1  **List of email addresses and ORCIDs for all authors**:

2  Daniele Mercatelli, daniele.mercatelli2@unibo.it, ORCID 0000-0003-3228-0580

3  Luca Triboli, luca.triboli@studio.unibo.it, ORCID 0000-0002-1261-0637

4  Eleonora Fornasari, eleonora.fornasari@ordingbo.it, ORCID 0000-0002-7636-085X

5  Forest Ray, forest.ray@zoho.com, ORCID 0000-0002-8655-7066

6  Federico M. Giorgi, federico.giorgi@unibo.it, ORCID 0000-0002-7325-9908

7

8  *coronapp*: A Web Application to Annotate and Monitor

9  SARS-CoV-2 Mutations

10  Daniele Mercatelli[1,#], Luca Triboli[1,#], Eleonora Fornasari[1], Forest Ray[2], Federico M.

11  Giorgi[1,*]

12  [1] *Department of Pharmacy and Biotechnology, University of Bologna, Bologna,*

13  *40126, Italy*

14  [2] *Department of Systems Biology, Columbia University Medical Center, New York*

15  *City, 10032, United States*

16  [#] Equal contribution.

17  [*] Corresponding author.

18  E-mail: federico.giorgi@unibo.it (Giorgi FM)

19

20  **Running title:** *Mercatelli D et al / coronapp – monitoring SARS-CoV-2 mutations*

21

22  Word number: 3531

23  Figure number: 3

24

25

## Abstract

27  The avalanche of genomic data generated from the SARS-CoV-2 virus requires the

28  development of tools to detect and monitor its mutations across the World. Here, we

29  present a webtool, *coronapp*, dedicated to easily processing user-provided

30  SARS-CoV-2 genomic sequences, in order to detect and annotate protein-changing

31  mutations. This results in an up-to-date status of SARS-CoV-2 mutations, both

32  worldwide and in user-selected countries. The tool allows users to highlight and

33  prioritize the most frequent mutations in specific protein regions, and to monitor their

34  frequency in the population over time.

35  The tool is available at http://giorgilab.dyndns.org/coronapp/ and the full code is

36  freely shared at https://github.com/federicogiorgi/giorgilab/tree/master/coronapp

42  **KEYWORDS:** COVID-19; SARS-CoV-2; mutations; web application

## Introduction

46 SARS-CoV-2 is a novel pathogenic enveloped RNA beta-coronavirus causing a

47 severe illness in human hosts known as coronavirus disease-2019 (COVID-19). The

48 predominant COVID-19 illness is a viral pneumonia, often requiring hospitalization

49 and in some cases intensive care [1]. With almost 6 million laboratory-confirmed

50 positive cases worldwide as of 31 May 2020 and an estimated case fatality rate across

51 204 countries of 5.2%, COVID-19 has become a global health challenge in only a few

52 months [2]. SARS-CoV-2 infection depends on the recognition of host angiotensin

53 converting enzyme 2 (ACE2), exposed on the cell surface in human lung tissues [3,4].

54 SARS-CoV-2 spike glycoprotein binds ACE2, mediating membrane fusion and cell

55 entry [5]. Upon cell entry, the virus subverts host cell molecular processes, inducing

56 interferon responses and eventually apoptosis [6].

57 To date, much effort has been made to develop therapeutic strategies to limit

58 SARS-CoV-2 transmission and replication, but no treatment or vaccine has proven

59 effective against the virus, and repurposing of approved therapeutic agents has been

60 the main practical approach to manage the emergency so far [7]. As viruses mutate

61 during replication, the emergence of SARS-CoV-2 sub-strains and the challenge of a

62 probable antigenic drift require attention, especially for vaccine development [8].

63 Although sequence analyses of SARS-CoV-2 have shown that genomic variability

64 is very low [9], new SARS-CoV-2 mutation hotspots are emerging due to the high

65 number of infected individuals across countries and to viral replication rates [10].

66 Three major SARS-CoV-2 clades known as clade G, V, and S have emerged, showing

67 a different geographical prevalence [10]. The most frequent mutation detected so far

68 defines the G clade and causes an aminoacidic change, aspartate (D) or glycine (G), at

69 position 614 (D614G) of the viral Spike protein [11].

70 Continual genomic surveillance should be considered to monitor the possible

71 appearance of viral subtypes characterized by altered tropism, or causing more

72 aggressive symptoms. Constant and widespread monitoring of mutations is also a

73  powerful means of informing drug development and global or local pandemic

74  management. The Global Initiative on Sharing All Influenza Data (GISAID) has

75  collected to date (31 May 2020) over 30,000 publicly accessible SARS-CoV-2

76  sequences. The GISAID effort has made it possible to compare genomes on a

77  geographical and temporal scale and an increasing number of laboratories have started

78  to sequence COVID-19 patient samples worldwide [13,14]. Several online tools have

79  been developed to monitor the evolution of the virus from a phylogenetic perspective,

80  such as Nextstrain [15], or to visualize epidemiological data such as number of cases

81  and deaths [16]. However, no tool currently exists to annotate user-provided

82  SARS-CoV-2 genomic sequences, which may derive from specific GISAID subsets

83  or from sequencing efforts of individual laboratories. Neither does any tool

84  specifically monitor the prevalence of specific SARS-CoV-2 mutations associated to

85  particular geographic regions or protein locations, nor their frequency in the

86  population over time.

87      To overcome these limitations, we have developed *coronapp*, a web application

88  with two purposes: real-time tracking of SARS-CoV-2 mutational status and

89  annotation of user-provided viral genomic sequences. Our tool enables users to easily

90  perform genomic comparisons and provides an instrument to monitor SARS-CoV-2

91  genomic variance, both worldwide and by uploading custom and locally produced

92  genomic sequences. The webtool is available at http://giorgilab.dyndns.org/coronapp/

93  and  the  full  source  code  is  shared  on  Github

94  https://github.com/federicogiorgi/giorgilab/tree/master/coronapp

95

96  **Results**

97  The  webtool  *coronapp*  is  available  at  the  website

98  http://giorgilab.dyndns.org/coronapp/ and it automatically provides the user with the

99  current status of SARS-CoV-2 mutations worldwide. The app also allows users to

100  annotate user-provided sequences (Figure 1 A). There are multiple functionalities of

101  *coronapp*, described in the following paragraphs.

102

**Current Status of SARS-CoV-2 mutational data**

104  A worldwide analysis is shown, generated using data from GISAID. Specifically, we

105  processed all SARS-CoV-2 complete (>29,000 sequenced nucleotides) genomic

106  sequences, excluding low-quality sequences (>5% undefined nucleotide "N") and

107  viruses extracted from non-human hosts.

108      The underlying database is updated weekly, and we provide the date of the last

109  version as a reference for studies based on the data provided. We indicate the number

110  of samples processed and the total number of mutational events detected (Figure 1 A).

111  We also show the number of distinct mutated loci. Currently, this number is slightly

112  below 11,000, meaning that less than half of the original Wuhan SARS-CoV-2

113  genome has been affected by mutations and/or sequencing errors (the full length of

114  the reference genome is 29,903 nucleotides, based on sequence id NC_045512.2).

115

**Mutation frequency in SARS-CoV-2 proteins**

117  We show the frequency of mutations along the length of every SARS-CoV-2 protein,

118  reporting in the X-axis the amino acid position and on the Y-axis its frequency, either

119  as number of observed samples carrying the mutation, the vase 10 logarithm of that

120  number, or the percentage over all sequenced samples. In the example in Figure 1 B,

121  we show the most frequent mutations affecting the viral Spike protein S,

122  distinguishing silent mutations and amino acid-changing mutations (including the

123  introduction of STOP codons). For Spike, the mutations appear to be evenly

124  distributed in frequency along the protein length, with the most frequent mutation

125  being the aforementioned D614G. Mouse-over functionality is provided to allow the

126  user to identify the selected mutation (N439K in Figure 1 B).

127

128    **The SARS-CoV-2 mutation table**

129    The user can visualize or download the full table of mutations on which the webtool

130    operates (Figure 2 A). This table is frequently updated and allows the user to specify a

131    worldwide or a country-specific dataset. The table also provides a Search function to

132    look for specific variants or sample ids, and it can be viewed online or downloaded in

133    full as a Comma-Separated Values (CSV) file.

134        The table shows every mutation in a specific geographical area, reporting:

135    • the GISAID sample ID (useful for cross-reference with the GISAID database

136        and other analyses based on it, e.g. Nexstrain).

137    • The country where the sample was collected.

138    • The position of the mutation, on the reference genome (refpos) and on the

139        sample (qpos).

140    • The sequence at the mutation site, on the reference genome (refvar) and on the

141        sample (qvar).

142    • The length of the sample genome (qlength); the reference genome is 29,903

143        nucleotides long.

144    • The protein affected by the mutation or, if the mutation is extragenic, the

145        denomination of the untranslated region (UTR), e.g. 5'UTR or 3'UTR.

146    • The effect of the mutation on the amino acid sequence of the protein (variant).

147        This uses the canonical mutational standard, indicating the original amino

148        acid(s), the position on the protein, and the mutated amino acid(s). An asterisk

149        (*) indicates a STOP codon, while the letters indicate amino acids in IUPAC

150        code. E.g. a mutation P315L indicates a leucine mutation (L) on the amino

151        acid location 315, normally occupied by a proline (P). Nucleotide mutations

152        can be silent, i.e. not yielding any aminoacidic change, e.g. the mutation

153        F106F, where the codon of phenylalanine 106 is affected but without changing

154        the corresponding amino acid. As in the previous column, mutations affecting

155        UTR regions are simply reported as the location of the nucleotide affected.

156 • The class of the mutation, of which there are currently 10 types:

157 o SNP: a change of one or more nucleotides, determining a change in
158 amino acid sequence.

159 o SNP_stop: a change of one or more nucleotides, yielding the generation
160 of one or more STOP codons.

161 o SNP_silent: a change of one or more nucleotides with no effect in
162 protein sequence.

163 o Insertion: the insertion of 3 (or multiples of 3) nucleotides, causing the
164 addition of 1 or more amino acids to the protein sequence.

165 o Insertion_stop: the insertion of 3 (or multiples of 3) nucleotides, causing
166 the generation of a novel STOP codon.

167 o Insertion_frameshift: the insertion of nucleotides not as multiples of 3,
168 causing a frameshift mutation.

169 o Deletion: the deletion of 3 (or multiples of 3) nucleotides, causing the
170 removal of 1 or more amino acids to the protein sequence.

171 o Deletion_stop: the removal of 3 (or multiples of 3) nucleotides, causing
172 the generation of a novel STOP codon.

173 o Deletion_frameshift: the deletion of nucleotides not as multiples of 3,
174 causing a frameshift mutation.

175 o Extragenic: a mutation affecting intergenic or UTR regions.

176 • The extended annotation of the protein region affected by the mutation (e.g.
177 "Spike" for "S" or "Predicted phosphoesterase, papain-like proteinase" for
178 NSP3, the Non-Structural Protein 3).

179 • The full name of the variant (varname), in the format
180 proteinName:AApositionAA, to allow for unique denomination of viral
181 proteome variants.

182

183 **Mutational overview**

184 The user is also provided with a general overview of the mutational status of the

185 selected country or the entire world (Figure 2 B). Six bar plots provide a summary and

186 highlights of the dataset, specifically:

187 • The most mutated samples, indicating which samples (in GISAID IDs) carry

188 the highest number of mutations

189 • The overall mutations per sample, indicating the distributions of mutations per

190 sample. It has been previously reported [10] that the current mode for

191 mutation number compared to the reference NC_045512.2 genome is 7.5.

192 • The most frequent events per class. Classes are the same as reported in the

193 mutation table and are described in the previous paragraph.

194 • The most frequent events per type. Individual mutation types are shown as

195 specific nucleotides events, e.g. cytosine to thymidine transitions (C>T),

196 guanosine to thymidine transversion (G>T) or even multinucleotide mutations

197 (e.g. GGG>AAC, observed in the Nucleocapsid protein). As reported before,

198 nucleotide transitions seem to be the most abundant SARS-CoV-2 type of

199 mutational event worldwide [11].

200 • The most frequent events, either in nucleotide coordinates or in aminoacidic

201 coordinates. Currently, the most frequent events are four mutations affecting

202 SARS-CoV-2 genomes belonging to clade G, which is the most sequenced

203 worldwide and predominant in Europe. These mutations are A23403G

204 (associated to the already mentioned D614G mutation in the Spike protein),

205 C3037T, C14408T and C241T.

206

207 **Analysis of mutations over time**

208 The *coronapp* webtool allows users to monitor the abundance and frequency of any

209 SARS-CoV-2 mutation in any country specified (Figure 3). Both plots in this section

210 report continuous dates on the X-axis, starting on the day of the first collected

211 SARS-CoV-2 genome available on GISAID: December 24, 2019.

212      The "abundance" plot reports on the Y-axis the number of samples carrying a

213    selected mutation in a particular day, in the specified country or worldwide. Since the

214    date reported is the collection date (not the submission date to the GISAID database),

215    there is usually a drop towards the right part of the plot, as there are fewer sequences

216    collected approaching the day of the analysis. The "frequency" plot on the other hand

217    normalizes the abundance of mutations by the total number of sequences generated on

218    each day. The plot currently shows a sharp increase in clade G-associated mutations

219    (e.g. S:D614G), as these mutations are most frequent in countries where sequencing is

220    more pervasive (e.g. United Kingdom).

221

222    **Annotation of user-provided SARS-CoV-2 genomic sequence.**

223    *coronapp* provides the user with the optional possibility of uploading one or more

224    SARS-CoV-2 genomic sequences, which can be complete or partial. The format of

225    the sequences is standard FASTA, and an example input FASTA containing 12

226    sequences is provided (Figure 1 A). The analysis is almost instantaneous and shows

227    an overall breakdown of the most mutated samples and most frequent mutations in the

228    dataset. Moreover, a full table of all detected mutations is provided: this can be

229    visualized and searched on the web browser or downloaded as a standard CSV file.

230    Finally, a mutation frequency plot is provided, allowing the user to visualize mutation

231    frequency in selected proteins.

232      The user can easily return to the worldwide status of the app by refreshing or

233    reopening the page.

234

235    **Discussion**

236    Our webtool *coronapp* provides a fast, simple tool to annotate user-provided

237    SARS-CoV-2 genomes and visualize all mutations currently present in viral

238    sequences collected worldwide. The results provided by this instrument can have

239    several applications. The main purpose of *coronapp* is to help medical laboratories at

240 the front lines of COVID-19 fight with the opportunity to quickly define the

241 mutational status of their sequences, even without dedicated bioinformaticians.

242  Additionally, it enables scientists to perform mutational co-variance analyses and

243 to identify present and future significant functional interactions between viral

244 mutations, as previously attempted for the influenza virus and the human

245 immunodeficiency virus (HIV) [17]. Another application is the identification of the

246 most frequent mutations in specific protein regions: for example, our tool can quickly

247 identify that the most frequent mutation in the Spike protein, D614G, lies outside the

248 known interaction domain with the human protein ACE2, which spans roughly

249 between Spike amino acids 330 and 530 [18].

250  A recently published structural model simulating the effect of the D614G mutation

251 on the 3D structure of the spike protein has suggested that this mutation may result in

252 a viral particle which binds ACE2 receptors less efficiently, due to the masking of the

253 host receptor binding site on viral spikes [12]. The same researchers have reported a

254 possible correlation of the D614G form with increased case fatality rates,

255 hypothesizing that this mutation may lead to a viral form which is better suited to

256 escape immunologic surveillance by eliciting a lower immunologic response [12].

257 The *coronapp* analysis highlighted in Figure 1 B shows that a mutation located within

258 the Spike/ACE2 interaction domain is the change of Asparagine (N) to a Lysine (K)

259 in position 439 of the Spike sequence; this mutation could affect the protein folding or

260 its affinity with ACE2, as Asparagine is less charged than the basic amino acid

261 Lysine.

262  One of *coronapp*'s key strengths is to help prioritize scientific efforts on specific

263 aminoacidic variations that could affect the efficacy of anti-viral strategies or the

264 development of a vaccine by tracking the most frequent mutations in the population.

265 A further novelty of *coronapp* is that it provides a mean to assess the growth or

266 decline of specific mutations over time, in order to identify possible viral adaptation

267 mechanisms.

268    We provide not only the webtool, but also all the underlying code for the

269    annotation and visualization steps on a public Github repository, in order to help other

270    computational scientists in the ongoing battle against COVID-19. Furthermore, the

271    *coronapp* structure and concept could be expanded to other current and future

272    pathogens as well (e.g. the seasonal influenza or HIV), in order to monitor the

273    mutational status across proteins, countries and time.

274

## Materials and methods

276    The webtool *coronapp* has been developed using the programming language R and is

277    based on a Shiny server (current version 1.4.0.2) running on R version 3.6.1. The app

278    is based on two distinct files, server.R and ui.R, managing the server functionalities

279    and the browser visualization processes, respectively. The results visualization utilizes

280    both basic R functions and Shiny functionalities; for tooltip functionality, *coronapp*

281    uses the R package *googleVis* v0.6.4, which provides an interface between R and the

282    Google visualization API [19].

283    The core of the annotation of the user-provided sequences rests in the NUCMER

284    (Nucleotide Mummer) alignment tool, version 3.1 [20]. Nucmer output is processed

285    by UNIX and R scripts provided in Github within the server.R file.

286

287

## Authors' contributions

DM drafted the manuscript and performed the mutational analysis and literature search. LT developed the user interface code and drafted the methodological parts of the manuscript. EF worked on graphical interface of the webtool. FR wrote the manuscript and performed literature search. FMG designed the study, developed the server code, finalized the manuscript and provided financial support. All authors tested the webtool and provided original contributions to its development. All authors read and approve the final manuscript.

## Competing interests

The authors have declared no competing interests.

## Acknowledgements

## References

[1] Guan W-J, Ni Z-Y, Hu Y, Liang W-H, Ou C-Q, He J-X, et al. Clinical Characteristics of Coronavirus Disease 2019 in China. N Engl J Med 2020;382:1708–20. https://doi.org/10.1056/NEJMoa2002032.

[2] Phua J, Weng L, Ling L, Egi M, Lim C-M, Divatia JV, et al. Intensive care management of coronavirus disease 2019 (COVID-19): challenges and recommendations. Lancet Respir Med 2020. https://doi.org/10.1016/S2213-2600(20)30161-2.

[3] Zhang H, Penninger JM, Li Y, Zhong N, Slutsky AS. Angiotensin-converting enzyme 2 (ACE2) as a SARS-CoV-2 receptor: molecular mechanisms and potential therapeutic target. Intensive Care Med 2020;46:586–90. https://doi.org/10.1007/s00134-020-05985-9.

[4] Guzzi PH, Mercatelli D, Ceraolo C, Giorgi FM. Master Regulator Analysis of the SARS-CoV-2/Human Interactome. J Clin Med 2020;9:982. https://doi.org/10.3390/jcm9040982.

[5] Ou X, Liu Y, Lei X, Li P, Mi D, Ren L, et al. Characterization of spike glycoprotein of SARS-CoV-2 on virus entry and its immune cross-reactivity with

321       SARS-CoV. Nat Commun 2020;11:1620.
322       https://doi.org/10.1038/s41467-020-15562-9.
323  [6]  Blanco-Melo D, Nilsson-Payant BE, Liu W-C, Uhl S, Hoagland D, Møller R, et
324       al. Imbalanced Host Response to SARS-CoV-2 Drives Development of
325       COVID-19. Cell 2020;181:1036-1045.e9.
326       https://doi.org/10.1016/j.cell.2020.04.026.
327  [7]  Tu Y-F, Chien C-S, Yarmishyn AA, Lin Y-Y, Luo Y-H, Lin Y-T, et al. A Review
328       of SARS-CoV-2 and the Ongoing Clinical Trials. Int J Mol Sci 2020;21.
329       https://doi.org/10.3390/ijms21072657.
330  [8]  Koyama T, Weeraratne D, Snowdon JL, Parida L. Emergence of Drift Variants
331       That May Affect COVID-19 Vaccine Development and Antibody Treatment.
332       Pathog Basel Switz 2020;9. https://doi.org/10.3390/pathogens9050324.
333  [9]  Ceraolo C, Giorgi FM. Genomic variance of the 2019□nCoV coronavirus. J Med
334       Virol 2020;92:522–8. https://doi.org/10.1002/jmv.25700.
335  [10] Mercatelli D, Giorgi FM. Geographic and Genomic Distribution of SARS-CoV-2
336       Mutations. Preprints; 2020. https://doi.org/10.20944/preprints202004.0529.v1.
337  [11] Chiara M, Horner DS, Gissi C, Pesole G. Comparative genomics suggests limited
338       variability and similar evolutionary patterns between major clades of
339       SARS-CoV-2. BioRxiv; 2020. https://doi.org/10.1101/2020.03.30.016790.
340  [12] Becerra-Flores M, Cardozo T. SARS-CoV-2 viral spike G614 mutation exhibits
341       higher case fatality rate. Int J Clin Pract 2020. https://doi.org/10.1111/ijcp.13525.
342  [13] Gudbjartsson DF, Helgason A, Jonsson H, Magnusson OT, Melsted P, Norddahl
343       GL, et al. Spread of SARS-CoV-2 in the Icelandic Population. N Engl J Med
344       2020. https://doi.org/10.1056/NEJMoa2006100.
345  [14] Fauver JR, Petrone ME, Hodcroft EB, Shioda K, Ehrlich HY, Watts AG, et al.
346       Coast-to-Coast Spread of SARS-CoV-2 during the Early Epidemic in the United
347       States. Cell 2020;181:990-996.e5. https://doi.org/10.1016/j.cell.2020.04.021.
348  [15] Hadfield J, Megill C, Bell SM, Huddleston J, Potter B, Callender C, et al.
349       Nextstrain: real-time tracking of pathogen evolution. Bioinformatics
350       2018;34:4121–3. https://doi.org/10.1093/bioinformatics/bty407.
351  [16] Max Roser EO-O Hannah Ritchie, Hasell J. Coronavirus Pandemic (COVID-19).
352       Our World Data 2020.
353  [17] Sruthi CK, Prakash MK. Statistical characteristics of amino acid covariance as
354       possible descriptors of viral genomic complexity. Sci Rep 2019;9:18410.
355       https://doi.org/10.1038/s41598-019-54720-y.
356  [18] Lan J, Ge J, Yu J, Shan S, Zhou H, Fan S, et al. Structure of the SARS-CoV-2
357       spike receptor-binding domain bound to the ACE2 receptor. Nature
358       2020;581:215–20. https://doi.org/10.1038/s41586-020-2180-5.
359  [19] Gesmann M, de Castillo D. Using the Google visualisation API with R. R J
360       2011;3:40–44.

361    [20]Delcher AL, Salzberg SL, Phillippy AM. Using MUMmer to Identify Similar

362        Regions in Large Sequence Sets. Curr Protoc Bioinforma 2003;00:10.3.1-10.3.18.

363        https://doi.org/10.1002/0471250953.bi1003s00.

364

365    **Figure legends**

366    **Figure 1 Overview of *coronapp***

367    **A**. Screenshot of the entry page of *coronapp* showing the basic tool description, the

368    interface to upload user-provided sequences and the overall summary of the mutations

369    detected worldwide. **B**. Common interface showing mutation frequency in

370    SARS-CoV-2 proteins, with occurrence of the mutation on the Y-axis and protein

371    coordinate on the Y-axis. Red dots indicate amino acid (aa)-changing mutations, and

372    blue dots indicate silent mutations. Tooltip functionality is also provided to identify

373    and quantify each mutation on mouse-over.

374

375    **Figure 2 Mutation table and overview in *coronapp***

376    **A**. Result table of *coronapp*, available both for worldwide-precomputed and

377    user-input analyses. A "download full table" button is provided to allow the user to

378    perform larger-scale analyses autonomously. **B**. Barplots showing the most mutated

379    samples, overall sample mutations and most frequent mutation events, classes and

380    types. This analysis is also available both for worldwide-precomputed and user-input

381    analyses.

382

383    **Figure 3 Analysis of mutations over time**

384    The final output of *coronapp*, showing the abundance of each user-specified mutation

385    in any user-specified country (or worldwide). The left graph indicates the absolute

386    amount of samples where the indicated mutation is detected. The right graph shows

387    the same data normalized by total number of samples, as the percentage of samples

388    sequenced in a specific day and carrying the mutation.

389

# A

**COVID-19 genome annotator** ☰

*coronapp* is a web application written in Shiny with a double purpose:

- Monitor SARS-CoV-2 worldwide mutations
- Annotate user-provided mutations

**Provide your own (multi)FASTA file**

| Browse... | No file selected |
|-----------|------------------|

Example input multiFASTA

The FASTA annotator will discover and annotate every mutation present in the uploaded SARS-CoV-2 genomic sequences, even partial. The GFF3 genome annotation file is available here

## Current Status of SARS-CoV-2 mutational data

updated May 30, 2020

Number of samples: 29668
Number of distinct mutated loci: 10458
Total number of mutational events: 203292

**Select Country:**
World ▾

**Select Protein:**
S ▾

☑ Log10

☐ Percentage

**Mutation frequency for protein S (Spike) in World**

# B



**Mutation frequency for protein S (Spike) in World**

■ aa change
■ silent

N439K
aa: **439**
occurrence: **2.111**
status: **aa change**

*Occurrence of event (Log10)*

*aa coordinate*

**A**

## Showing results for World

⬇ Download Full table (CSV format)

Show [10 ▾] entries

Search: [＿＿＿＿＿]

| sample | country | refpos | refvar | qvar | qpos | qlength | protein | variant | varclass | annotation | varname |
|--------|---------|--------|--------|------|------|---------|---------|---------|----------|------------|---------|
| EPI_ISL_415706 | Switzerland | 4 | A | T | 4 | 29903 | 5'UTR | 4 | extragenic | | 5'UTR:4 |
| EPI_ISL_415706 | Switzerland | 241 | C | T | 241 | 29903 | 5'UTR | 241 | extragenic | | 5'UTR:241 |
| EPI_ISL_415706 | Switzerland | 3037 | C | T | 3037 | 29903 | NSP3 | F106F | SNP_silent | Predicted phosphoesterase, papain-like proteinase | NSP3:F106F |
| EPI_ISL_415706 | Switzerland | 14408 | C | T | 14408 | 29903 | NSP12b | P314L | SNP | RNA-dependent RNA polymerase, post-ribosomal frameshift | NSP12b:P314L |
| EPI_ISL_415706 | Switzerland | 15324 | C | T | 15324 | 29903 | NSP12b | N619N | SNP_silent | RNA-dependent RNA polymerase, post-ribosomal frameshift | NSP12b:N619N |
| EPI_ISL_415706 | Switzerland | 23403 | A | G | 23403 | 29903 | S | D614G | SNP | Spike | S:D614G |
| EPI_ISL_416497 | France | 4 | | | | 29862 | 5'UTR | 4 | extragenic | | 5'UTR:4 |
| EPI_ISL_416497 | France | 241 | C | T | 241 | 29862 | 5'UTR | 241 | extragenic | | 5'UTR:241 |
| EPI_ISL_416497 | France | 2416 | C | T | 2416 | 29862 | NSP2 | Y537Y | SNP_silent | Non-Structural protein 2 | NSP2:Y537Y |
| EPI_ISL_416497 | France | 3037 | C | T | 3037 | 29862 | NSP3 | F106F | SNP_silent | Predicted phosphoesterase, papain-like proteinase | NSP3:F106F |

Showing 1 to 10 of 192,208 entries
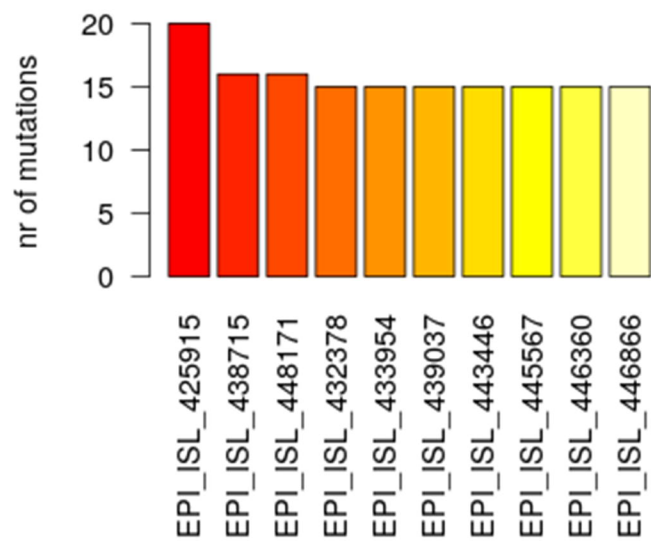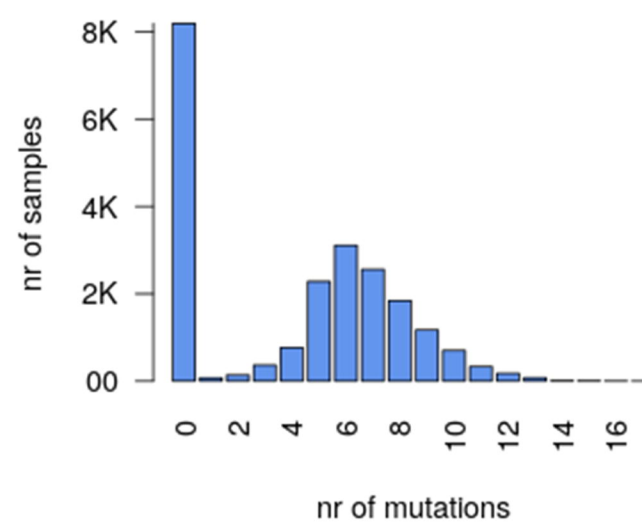
Previous  1  2  3  4  5  …  19221  Next

**B**

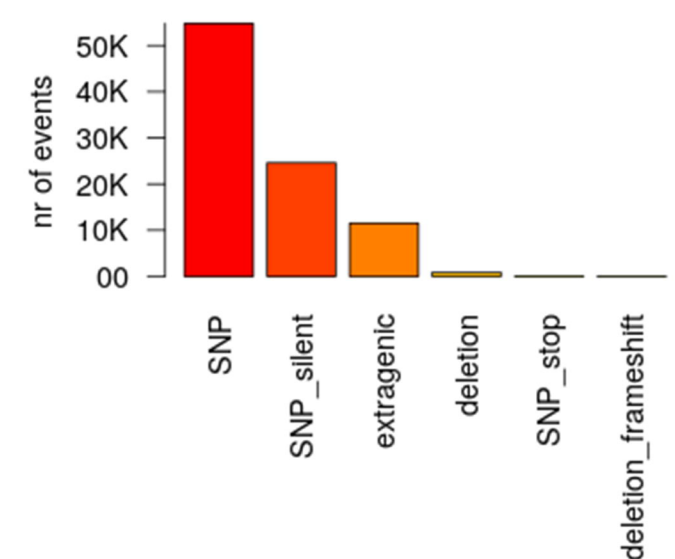## Mutational overview for United Kingdom

# Analysis of mutations over time

**Select Country:**

World ▼

**Select Mutation:**

S:D614G ▼

## S:D614G abundance in World



## S:D614G frequency in World