

1 Exploring the coronavirus epidemic 2 using the new WashU Virus Genome 3 Browser

4
5 Jennifer A. Flynn^{1*}, Deepak Purushotham^{1*}, Mayank NK Choudhary^{1*}, Xiaoyu Zhuo^{1*},
6 Changxu Fan^{1*}, Gavriel Matt^{1*}, Daofeng Li^{1†} and Ting Wang^{1,2†}

7
8 * These authors contributed equally to this work.

9 † These authors jointly supervised this work. Co-corresponding author emails: dli23@wustl.edu
10 and twang@genetics.wustl.edu

11 ¹The Edison Family Center for Genome Sciences & Systems Biology, Department of Genetics,
12 Washington University, 4515 McKinley Avenue, Campus Box 8510, St. Louis, MO 63110, USA

13 ²McDonnell Genome Institute, Washington University School of Medicine, St. Louis, MO 63108,
14 USA

15 Abstract

16 Since its debut in mid-December, 2019, the novel coronavirus (2019-nCoV) has rapidly spread
17 from its origin in Wuhan, China, to several countries across the globe, leading to a global health
18 crisis. As of February 7, 2020, 44 strains of the virus have been sequenced and uploaded to
19 NCBI's GenBank [1], providing insight into the virus's evolutionary history and pathogenesis.
20 Here, we present the WashU Virus Genome Browser, a web-based portal for viewing virus
21 genomic data. The browser is home to 16 complete 2019-nCoV genome sequences, together
22 with hundreds of related viral sequences including severe acute respiratory syndrome
23 coronavirus (SARS-CoV), Middle East respiratory syndrome coronavirus (MERS-CoV), and
24 Ebola virus. In addition, the browser features unique customizability, supporting user-provided
25 upload of novel viral sequences in various formats. Sequences can be viewed in both a track-
26 based representation as well as a phylogenetic tree-based view, allowing the user to easily
27 compare sequence features across multiple strains. The WashU Virus Genome Browser
28 inherited many features and track types from the WashU Epigenome Browser, and additionally
29 incorporated a new type of SNV track to address the specific needs of viral research. Our Virus
30 Browser portal can be accessed at <https://virusgateway.wustl.edu>, and documentation is
31 available at <https://virusgateway.readthedocs.io/>.

32 Introduction

33 On December 12, 2019, the first case of a novel coronavirus (2019-nCoV) was reported in
34 Wuhan, China, and by February 6, 2020, the virus spread to 24 additional countries, infecting
35 more than 27,000 individuals and resulting in 565 fatalities, according to the World Health
36 Organization (WHO) [2]. The 2019-nCoV is a member of the *Betacoronavirus* genus, which is
37 one of four genera of coronaviruses of the subfamily Orthocoronavirinae in the family
38 Coronaviridae, of the order Nidovirales [3, 4]. The species in this genus are enveloped, contain
39 a positive single-stranded RNA genome, and are of zoonotic, likely bat, origins [5]. 2019-nCoV
40 is one of the largest RNA virus genomes varying from 27kb to 32kb in size, with this particular

41 strain ringing in at 29,903 bps long [6]. The virus is one of 7 coronaviruses known to infect
42 humans, and along with the severe acute respiratory syndrome coronavirus (SARS-CoV) and
43 the Middle East respiratory syndrome coronavirus (MERS-CoV), 2019-nCoV is one of the
44 species responsible for severe respiratory distress in humans as well as other animals [4]. In an
45 effort to better understand the pathogenesis of this family of viruses, several groups have
46 sequenced individual strains, providing a powerful resource hosted by NCBI.

47
48 The WashU Epigenome Browser is a powerful tool for visualizing multiple functional genomic
49 datasets and data types simultaneously [5-8]. The general layout of the Epigenome Browser
50 displays the genome on the x-axis, and individual tracks encompassing many different varieties
51 can be loaded and viewed in the context of the genome and accompanying metadata. Recent
52 updates to the browser have incorporated new functionality, including live browsing, greatly
53 enhancing its functionality [5]. With this powerful tool in-hand, we sought to adapt the browser
54 for use of visualizing viral genomes, to support more efficient research and more rapid
55 knowledge dissemination in response to the recent 2019-nCoV outbreak. To accomplish this,
56 we created the WashU Virus Genome Browser, adapted from the WashU Epigenome Browser.
57 The Virus Genome Browser houses reference genomes for 2019-nCoV, MERS, SARS, and
58 Ebola virus, along with several annotation tracks including gene annotation, putative antibody-
59 binding epitopes, CG density, and sequence diversity. Complete genomes of individual strains
60 of each virus species (16, 551, 332, and 1574, respectively as of February 7, 2020, and
61 periodically updated) are available as a database for instant viewing on the Virus Browser via
62 multiple track types designed to display pairwise comparison to the references. Additionally, we
63 aligned the genomes of all available strains in the database and generated a phylogenetic tree
64 for each virus species that allows the user to directly select strains from the tree and view as
65 tracks in the genomic display. In addition to all track types supported by the Epigenome
66 Browser, we designed a new SNV track type to display sequence variation. Users can upload
67 their own alignment results from any aligner and display them as SNV tracks on the browser.

68
69 The functionality of the Virus Browser is not limited to the 4 species currently housed. Users can
70 upload their own reference genome in FASTA format and display tracks in the context of the
71 user-specified reference. While maintaining the same functionality as that of the Epigenome
72 Browser and providing novel functionality to aid specifically in viral genome research, we hope
73 that the Virus Browser may facilitate research against new epidemic viruses.

74

75 Materials and Methods

76 Reference sequences, additional strains, and gene annotations:

77 Genomic sequences of all viral strains were downloaded as FASTA files from NCBI
78 [Supplementary Table 1]. All available sequences as of January 31, 2020, for 2019-nCoV,
79 MERS, SARS, and Ebola were downloaded (n=16, 551, 332, and 1574, respectively). The
80 reference genomic sequence of the selected virus (2019-nCoV: NC_045512.2; MERS:
81 NC_019843.3; SARS: NC_004718.3, Ebola: KM034562.1) is automatically displayed as a color
82 coded track when opening the genomic track browser viewing format. Genic annotations of
83 reference genomes were downloaded as GFF3 files from NCBI and converted to refBed format
84 for viewing on the browser.

85 Sequence alignment and tree generation:

86 The genomes of all individual strains of each virus were aligned to the reference genome using
87 the pairwise alignment tool stretcher [9] with parameters “-gapopen 16 -gapextend 4”. To
88 generate the phylogenetic trees, we used the MAFFT program, employing the fast option to
89 align individual strains of each viral genome to its reference [10, 11]. Phylogenetic trees were
90 built using FastTree with the GTR model [12, 13].

91 Data Tracks:

92 Genome Comparison Track:

93 We adopted the genome comparison tracks from the WashU Epigenome Browser. Any pairwise
94 alignment results in markx3 or FASTA format can be converted with our publicly accessible
95 script “aligned_fa_2_genomealign.py” [14] and directly displayed as genome comparison tracks
96 on the Virus Browser.

97 SNV Track:

98 We developed the SNV track type to display sequence variation of individual strains relative to
99 their reference. Variations from the reference genome, including mismatches and deletions, are
100 displayed with customizable colors. Insertions compared to the reference genome can be
101 expanded upon selecting to show the nucleotides inserted. When viewing large regions, such as
102 the whole genome, it is not possible to display all individual variation events. Therefore, the
103 frequency of variation events is also displayed in a “density mode” where a high value over a
104 region signifies multiple sequence variation events within the region.

105 Congeneric (or Closely-related) Immune Epitope Locations:

106 We wrote a text processing utility to import antibody-binding epitopes curated by the Immune
107 Epitope Database and Analysis Resource (IEDB) for MERS-CoV and SARS-CoV [15].
108 Subsequently, we used tblastn to align linear epitopes to the Wuhan seafood market pneumonia
109 virus isolate Wuhan-Hu-1 (Taxonomy ID: 2697049; NCBI:txid2697049). We found 955 out of
110 2,817 linear epitopes identified in SARS had at least 1 “hit” in the 2019-nCoV genome
111 [Supplementary Data 1]. Three epitopes have 2 “hits” each. However, the secondary hit is on
112 the negative strand with very low percent identity (37.5% to 53.8%) to the 2019-nCoV genome
113 and are hence filtered out as 2019-nCoV is a (+) ssRNA virus. Similarly, we found 1 hit out of 38
114 linear epitopes identified in MERS. We also provide scripts [14] that can be used to obtain a
115 quick overview of the similarity of linear epitopes identified in other viruses in databases like
116 IEDB. These tracks can provide researchers preliminary data to support exploratory analyses
117 pertaining to the immunogenicity of 2019-nCoV—an actively explored vertical of 2019-nCoV
118 research.

119 GC Density Track:

120 GC density tracks were created for each reference genome, displaying the percentage of G
121 (guanine) and C (cytosine) bases in 5-bp windows.

122 Sequence Diversity Track and Shannon Track:

123 In order to display a measure of sequence conservation across the genome, we calculated the
124 percentage of each of the 4 nucleotides at each position in the genome across all strains for a
125 given virus species. The resulting bed tracks display the percentages each nucleotide
126 comprises across all strain for each genomic position. We also calculated Shannon entropy for
127 each position along the genome using the percentages of each of the 4 nucleotides. A high
128 Shannon entropy at a position signifies that the 4 possible nucleotides are equally likely across
129 all strains of this virus, and thus the position is likely divergent. A low Shannon entropy at a
130 position means that the identity of the nucleotide at this position is highly conserved across all
131 strains. The entropy() function of the R package “entropy” was used for calculations.

132 Resources for User-Defined Bed and Categorical Tracks:

133 In addition to our housed data tracks, we also offer scripts (“publicParseAlignment.py”,
134 “publicAlignment.py”, and “publicConvertMarkx3.py”) to convert any markx3 or FASTA-
135 formatted alignment into displayable bed and categorical formats, and a script
136 (“publicJsonGen.py”) to generate a json file for uploading multiple data files together for display
137 [<https://github.com/debugpoint136/WashU-Virus-Genome-Browser>]. A default color code for
138 sequence variation is also included in the script.

139 Results

140 Organization of the Virus Genome Browser

141 The WashU Virus Genome Browser houses consensus reference genomic sequences for 4
142 different pathogenic virus species: 2019-nCoV, MERS, SARS, and Ebola, as well as a
143 comprehensive set of genome assemblies for the individual strains of each virus (16, 551, 332,
144 and 1574, respectively). When users first navigate to the WashU Virus Browser and select
145 “Browse Data”, they are directed to a page with several customizable options, including a drop-
146 down menu from which they may choose a reference genome [Figure 1]. Corresponding with
147 the reference genome selected, a metadata table is displayed containing sortable features such
148 as species, strain, isolate, isolation source, host, country, and collection date, to allow for quick
149 and easy sorting of individual strains. The user may select viral isolates from the metadata table
150 to be visualized in one of our two displayable platforms: the track view (green arrow, Figures 2
151 and 3) or the phylogenetic tree view (orange arrow, Figures 4 and 5).

152 The Track View

153 The track view option has a standard genome browser layout similar to that of the WashU
154 Epigenome Browser, in which a reference genome sequence is visualized as a sliding window.
155 Various annotation data tracks are hosted on the browser and can be loaded for visualization in
156 a genomic context. For each virus, we downloaded publicly available annotations of the
157 reference genome and converted these annotations into refBed tracks that can be visualized in
158 the genome browser. Likewise, immune epitopes identified in SARS were aligned to the 2019-
159 nCoV reference [Materials and Methods], and a track displaying their coordinates in 2019-nCoV
160 is provided. GC-density tracks were also created for each reference genome, and display the
161 percentage of Gs (Guanines) and Cs (Cytosines) per 5bp window. An entropy track [Materials
162 and Methods] showing the degree of sequence diversity at each position and a diversity track
163 [Materials and Methods] showing the percentage of each of the 4 nucleotides at each position
164 across all strains of the given virus species are also included in the database. In addition to

165 hosting 4 virus species reference genomes, The Virus Genome Browser also supports
166 displaying user-specified genomes provided in FASTA format, as shown in the top left part of
167 Figure 2A, under the browser logo.

168
169 The WashU Virus Browser supports a “zoomed-out” view of the entire viral genome. The
170 zoomed-out view can help the user quickly determine the regions of interest that have high
171 frequencies of variation from the reference (SNV track), and also the regions with high
172 nucleotide diversity among all strains (Shannon tracks) [Figure 2A]. Figure 2A illustrates a
173 genome-level browser view of the 2019-nCoV reference genome and 2 SARS strains, each
174 aligned to the SARS reference genome (AY278488.2 = BJ01, DQ071615.1 = Bat rp3,
175 NC_045512.2 = 2019-nCoV). Sequence variation displayed in density mode [Materials and
176 Methods] shows that the divergence between the 2019-nCoV reference genome (red) and the
177 SARS reference genome is higher than the divergence between the two additional SARS
178 strains (green) and the SARS reference genome. For AY278488.2, the variation from reference
179 is mainly confined to the beginning of the genome, while the remainder of the genome is
180 relatively consistent with the reference. However, for DQ071615.1 (bat-derived), the 5’ end of
181 gene S displays high variation from the reference genome. Likewise, the SARS Shannon track
182 shows that the SARS genome is highly diverse across different strains at gene S.

183
184 Once a region of interest is identified, the standard magnification tool of the browser can be
185 used to quickly zoom into the region [Figure 2A]. Upon zooming in, a genome comparison track
186 can be used to inspect variations from the reference genome, particularly useful for comparing
187 cross-species alignments and viewing structural variations [Figure 2B]. The genome comparison
188 track is adopted from the Epigenome Browser. The top navy-colored horizontal bar represents
189 the reference genome loaded (SARS in the case of Figure 2B) and the bottom purple-colored
190 horizontal bar represents the sequence being aligned to the reference (the 2019-nCoV
191 reference sequence, NC_045512.2, in this case). Insertions and deletions are represented as
192 gaps in either the reference or the query. Matches are represented by black lines linking the 2
193 genomes while mismatches are distinguished by omission of the black bar. When the user
194 hovers over a specific nucleotide, the alignment details around that specific nucleotide are
195 shown.

196
197 Upon further magnification, regions can be inspected on a nucleotide level. Mismatches,
198 insertions, and deletions are color-coded in the SNV tracks and stretches of grey signify
199 positions matching the reference [Figure 2C]. Detailed information, such as inserted
200 nucleotides, is displayed upon clicking. When zoomed into individual nucleotides, as shown in
201 Figure 2C, The diversity bed track shows the percentage of each nucleotide across all strains of
202 SARS at the specific position.

203
204 The versatility of the WashU browser framework makes it possible to adapt the browser to
205 address various questions of interest. Figure 3 demonstrates the utility of using the browser for
206 immune epitope conservation discovery. We recapitulated Zhou et al.’s [16] alignment results of
207 two SARS strains to the reference 2019-nCoV nucleocapsid protein sequence [Figure 3A, 3B].
208 Upon inspection of the region, we could directly observe that many immune epitopes are
209 conserved between SARS and 2019-nCoV [Figure 3C]. The user can identify the amino acid
210 sequence of an epitope by simply clicking the track.

211
212 Encouraged by the high sequence similarity between SARS-CoV and the 2019-nCoV reference
213 strain (NCBI:txid2697049), we mined the list of experimentally identified linear epitopes from T-
214 cell, B-cell and MHC-ligand assays from IEDB [15]. We identified a list of 320 high-confidence

215 linear epitopes [Supplementary Table 2] whose amino acids are identical to predicted translated
216 products from the 2019-nCoV reference strain. These provide a catalogue of epitopes for
217 researchers testing immune targets that can potentially elicit T-cell, B-cell and antibody
218 response to 2019-nCoV.

219
220 We also provide these as an annotated bed track to the reference 2019-nCoV genome. Along
221 with the individual strains' SNV tracks, the epitope tracks can provide a quick, intuitive and
222 visual resource to guide prioritization of experimental resources towards developing diagnostics
223 and therapeutics against 2019-nCoV. The value of our novel SNV tracks will only increase as
224 additional strains are sequenced, helping us better understand the evolving 2019-nCoV genome
225 and prioritize epitopes.

226

227 The Phylogenetic Tree View

228 The second viewing option offered by the WashU Virus Genome Browser is a "tree" format, in
229 which the evolutionary relationships of different viral isolates can be visualized as a
230 phylogenetic tree [17]. When the user navigates to the data page of the browser, and selects
231 "Tree View" [Figure 1], all viral genomes hosted on the browser for the selected virus species
232 are displayed in the form of a right-aligned phylogenetic tree, where solid lines indicate branch
233 lengths [Figure 4]. To the right of the tree is a metadata heatmap displaying strain-specific
234 details such as isolate, isolation source, host, country, and collection date. Additionally, if the
235 user added any individual tracks to their cart from the main page, those selected will display a
236 checkmark to the right, allowing the user to easily see where their strains of interest lie among
237 all other strains.

238

239 In addition to the right-aligned tree view, the browser also supports a more traditional left-
240 aligned linear tree view and a radial view. The left-aligned tree view displays branch lengths
241 indicating relatedness of isolates [Figure 5A]. We noticed that in each virus type, several
242 individual strains maintained high sequence similarity, resulting in several short branch lengths
243 and a long vertical tree. In order to improve visualization, we also created a radial tree view
244 [Figure 5B].

245

246 Discussion

247 Maps help us understand the world around us and navigate it. Moreover, they play a critical role
248 in disaster management during disease outbreaks. Herein, we describe the first genetic
249 mapping, exploration, and visualization tool from the WashU Epigenome Browser team that is
250 specifically dedicated to viral genomes. We provide reference genome maps and genomic
251 datasets related to 4 viral disease outbreaks: SARS (2002-03), MERS (2012), Ebola (2014-16)
252 and the latest nCoV (2019-20). More importantly, we not only present publicly available
253 information in the format of easily accessible data tracks, but also offer a platform with high
254 customizability and flexibility where individual investigators and teams can upload and visualize
255 their own genomic datasets in a plethora of formats. In this report, we have demonstrated using
256 the Virus Browser to 1) quickly and intuitively compare multiple viral genomes and study the
257 viral genome at multiple levels [Figure 2, Figure 4, Figure 5]; and 2) combine viral genome
258 information with other functional genomic information (amino acid sequence and putative

259 immune epitope locations, as shown Figure 3) through multiple track types the browser
260 supports, and identify potential therapeutic targets.

261
262 We expect that the WashU Virus Browser can support research related to the latest novel
263 Coronavirus outbreak of 2019-20, and hope that this tool helps accelerate research to further
264 our understanding of 2019-nCoV and aid in the development of therapeutics. In addition, our
265 platform supports the study of any user-specified viral genome, and can be expanded to other
266 viral research.

267
268 To aid in the battle against this crisis, we are releasing the browser at first moment. The browser
269 is still under active construction and is constantly being updated. General feedback, suggestions
270 for additional tracks, and bug reports may be sent to the WashU Virus Genome Browser team
271 by opening an issue request at [https://github.com/debugpoint136/WashU-Virus-Genome-
272 Browser/issues](https://github.com/debugpoint136/WashU-Virus-Genome-Browser/issues).

273
274
275

276 References

- 277
278
- 279 1. **NCBI GeneBank** [<https://www.ncbi.nlm.nih.gov/genbank/2019-ncov-seqs/>]
 - 280 2. **Novel Coronavirus (2019-nCoV) Situation Report - 17**
281 [[https://www.who.int/docs/default-source/coronaviruse/situation-reports/20200206-
282 sitrep-17-ncov.pdf?sfvrsn=17f0dca_4](https://www.who.int/docs/default-source/coronaviruse/situation-reports/20200206-sitrep-17-ncov.pdf?sfvrsn=17f0dca_4)]
 - 283 3. **International Committee on Taxonomy of Viruses (ICTV)**
284 [<https://talk.ictvonline.org/taxonomy/>]
 - 285 4. Cui J, Li F, Shi ZL: **Origin and evolution of pathogenic coronaviruses.** *Nat Rev Microbiol*
286 2019, **17**(3):181-192.
 - 287 5. Li D, Hsu S, Purushotham D, Sears RL, Wang T: **WashU Epigenome Browser update**
288 **2019.** *Nucleic Acids Res* 2019, **47**(W1):W158-W165.
 - 289 6. Zhou X, Li D, Zhang B, Lowdon RF, Rockweiler NB, Sears RL, Madden PA, Smirnov I,
290 Costello JF, Wang T: **Epigenomic annotation of genetic variants using the Roadmap**
291 **Epigenome Browser.** *Nature biotechnology* 2015, **33**(4):345-346.
 - 292 7. Zhou X, Lowdon RF, Li D, Lawson HA, Madden PA, Costello JF, Wang T: **Exploring long-**
293 **range genome interactions using the WashU Epigenome Browser.** *Nat Methods* 2013,
294 **10**(5):375-376.
 - 295 8. Zhou X, Maricque B, Xie M, Li D, Sundaram V, Martin EA, Koebe BC, Nielsen C, Hirst M,
296 Farnham P *et al*: **The Human Epigenome Browser at Washington University.** *Nat*
297 *Methods* 2011, **8**(12):989-990.
 - 298 9. Myers EW, Miller W: **Optimal alignments in linear space.** *Comput Appl Biosci* 1988,
299 **4**(1):11-17.

- 300 10. Katoh K, Misawa K, Kuma K, Miyata T: **MAFFT: a novel method for rapid multiple**
301 **sequence alignment based on fast Fourier transform.** *Nucleic Acids Res* 2002,
302 **30(14):3059-3066.**
- 303 11. Katoh K, Standley DM: **MAFFT multiple sequence alignment software version 7:**
304 **improvements in performance and usability.** *Mol Biol Evol* 2013, **30(4):772-780.**
- 305 12. Price MN, Dehal PS, Arkin AP: **FastTree: computing large minimum evolution trees with**
306 **profiles instead of a distance matrix.** *Mol Biol Evol* 2009, **26(7):1641-1650.**
- 307 13. Price MN, Dehal PS, Arkin AP: **FastTree 2--approximately maximum-likelihood trees for**
308 **large alignments.** *PLoS One* 2010, **5(3):e9490.**
- 309 14. **Virus Browser Source Code** [[https://github.com/debugpoint136/WashU-Virus-Genome-](https://github.com/debugpoint136/WashU-Virus-Genome-Browser)
310 [Browser](https://github.com/debugpoint136/WashU-Virus-Genome-Browser)]
- 311 15. Vita R, Mahajan S, Overton JA, Dhanda SK, Martini S, Cantrell JR, Wheeler DK, Sette A,
312 Peters B: **The Immune Epitope Database (IEDB): 2018 update.** *Nucleic Acids Res* 2019,
313 **47(D1):D339-D343.**
- 314 16. Zhou P, Yang XL, Wang XG, Hu B, Zhang L, Zhang W, Si HR, Zhu Y, Li B, Huang CL *et al*: **A**
315 **pneumonia outbreak associated with a new coronavirus of probable bat origin.** *Nature*
316 2020.
- 317 17. Shank SD, Weaver S, Kosakovsky Pond SL: **phylotree.js - a JavaScript library for**
318 **application development and interactive data visualization in phylogenetics.** *BMC*
319 *Bioinformatics* 2018, **19(1):276.**
320

321 Acknowledgements

322 We thank doctors, nurses, investigators, and all other people fighting on the front line against
323 this viral outbreak, and we sincerely hope that this tool will aid in this battle.

324 Author Contribution:

325 Conceptualization, T.W. Web development, D.L and D.P. SNV track development, J.F. and C.F.
326 Immune epitope analysis, M.C. Data download, metadata generation and annotation, G.M.
327 Sequence alignments and tree generation, X.Z. Manuscript preparation, J.F, C.F, M.C, G.M,
328 T.W.
329

330 Author Support:

331 J.F. is supported in part by the Siteman Cancer Center Precision Medicine Pathway.
332 X.Z. is supported in part by 5R25DA027995.
333 TW is supported by NIH grants R01HG007175, U24ES026699, U01CA200060,
334 U01HG009391, and U41HG010972, and by the American Cancer Society Research Scholar
335 grant RSG-14-049-01-DMC.
336

337 Figure Captions

338 **Figure 1:** Screenshot of the WashU Virus Genome Browser data page. This view demonstrates
339 several customizable features of the browser, including which genome reference to use, which
340 data tracks to select based on several metadata features, and which browser view to use:
341 “genomic” view (green arrow) or phylogenetic tree view (orange arrow).
342

343 **Figure 2:** Illustration of genomic-level and nucleotide-level track views. A: “zoomed out” track
344 view of the entire genome. 2019-nCoV reference genome (shown in red, NC045512.2) and 2
345 SARS strains (shown in green, DQ071615.1 and AY278488.2) are aligned to the SARS
346 reference genome (NC_004718.3). The box in the top left corner allows users to upload and use
347 any sequence in FASTA format as the reference genome. The shaded vertical bar
348 demonstrates the user’s ability to select a region by mouse for further magnification. B:
349 “Zoomed in” view of the sequence flanking the 5’ end of the S protein. C: A further “zoomed in”
350 view to the level of individual nucleotides. Stretches of grey indicate matching while variations
351 are color coded.
352

353 **Figure 3:** Alignment of the genomic region encoding the nucleocapsid protein. A: 2 SARS
354 strains (DQ071615.1 and AY278488.2) and 5 2019-nCoV strains (MN938384.1, MN975262.1,
355 MN985325.1, MN988668.1, and MN988669.1) are aligned to the 2019-nCoV reference. The
356 region encoding the nucleocapsid protein is shown. Putative SARS immune epitopes [Materials
357 and Methods] are displayed in “density mode”.
358 B: A zoomed-in view of A (orange box), displaying the first 9 amino acids of the reference.
359 Results show a “TCA” insertion in the AY278488.2 alignment between positions 28294 and
360 28295 of the 2019-nCoV reference sequence, which is not present in DQ071615.1. These
361 results are consistent with the results reported in Extended Data Figure 5 of Zhou et al. [16]. C:
362 A zoomed-in view of A (purple box), displaying a region conserved between SARS and 2019-
363 nCoV, overlapping several putative immune epitopes.
364

365 **Figure 4:** Screenshot of a linear, right-aligned tree view displaying all housed 2019-nCoV
366 sequences with accompanying metadata. Solid lines signify distance.
367

368 **Figure 5:** A: Screenshot of a linear, left-aligned phylogenetic tree view, displaying all 2019-
369 nCoV strains hosted by the browser. B: Screenshot of a radial tree view for all 2019-nCoV
370 strains.
371

372

WashU Virus Genome Browser

nCov ▼
Reference

TREE VIEW

DATA

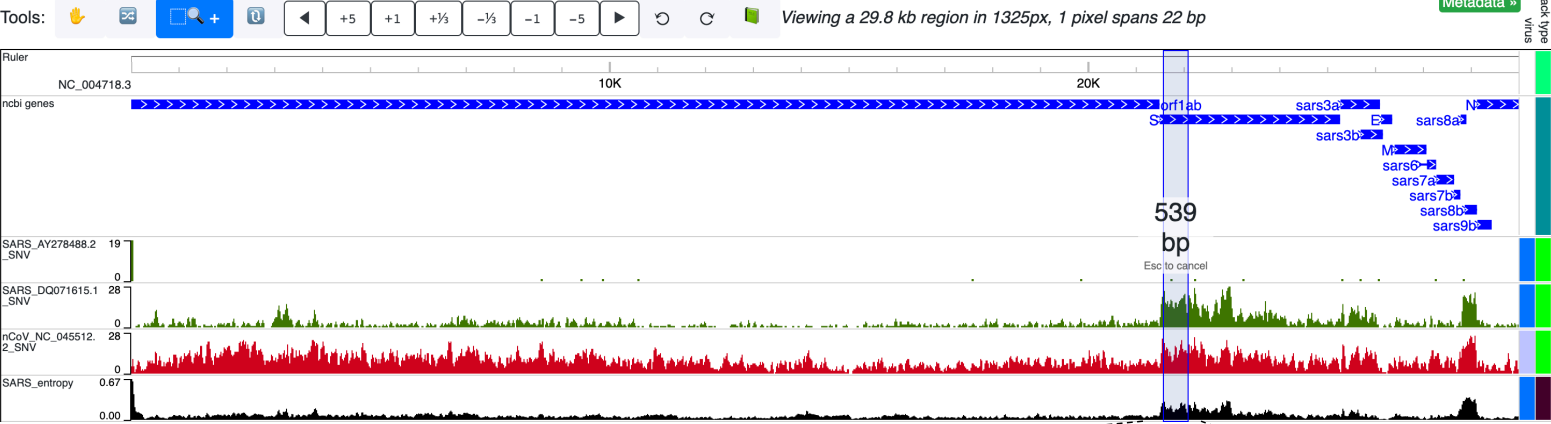
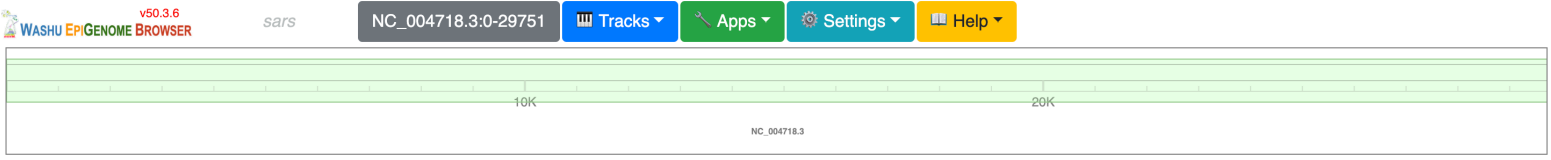
0 FILES

Show Browser View

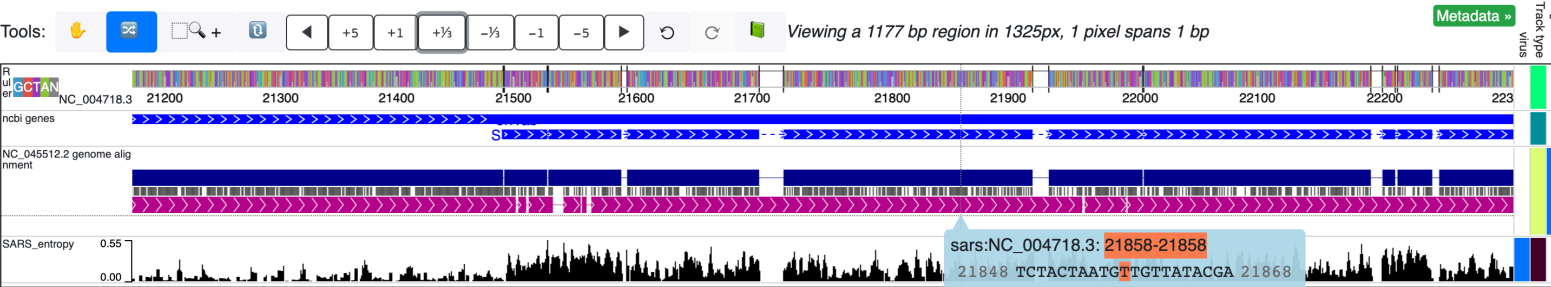
ID	Accession	Isolate	Molecule Type	Country	Collection Date
1	NC_045512.2	Wuhan-Hu-1	genomic RNA	China	Dec-2019
2	MN938384.1	2019-nCoV_HKU-SZ-002a_2020	genomic RNA	China: Shenzhen	Jan-2020
3	MN975262.1	2019-nCoV_HKU-SZ-005b_2020	genomic RNA	China	Jan-2020
4	MN985325.1	2019-nCoV/USA-WA1/2020	genomic RNA	USA	19-Jan-2020
5	MN988713.1	2019-nCoV/USA-IL1/2020	genomic RNA	USA: Illinois	21-Jan-2020
6	MN994467.1	2019-nCoV/USA-CA1/2020	genomic RNA	USA: CA	23-Jan-2020
7	MN994468.1	2019-nCoV/USA-CA2/2020	genomic RNA	USA: CA	22-Jan-2020
8	MN997409.1	2019-nCoV/USA-AZ1/2020	genomic RNA	USA: AZ	22-Jan-2020
9	MN988668.1	2019-nCoV WHU01	genomic RNA	China	02-Jan-2020
10	MN988669.1	2019-nCoV WHU02	genomic RNA	China	02-Jan-2020
11	MN996527.1	WIV02	genomic RNA	China: Wuhan	30-Dec-2019
12	MN996528.1	WIV04	genomic RNA	China: Wuhan	30-Dec-2019
13	MN996529.1	WIV05	genomic RNA	China: Wuhan	30-Dec-2019
14	MN996530.1	WIV06	genomic RNA	China: Wuhan	30-Dec-2019
15	MN996531.1	WIV07	genomic RNA	China: Wuhan	30-Dec-2019
16	MT007544.1	Australia/VIC01/2020	genomic RNA	Australia: Victoria	25-Jan-2020

A

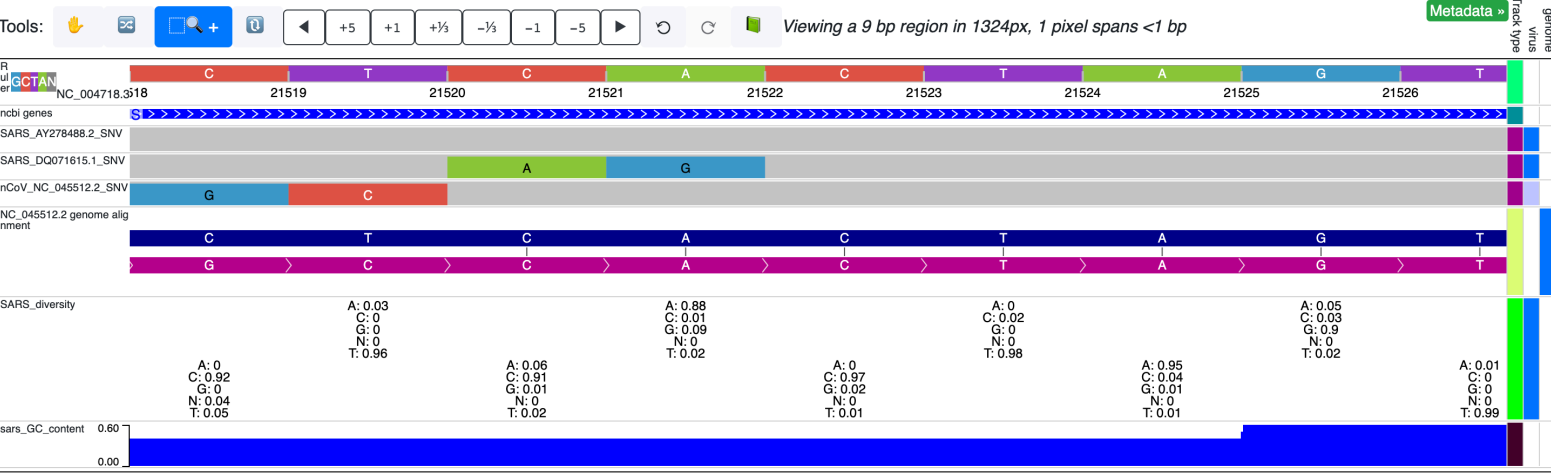
Use custom FASTA file as reference?

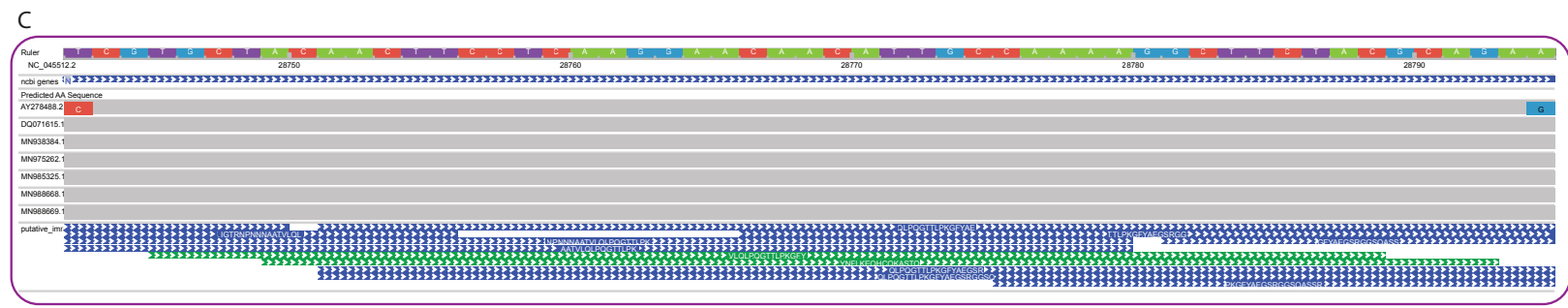
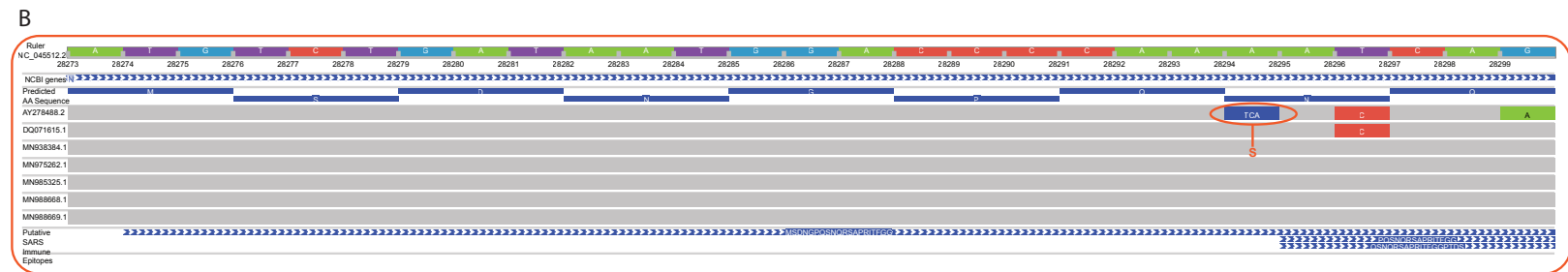
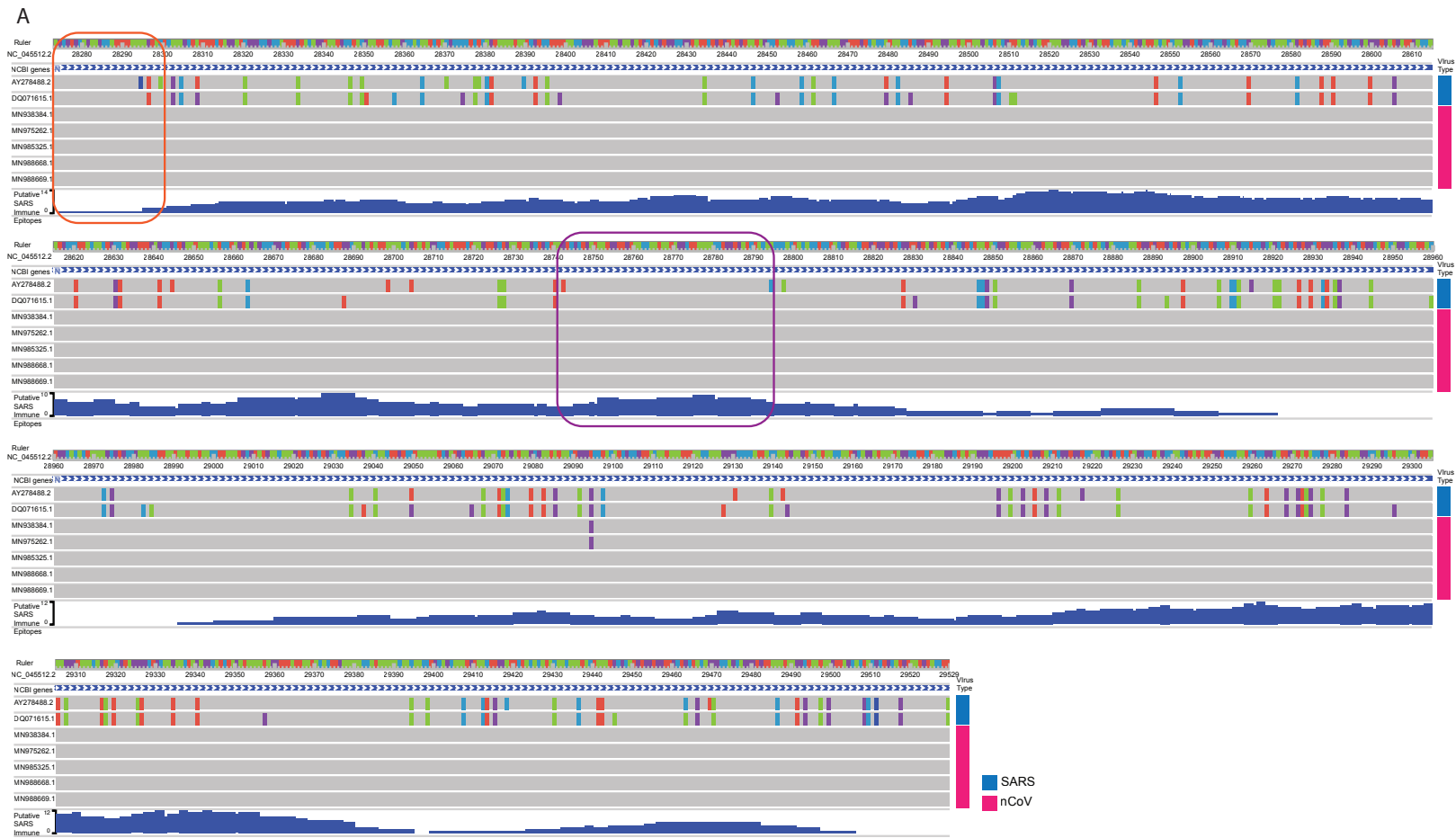


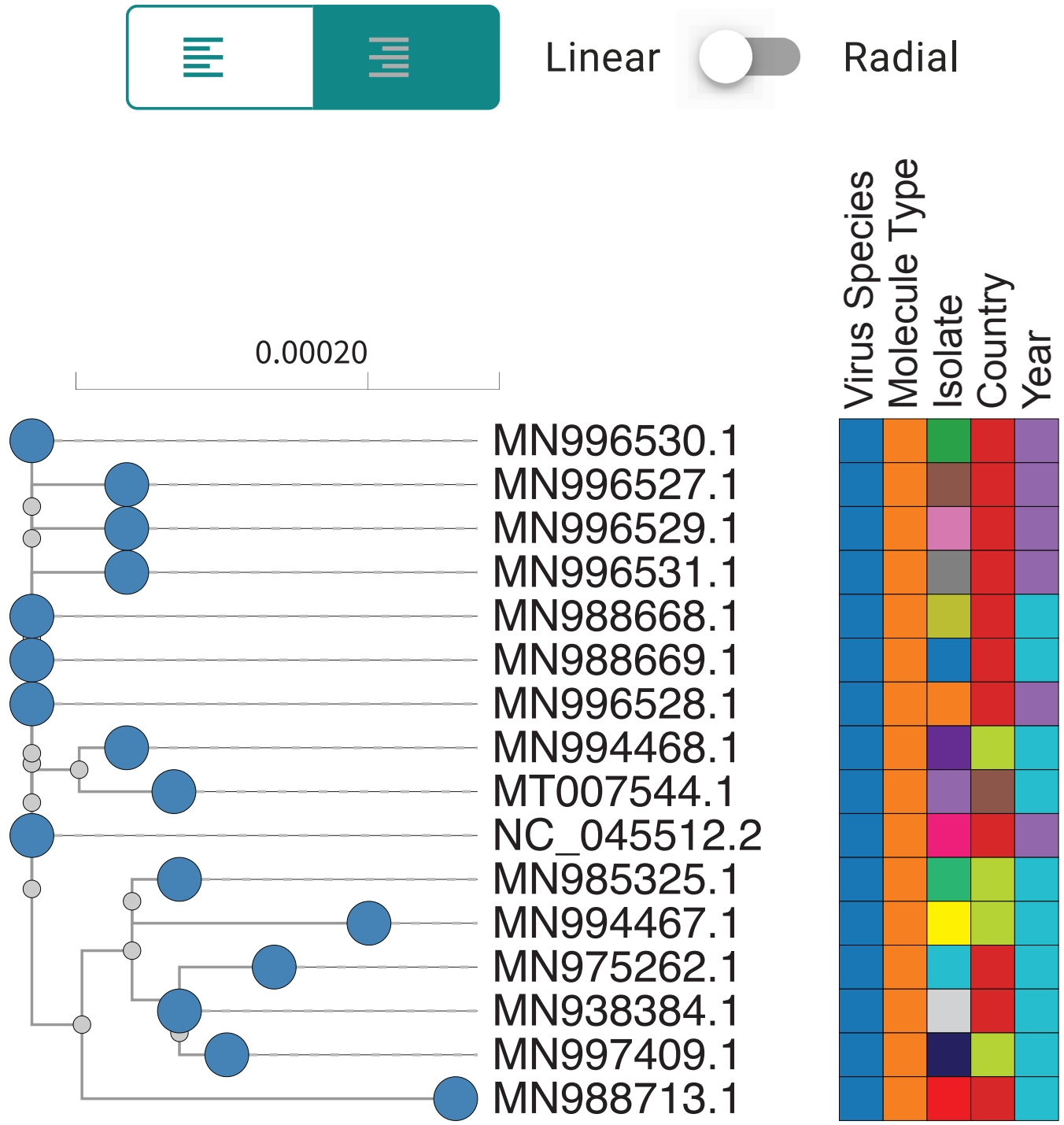
B



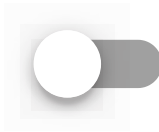
C





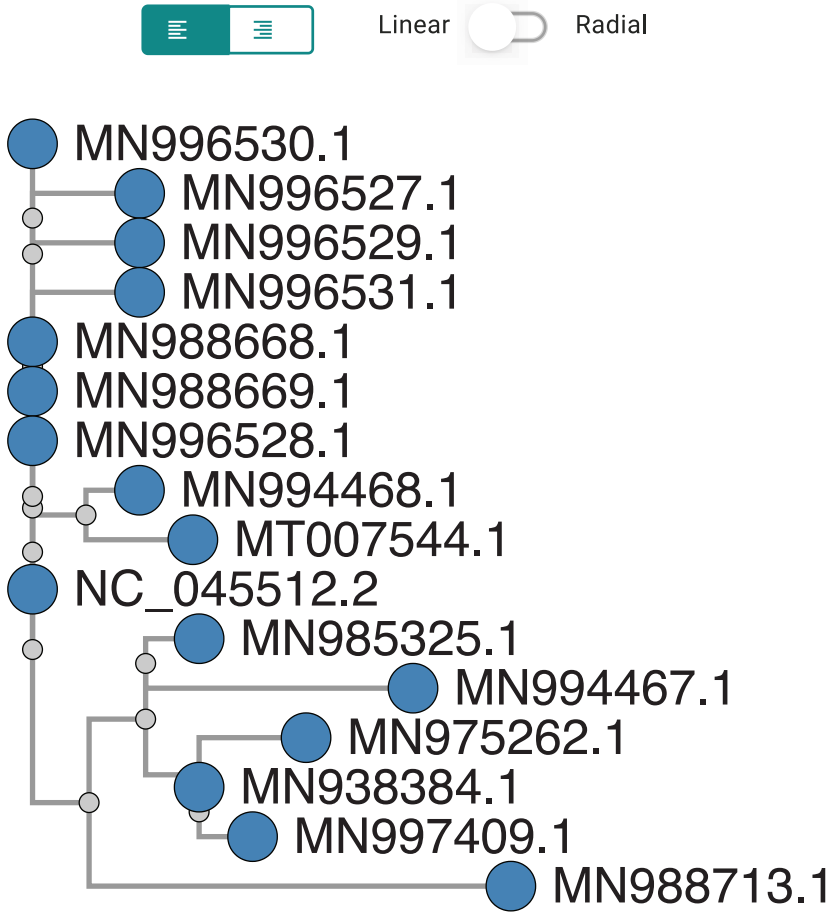


Linear



Radial

A



B

