 268 *Mutation Detail* to see expanded information related to the mutations in the that
 269 nucleotide range. We also provided links “Try your luck in Covariant” for each non-
 270 synonymous substitution as some of the frequent mutations (mostly in the S gene) can
 271 be reflected in that database. Below the *Substitution Frequency* table there is a
 272 histogram that displays SNV frequency across the nucleotide range selected.

273 Also on this page is SHAPE data that may help inform the user why certain regions
 274 may be conserved due to the secondary structure constraints. When the size of the
 275 searched region is large, the SHAPE Data is displayed in intervals where the SHAPE
 276 value is averaged across that region. If the size of the searched interval is under 100

277 nucleotides, each position and the SHAPE value is displayed individually. If the
278 shape value is above 0.5 it is displayed in blue indicating a high reactivity while
279 below 0.5 is displayed in red and indicates a low reactivity. We selected several
280 SHAPE datasets and displayed them in separate graphs. These data make up our
281 *Mutation Search* page and it is visualized on the page using CanvasJS (22).

282 **Repeats in the Genome**

283 The *Repeat Page* [Fig.2] allows the user to search the SARS-CoV-2 reference
284 sequence for repeats of size 6 nucleotides or greater. Displayed on this page is the
285 genome schematic with proteins coloured distinctly. When a repeat is found red lines
286 appear on the genome indicating repeat locations, and a table displaying the
287 coordinates of the repeats as well as which protein they appear in is displayed. Also
288 available are repeats, which are super-strings containing the searched repeat; these are
289 deemed super-repeats. For example, the repeat AACAGGA is a super-repeat of
290 AACAGG as the former is a super-string of (i.e., contains) the latter. Clicking on
291 these super-repeats brings the user to a *Repeat Page* for the super-repeat (with their
292 respective super-repeats, if available). For the default search on the *Repeat Page* and
293 a clear biological example, we provide the minimal repeat of the transcription
294 regulatory sequence (TRS) from the SARS-CoV-2 virus(5), with all locations of
295 canonical TRS visualized throughout the genome for the user.



296

297 Figure 2. Visualization of the leader sequence repeat ACGAAAC across the genome.

298 Case Study

299 As stated previously there are many databases tracking the waves of VOCs and their
 300 typical mutations. The virus continues to evolve, and even the general public is made
 301 aware of new substitutions in the best-annotated spike protein. When new mutations
 302 appear, it is important to be able to quickly identify where they occur and analyse
 303 their effects by detecting genome features nearby. Furthermore, substitutions take
 304 place not only the spike protein, yet those affecting other parts of the genome are
 305 typically ignored in the databases and many analyses.

306 In contrast, SARSNTdb could be an excellent starting point for such quick evaluation.
 307 Consider the mutation C28311T, found in Omicron. Let us first go to our *Genome*
 308 *Search* page and input the coordinate 28311 and find it is part of two overlapping
 309 genes, encoding the Nucleocapsid protein as well as ORF9b. In N it is part of the N-

310 terminal arm/Intrinsically disordered region. In ORF9b we see it is part of the site
311 that interacts with NEMO. We also see that it is a part of some common repeats and
312 has intragenomic interactions at the 5' end of the protein as well as a region 200nt
313 away that it binds with. These close intragenomic interacting regions could form
314 pockets that are often the targets of RNA based therapeutics (11). Clicking *View*
315 *Details* for ORF9b, we find its function and see it suppresses the innate immune system
316 through regulating Mitochondrial Antiviral Signalling pathways (7). In comparing it
317 to SARS-CoV we find this domain is not well conserved with only 63% similarity
318 overall. In the paper linked via the domain coordinates in the table, we find that this
319 region, when deleted, resulted in a loss of function of the protein and its interaction
320 with NEMO (7). If this nucleotide change results in a non-synonymous mutation, it
321 could affect the function of the protein. By clicking *Mutations* on the table, we are
322 brought to the mutation page showing the NEMO interaction region's mutation
323 frequencies, SHAPE scores and, if we click *Detail*, a breakdown of that specific
324 mutation has been found. If it has been found the detail page will also show the type
325 of variant it creates. In this case the SNP has been found previously and changes a
326 proline to a serine. We find it has been found in >1500 samples but the link "Try
327 your luck in Covariant" will return Error 404, as this variant is not yet annotated
328 there. In addition, the SHAPE score of this nucleotide is low according to all datasets,
329 indicating that it may be prone to forming secondary structures within the RNA.
330 Overall, with all such results about this mutation we can conclude that it should be
331 monitored as it has been persisting over time and now, with Omicron spreading
332 rapidly, may be gaining increased prevalence. This mutation could affect the ability
333 of ORF9b to suppress the innate immune system through interacting with NEMO and
334 its effects should be explored further.

335 **Conclusion**

336 SARSdb is a database for users of varying levels of knowledge about virology or
337 genomics. It provides nucleotide-level functional information about various aspects of
338 the SARS-CoV-2 genome. It features a quick and easy coordinate-based search for
339 SARS-CoV-2 gene and protein functions, mutations found in patient samples,
340 structural and sequence elements of the virus RNA and several other features. We
341 reviewed, analysed, and provided visualization for data that could help users to better
342 understand the virus, and to do this rapidly. We will continue to add mutation data as
343 we process other samples with GROM.

344

345 **CONFLICT OF INTEREST**

346 The Authors declare no conflict of interest

347

348 **FUNDING**

349 The work in A.G.'s lab is supported by the National Science Foundation [MCB-
350 2027611 to A.G.] and National Institutes of Health [R15CA220059 to A.G.].

351 **REFERENCES**

- 352 1. Wu F, Zhao S, Yu B, Chen YM, Wang W, Song ZG, et al. A new coronavirus
353 associated with human respiratory disease in China. *Nature*. 2020;579(7798):265-9.
354 2. Cucinotta D, Vanelli M. WHO Declares COVID-19 a Pandemic. *Acta Biomed*.
355 2020;91(1):157-60.
356 3. Karim SSA, Karim QA. Omicron SARS-CoV-2 variant: a new chapter in the
357 COVID-19 pandemic. *The Lancet*. 2021;398(10317):2126-8.
358 4. Finkel Y, Mizrahi O, Nachshon A, Weingarten-Gabbay S, Morgenstern D, Yahalom-
359 Ronen Y, et al. The coding capacity of SARS-CoV-2. *Nature*. 2021;589(7840):125-30.
360 5. Kim D, Lee JY, Yang JS, Kim JW, Kim VN, Chang H. The Architecture of SARS-
361 CoV-2 Transcriptome. *Cell*. 2020;181(4):914-21.e10.
362 6. Manfredonia I, Nithin C, Ponce-Salvatierra A, Ghosh P, Wirecki TK, Marinus T, et
363 al. Genome-wide mapping of SARS-CoV-2 RNA structures identifies therapeutically-relevant
364 elements. *Nucleic Acids Research*. 2020;48(22):12436-52.
365 7. Wu J, Shi Y, Pan X, Wu S, Hou R, Zhang Y, et al. SARS-CoV-2 ORF9b inhibits
366 RIG-I-MAVS antiviral signaling by interrupting K63-linked ubiquitination of NEMO. *Cell*
367 *Reports*. 2021;34(7):108761.
368 8. Hodcroft EB. CoVariants: SARS-CoV-2 Mutations and Variants of Interest. 2021
369 [Available from: <https://covariants.org/>.
370 9. Khare S, Gurry C, Freitas L, Schultz MB, Bach G, Diallo A, et al. GISAID's Role in
371 Pandemic Response. *China CDC Wkly*. 2021;3(49):1049-51.
372 10. Cao Y, Wang J, Jian F, Xiao T, Song W, Yisimayi A, et al. Omicron escapes the
373 majority of existing SARS-CoV-2 neutralizing antibodies. *Nature*. 2022;602(7898):657-63.
374 11. Sun L, Li P, Ju X, Rao J, Huang W, Ren L, et al. In vivo structural characterization of
375 the SARS-CoV-2 RNA genome identifies host proteins vulnerable to repurposed drugs. *Cell*.
376 2021;184(7):1865-83.e20.
377 12. Gobeil SMC, Henderson R, Stalls V, Janowska K, Huang X, May A, et al. Structural
378 diversity of the SARS-CoV-2 Omicron spike. *Molecular Cell*. 2022;82(11):2050-68.e6.
379 13. Database resources of the National Center for Biotechnology Information. *Nucleic*
380 *Acids Res*. 2016;44(D1):D7-19.
381 14. Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*.
382 2018;34(18):3094-100.
383 15. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence
384 Alignment/Map format and SAMtools. *Bioinformatics*. 2009;25(16):2078-9.
385 16. Smith SD, Kawash JK, Grigoriev A. Lightning-fast genome variant detection with
386 GROM. *GigaScience*. 2017;6(10):gix091.
387 17. Okonechnikov K, Golosova O, Fursov M, team tU. Unipro UGENE: a unified
388 bioinformatics toolkit. *Bioinformatics*. 2012;28(8):1166-7.
389 18. Zheng W, Zhang C, Li Y, Pearce R, Bell EW, Zhang Y. Folding non-homologous
390 proteins by coupling deep-learning contact maps with I-TASSER assembly simulations. *Cell*
391 *Reports Methods*. 2021;1(3):100014.
392 19. Yang SL, DeFalco L, Anderson DE, Zhang Y, Aw JGA, Lim SY, et al.
393 Comprehensive mapping of SARS-CoV-2 interactions in vivo reveals functional virus-host
394 interactions. *Nature Communications*. 2021;12(1):5113.
395 20. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al.
396 BLAST+: architecture and applications. *BMC Bioinformatics*. 2009;10:421.
397 21. Thompson JD, Higgins DG, Gibson TJ. CLUSTAL W: improving the sensitivity of
398 progressive multiple sequence alignment through sequence weighting, position-specific gap
399 penalties and weight matrix choice. *Nucleic Acids Res*. 1994;22(22):4673-80.
400 22. Inc F. Canvas.js. non-commercial 3.6.6 ed: Fenopix Inc.

401