

1 **Rapid detection of inter-clade recombination in SARS-CoV-2 with**

2 **Bolotie**

3 Ales Varabyou ^{1,2**†}, Christopher Pockrandt ^{1,3**†}, Steven L. Salzberg ^{1,2,3,4}, and Mihaela Pertea
4 ^{1,2,3*}

5 1 Center for Computational Biology, Johns Hopkins University, Baltimore, MD, 21211, USA

6 2 Department of Computer Science, Johns Hopkins University, Baltimore, MD, 21218, USA

7 3 Department of Biomedical Engineering, Johns Hopkins University, Baltimore, MD, 21218, USA

8 4 Department of Biostatistics, Johns Hopkins University, Baltimore, MD, 21205, USA

9 † These authors contributed equally to this work.

10 *Corresponding author.

11 Email: ales.varabyou@jhu.edu, pockrandt@jhu.edu, mpertea@jhu.edu

12 **Running title:** SARS-CoV-2 inter-clade recombination

13 **Keywords:** SARS-CoV-2, COVID-19, recombination, coronavirus

14 **Abstract**

15 The ability to detect recombination in pathogen genomes is crucial to the accuracy of
16 phylogenetic analysis and consequently to forecasting the spread of infectious diseases and to
17 developing therapeutics and public health policies. However, in case of the SARS-CoV-2, the low
18 divergence of near-identical genomes sequenced over a short period of time makes
19 conventional analysis infeasible. Using a novel method, we identified 225 anomalous SARS-CoV-
20 2 genomes of likely recombinant origins out of the first 87,695 genomes to be released, several
21 of which have persisted in the population. Bolotie is specifically designed to perform a rapid
22 search for inter-clade recombination events over extremely large datasets, facilitating analysis
23 of novel isolates in seconds. In cases where raw sequencing data was available, we were able to
24 rule out the possibility that these samples represented co-infections by analyzing the
25 underlying sequence reads. The Bolotie software and other data from our study are available at
26 <https://github.com/salzberg-lab/bolotie>.

27

28 **Introduction**

29 Since the beginning of 2020, the COVID-19 pandemic caused by a newly emerged strain of a
30 betacoronavirus, SARS-CoV-2, has been responsible for over 950,000 deaths and over 30 million
31 infections to date (Dong *et al.* 2020). The strain has been hypothesized to have emerged as a
32 result of a recombination event between strains of betacoronavirus endemic to certain species
33 of bats and pangolins (Zhang *et al.* 2020), although its precise origin is not yet known. To date,
34 the genetic diversity of the SARS-CoV-2 has been increasing slowly compared to other RNA
35 viruses, with 5 to 7 major circulating clades being identified based on multiple variants common
36 to large numbers of isolates in the GISAID database (Shu and McCauley 2017; Hadfield *et al.* 2018).
37 This relative stability of the genetic content of the circulating forms of the virus is promising for
38 the development of vaccines and therapeutics, as well as general understanding of the biology
39 and pathology of SARS-CoV-2.

40 However, as with other RNA viruses, coronaviruses are known to undergo mutations at high
41 rates (Drake and Holland 1999). Inter- and intra-host recombinations are also well-studied and
42 occur frequently (Su *et al.* 2016). As more mutations and lineages of SARS-CoV-2 get fixed in the
43 population, a recombination event caused by a co-infection of a single patient with particles of
44 distinct clades may lead to emergence of novel lineages, posing risks to the efficacy of future
45 treatments. In fact, several accounts of recombination events in SARS-CoV-2 have been
46 reported in recent months (VanInsberghe *et al.* 2020). As such, rapid and consistent
47 surveillance of the sequenced genomes of SARS-CoV-2 for both novel mutations and
48 recombinations is critical to the development of effective treatments and vaccines (Demir *et al.*
49 2020).

50 Multiple computational methods have been developed to detect recombination in microbial
51 genomes and have been used in studies of HIV-1 mutagenesis, bacterial evolution, and other
52 applications (Posada 2002). Some popular methods, such as 3seq, analyze every possible triplet
53 of a set of genomes and statistically evaluate the triplets (Lam *et al.* 2018). Other methods, like
54 PhiPack, are designed to work for low-divergence genomes but still require a significant number
55 of variants (1-5%) to perform statistical analysis (Bruen *et al.* 2006). Another limitation of some of

56 these methods is that the algorithms are computationally intensive and require significant
57 resources and time to perform analysis, particularly for the amount of data (nearly 300,000
58 complete high-coverage genomes to date) that has been generated for SARS-CoV-2. There exist
59 over 4 quadrillion unique triplets of sequences for the currently available SARS-CoV-2 genomes,
60 and a similarity analysis for each triplet would be computationally infeasible. A more efficient
61 approach is necessary if we want to be able to detect recombinants in a realistic amount of
62 time.

63 In this work we present Bolotie, a new algorithm designed to conduct mutational analysis and
64 to detect recombinant forms and other anomalies among a very large set of viral sequences.
65 The methods presented are also designed such that novel sequences can be analyzed efficiently
66 without the need to rerun the entire protocol.

67 We applied Bolotie to search for recombination events in 87,695 complete genomes of SARS-
68 CoV-2 currently available in the GISAID database (Shu and McCauley 2017). In our analysis we
69 identified multiple unique cases of recombination between 4 prominent clades of the virus.
70 Several of the identified recombination events, including some previously reported
71 (VanInsberghe *et al.* 2020), appear in multiple isolates, suggesting transmission in the
72 population.

73 Lastly, we propose a methodology for distinguishing true recombination events from cases of
74 mis-assembly of isolates from a host co-infected with several distinct lineages of a pathogen.
75 The proposed method can be applied by to verify future SARS-CoV-2 genomes prior to database
76 submission.

77 **Materials and Methods**

78 In this work we present Bolotie, a collection of methods that enables rapid alignment, variant
79 calling, inter-clade recombination detection, and parent sequence search for large sets of
80 assembled viral genomes. Our method is robust even when the divergence of collected
81 genomes is very small, and it is designed for ultra-fast detection of anomalies. Because Bolotie
82 utilizes a probability index that can be used for a one-against-all analysis, alignments and

83 indices do not have to be recomputed when evaluating novel sequences as recombinants,
84 which greatly increases the efficiency of single-sequence analysis.

85 **Data**

86 304,811 complete high-coverage genomes were obtained from GISAID (Shu and McCauley
87 2017) using the provided interface. The detailed analysis and evaluation of the software
88 presented in this work is based on a subset of this dataset of 87,695 genomes collected before
89 October 2020. SARS-CoV-2 reference genome isolate Wuhan-Hu-1 (GenBank, accession no.
90 MN908947) was obtained from NCBI and used to guide the alignment, variant calling and
91 consensus sequence generation in the protocol.

92 Clade assignments for all available genomes were obtained from GISAID (Shu and McCauley
93 2017). Lineages S, L, V and O were grouped together, similarly to NextStrain (Hadfield *et al.*
94 2018), and based on the close distances observed in our independent phylogenetic analyses.
95 Mappings between clade IDs in our analysis, those defined by GISAID, and those defined by
96 NextStrain are provided in Supplementary Table 1.

97 To test genomes for the possibility that the putative recombinant isolates might represent a co-
98 infection with two or more distinct SARS-CoV-2 isolates, we searched the SRA database of raw
99 read data using all GISAID and lab-assigned identifiers.

100 **Computing Alignments**

101 A global alignment method to the reference using the implementation from the KSW2 library (Li
102 2018; Suzuki and Kasahara 2018) was developed to facilitate efficient and parallelized variant
103 calling for each query sequence. Because sequence divergence is very low in the collection of
104 SARS-CoV-2 genomes, the reference-guided approach used instead of the conventional multiple
105 sequence alignment is 1) significantly faster, 2) allows simple addition of new sequences, and 3)
106 does not cause explosions in gaps at 3' and 5' ends of the viral genomes. Alignment was
107 performed with a DNAFULL scoring matrix, a gap penalty of 12 and gap extension penalty of 4.

108 **Computing Consensus Sequences**

109 After the pairwise alignment step, we construct for every genome in the input set a consensus
110 sequence by substituting reference genome alleles for the high-frequency variants called by the
111 aligner. Not only does this approach allow us to filter variants used to search for recombinants,
112 but also produces a set of genomes with a standardized set of coordinates which can be used as
113 a multiple sequence alignment (MSA) for phylogenetic analysis.

114 Because the 3' and 5' ends of viral genomes are notoriously difficult to assemble correctly, we
115 chose to force the first and last 200 bases in the consensus sequence of each genome to be
116 identical to the reference genome in accordance with previous studies (VanInsberghe *et al.*
117 2020).

118 Next, since Bolotie is designed to work for genomes with very few variants it is particularly
119 sensitive to ambiguous nucleotides. To avoid biases caused by uncalled bases, any such
120 instances were treated as an unknown base (N). Furthermore, to avoid bias in our predictions
121 for all sequences in the dataset we replace nucleotides at a position of low-frequency variants
122 with the reference allele making such positions equally probable for any clade (clade neutral).
123 We define a low-frequency SNV as one that has fewer than 100 genome sequences that differ
124 from the reference sequence at that position.

125 Because we did not use structural variants when constructing consensus sequences, the final
126 collection of filtered genomes represents a MSA. This not only allows us to use it directly for
127 phylogenetic analysis but also to compute distances among sequences efficiently.

128 **Identifying Recombinants**

129 Based on the provided clade information for each sequence, a model is created to evaluate
130 each sequence as a potential recombinant. From the MSA the conditional probability for each
131 nucleotide position is computed; i.e., the probability of this position belonging to clade C_i given
132 the nucleotide b . However, to account for differences in clade sizes, we multiply the base
133 counts at each position for each clade with the reciprocal of the number of sequences in that
134 clade. The algorithm also ensures that any ambiguous character is assigned a neutral
135 conditional probability of $\frac{1}{|C_i|}$. For every position in the MSA we now have the normalized

136 conditional probabilities of that base position belonging to a certain clade given the base
 137 observed in the consensus sequence.

138 Let $C = c_1, c_2, \dots, c_k$ be the clades (in this paper: $k=4$) with $c_i \in C$. We define two sets of
 139 bases $B_4 = \{A, C, G, T\}$ and B_{16} being all IUPAC characters. $n_{c_i,b}$ denotes the number of
 140 sequences in clade c_i that have base $b \in B_4$ at the position of interest. Now the weighted
 141 conditional probability $Pr(c_i|b)$ of the sequence belonging to clade c_i given the observed base
 142 $b \in B_{16}$ at a certain position is defined as:

$$143 \quad Pr(c_i | b) = \begin{cases} \frac{n_{c_i,b} * w_{c_i}}{\sum_{c_j \in C} n_{c_j,b} * w_{c_j}}, & b \in B_4 \\ \frac{1}{|C|}, & \text{otherwise (i. e., } b \text{ is ambiguous)} \end{cases} \quad (1)$$

$$144 \quad w_{c_i} = \frac{1}{\sum_{b \in B_4} n_{c_i,b}} \quad (2)$$

145 Now we determine for each sequence whether it might be a recombinant. To do this we need
 146 to determine for each base the most likely clade. This can be modeled as an HMM and solved
 147 efficiently (in both time and space) with the Viterbi algorithm in $O(nc^2)$ time, where n is the
 148 sequence length and c is the number of clusters (Forney 1973). The Viterbi algorithm
 149 determines the most likely sequences of clades that should be assigned to the positions in the
 150 genome, by finding the sequence of assignments that maximizes the posterior probability
 151 across the space of all possible clade assignments. This probability depends on the prior
 152 probabilities computed for each position to be in a clade, and the probability of switching from
 153 one clade to another. We assigned a very small value to this probability (0.0001, shown in
 154 Figure 1), on the assumption that a recombinant is an unlikely event. We empirically tried
 155 several larger and smaller values and found that they made almost no difference in the results.

156 At each position of the MSA, each state is representing one clade with the conditional
 157 probabilities for every four bases. The Viterbi algorithm then finds the path with the highest
 158 likelihood given the sequence of an isolate. The model is identical for every isolate and only
 159 depends on the MSA. Figure 1 shows the approach, but only illustrates the conditional
 160 probability of the observed base of the sequence in each state.

161 This method cannot detect the exact breakpoint location for a recombinant, because it relies on
162 discrete SNV differences between clades; thus, it can only narrow down the breakpoint to the
163 region between two clade-specific SNVs.

164 **Searching for Closest Parental Genomes**

165 Previous methods for recombination detection examine triplets of sequences to detect a
166 potential recombinant genome and its parents. Unlike such methods, our algorithm relies on
167 conditional probabilities computed for clades – a strategy that provides additional statistical
168 power and also reduces the complexity of the problem.

169 To identify potential parents of recombinant sequences we have implemented a dedicated
170 method. For each recombined segment of the sequence, as inferred by our algorithm, we
171 compute a Kimura distance matrix (Kimura 1980) to other sequences in the clade of the
172 corresponding segment. Kimura distance is a distance metric that scores transitions (A <-> G
173 and C <-> T) differently than transversions (interchange of purine for pyrimidine bases). To be
174 more “accurate” we also include low-frequency SNVs that we neglected in the previous step in
175 the consensus sequences. All sequences with the lowest distance score are reported as most
176 likely parents that have contributed the segment in the recombination event.

177 **Investigating sequencing data**

178 To test our recombination candidates for signs of co-infection, we aligned available reads with
179 Bowtie2 (Langmead and Salzberg 2012) against the Wuhan-Hu-1 reference genome using the “-
180 -very-sensitive-local” option and otherwise default parameters. The mappings were further
181 sorted and indexed using samtools (Li *et al.* 2009). Counts for individual nucleotides were
182 obtained using bam-readcount software (<https://github.com/genome/bam-readcount>) and
183 positions with high conditional probabilities were extracted and summarized.

184 **Phylogenetic analysis**

185 To test how well information is preserved in our consensus sequences, we obtained a pre-
186 computed tree from NextStrain (Hadfield *et al.* 2018) which contained a total of 4,494
187 representative genomes chosen by NextStrain. Only 4,039 genomes from those in the

188 NextStrain tree were available on GISAID at the time when we obtained genome assemblies for
189 our analysis.

190 First, we re-built the tree using the set of 4,039 genomes using the general time-reversible
191 (GTR) model as used by the NextStrain platform and allowing IQ-TREE (Minh *et al.* 2020) to
192 automatically choose the precise model using its ModelFinder package (Kalyaanamoorthy *et al.*
193 2017). The final tree was generated under the GTR model with 1000 rounds of bootstrapping.
194 The same approach was taken to build the phylogenetic tree that included all identified
195 anomalous sequences.

196 **Results**

197 We first aligned all genomes to the Wuhan-Hu-1 reference genome (GenBank accession
198 [MN908947](#)), from which we detected 84,322 single-nucleotide variants (SNVs) at 29,503 sites.
199 After removing all variants that appear in fewer than 100 sequences, we retained a set of 659
200 SNVs at 411 unique sites. Alignments also revealed 1,349 unique structural variants (934
201 deletions and 415 insertions). While 2 deletions and 1 insertion were present in over 100
202 genomes, for the purpose of computational efficiency we did not consider them further.

203 Using Bolotie on the set of well-supported variants to search for recombination events between
204 sequences in the 4 major clades of SARS-CoV-2, we identified 225 possibly recombinant
205 genomes. Figure 2 illustrates that many of the identified recombination events were
206 represented by a single genome. However, several lineages with near-identical sequences were
207 observed. In Figure 2 these lineages appear as broad red bands with a high density of outgoing
208 arcs. Additionally, several smaller groups of near-identical recombinant signatures have been
209 observed in which genomes differed by one or two variants. Those genomes were often found
210 to be neighbors in the computed maximum likelihood (ML) trees and often had the same or
211 neighboring inferred parental sequences.

212 Of the 225 recombinant genomes, 109 were labeled in the original GISAID data as belonging to
213 clade #0, 111 in clade #1 and 5 in clade #2 (Figure 2). Recombination events happened between
214 members of all 4 clades, with 171 parental genomes identified in clade #0, 41 parental genomes
215 in clade #1, 148 in clade #2 and 90 in clade #3. Additionally, 15 out of 225 potential

216 recombinants were found present in the set of 4,039 representative genomes used by
217 NextStrain (Supplementary Table 2).

218 Of the 225 identified recombinants, most of the recombinant signatures had 1 or 2 breakpoints
219 like the ones shown in Figure 3A, 3B and 3C. However, at least 6 genomes including the one
220 depicted in Figure 3D exhibited more complex patterns of mosaicism with 3 breakpoints.

221 While in majority of cases the path-finding algorithm of Bolotie relied on clade-defining variants
222 with high conditional probabilities > 0.9 , several positions exhibited an inverse pattern and
223 were also useful in the analysis. For example, in Figure 3A (further detailed in Table 1),
224 mutations of cytosine (C) to thymine (T) at position 14,407 and adenine (A) to guanine (G) at
225 position 23,402 are not characteristic of clade 1 (yellow) since they have the same conditional
226 probability of ~ 0.33 of defining clades 2 and 3. However, these positions are informative in an
227 inverse way, namely that observing a C and an A at those positions is very unlikely if the
228 sequence originated in clade 0. This adds additional evidence to the recombinant origins of the
229 genomes.

230 Even though most anomalous sequences reported by Bolotie had clean separations between
231 two parental clades, some mutational signatures contained admixtures in conditional
232 probabilities from other clades. As illustrated in Figure 3C, at position 240 the conditional
233 probability has a greater affinity towards a blue clade. A signature like that could be indicative
234 of a random mutation, sequencing or assembly artifacts, and Bolotie resolved the parental
235 clades in a seemingly parsimonious way.

236 Finally, we have applied the same analysis to the most recent set of genomes available through
237 GISAID. As of February 2021 a total of 304,811 complete high-coverage genomes were available
238 on GISAID with classifications into 5 clades. We analyzed new data with Bolotie using the same
239 set of parameters, except for increasing the minimum allele frequency from 100 to 500
240 genomes to account for the larger total dataset size. Despite identifying 775 anomalous
241 sequences in the data, the rate at which anomalous sequences are being sequenced and
242 deposited on GISAID remains constant (Supplementary Figure 2, Supplementary Tables 4,6) and
243 no widely spread clusters of recombinant genomes were identified.

244 **Sequencing and assembly artifacts**

245 Although recombination has been extensively observed within and between members of the
246 coronaviridae family, such observations were in the past characterized based on years of
247 accumulated variants in a relatively small collection of genomes. The bulk of the SARS-CoV-2
248 sequences currently deposited in GISAID and GENBANK were collected and sequenced between
249 the months of March and May of 2020 and show very high degree of similarity, which
250 combined with the large number of available genomes, makes evolutionary analysis very
251 challenging.

252 Another complication is that an apparent recombinant strain might instead be the result of a
253 co-infection. Suppose a sample was collected from a patient co-infected with two distinct
254 lineages of the virus, where one lineage contains two SNVs while the other lineage contains two
255 distinct SNVs at the same positions. Such a sample would be amplified and sequenced in a
256 single batch. Reads representing both alleles would be provided to an assembler to produce the
257 final genome. Due to differences in coverage of each allele, it is possible that an assembly
258 algorithm would produce a sequence with alleles from both clades, creating an artifactual
259 recombinant. Re-analysis of raw sequencing reads by mapping them against the genome should
260 reveal such artifacts, because reads from a patient co-infected by multiple clades would reveal
261 both alleles at the corresponding sites.

262 To evaluate this hypothesis, we obtained sets of raw sequence reads deposited at NCBI/SRA or
263 ENA for some of the recombinant genomes identified with Bolotie. Unfortunately, GISAID does
264 not require authors to submit raw data, and only a limited number of submitters have placed
265 their data in public archives with corresponding GISAID identifiers. Thus, we were only able to
266 recover a limited number of datasets for our analysis.

267 Reads for the EPI_ISL_439137 isolate (Figure 3A) were obtained from the European
268 Bioinformatics Institute's ENA database. As summarized in Table 1, all positions with high
269 conditional affinities for a clade had a homogenous composition, indicating that the data did
270 not derive from two distinct isolates, but instead it was likely a single isolate containing variants
271 from two parental lineages.

272 Searching the list of recombinant genomes for similar signatures, we identified another isolate,
273 EPI_ISL_489588 also from Scotland, dated one week earlier, which contained the same variants.
274 Another isolate, EPI_ISL_510303 from Spain, had all but one variant (at position 28,143)
275 matching the recombinant signature of the isolates from Scotland. Given the rapid mutation
276 rate in RNA viruses, it is possible that an independent mutation occurred at that position, or
277 that the reference allele is an assembly artifact, however we could not find raw data
278 corresponding to the Spanish isolate and were unable to further investigate possible reasons
279 for the missing variant.

280 **Estimation of false positives**

281 To provide a simple estimate of the false discovery rate of Bolotie, we generated a large set of
282 simulated SARS-CoV-2 genomes that contained no recombination events. First, we computed a
283 consensus genome for each of the 4 main clades shown in Figure 2. Then for each of the four
284 genomes we generated 25,000 descendants by incorporating minor alleles at the frequencies
285 found in the data, as well as random mutations at the rate of 6×10^{-4} , as reported by others (van
286 Dorp *et al.* 2020). We then ran Bolotie on these 100,000 genomes, and found only 4 false
287 positive results. Further details can be found in the Supplementary Methods.

288 **Performance**

289 Complete analysis of the 87,695 genomes using Bolotie including alignment and index
290 construction took a total of ~5.5 hours using 36 threads on two Intel Xeon E5-2680 v2
291 processors with 10 cores each. Using the conditional probability table provided with the
292 software, analyzing a single additional genome takes on average ~30 seconds.

293 **Discussion**

294 The method and experiments presented in this work demonstrate that recombination has
295 occurred between the four existing major clades of SARS-CoV-2. While some of the inferred
296 events may be homoplasies or technical artifacts, our analysis shows that at least some of the
297 genomes likely represent true cases of recombination.

298 Of the 225 recombination events identified in our analysis, the majority were represented by
299 single isolates, suggesting that the event was not established in the population. However,
300 because two-thirds of the available genomes were sequenced between late March and early
301 May, it is possible that more data will reveal additional recombinant lineages.

302 The 225 inferred recombinant genomes comprise less than 1% of all sequences analyzed. It is
303 possible that many more anomalous genomes could be detected by lowering the variant
304 frequency threshold. For example, if we require 50 sequences to confirm a variant rather than
305 100, the number of informative sites increases more than two-fold from 411 to 996, possibly
306 allowing detection of events that are rarer, such as those which involve smaller emerging
307 lineages within the 4 clades. However, due to decreased specificity such an approach might
308 require stricter manual inspection as the false positive rate is expected to increase
309 substantially.

310 While overall all events detected by Bolotie passed manual verification, a possibility of mis-
311 assembly in cases of co-infection by particles from different clades could also explain the
312 presence of multiple clade-specific variants within a single genome. Although we were unable
313 to obtain raw read data for all recombinants, our analysis of the EPI_ISL_439137 isolate (Figure
314 3A, Table 1) shows that at least in one case the recombinant origins of the genome can be
315 validated. Several other recombinants identified in our analysis (EPI_ISL_468407,
316 EPI_ISL_452334, EPI_ISL_475584, EPI_ISL_464547) have also been previously identified by
317 other groups (VanInsberghe *et al.* 2020).

318 Due to the differences in library preparation, sequencing technology and assembly protocols,
319 the need for raw data and independent validation is very high. We urge researchers to submit
320 raw sequencing data so that any future studies can verify their findings, not only when studying
321 recombination events, but also individual rare variants, transmission patterns, clade prevalence
322 in different populations, etc.

323 Because our method relies heavily on the accuracy of variant calls, we sought to compare how
324 well available phylogenetic trees agree with trees built using consensus sequences constructed
325 from the alignment data we obtained. The trees shown in figures 4A and 4B are very similar,

326 confirming that consensus sequences constructed by Bolotie preserve essential information
327 sufficient for accurate phylogenetic analysis. In our analysis, Bolotie identified 15 out of 4,039
328 sequences used by NextStrain (Hadfield *et al.* 2018) as recombinants (Supplementary Table 2).
329 Minor differences between the NextStrain tree and the tree computed from Bolotie consensus
330 sequences are to be expected since consensus sequences have 200 bases replaced with the
331 reference at both 3' and 5' ends and do not include any structural variants. Additionally, since
332 NextStrain tree includes 455 isolates submitted to GISAID after we downloaded our latest set,
333 those sequences are also expected to slightly alter the topology. Lastly, differences in software
334 versions and randomized methods inherent in the tree-building software are expected to
335 produce trees with minor differences on each iteration.

336 On the other hand, the introduction of 225 recombinant or otherwise anomalous genomes
337 produced a mildly distorted tree with multiple outliers, shown in Figure 4C and Figure 2. In both
338 illustrations most of the recombinant genomes were assigned a clade different from the non-
339 recombinant neighbors by GISAID. Furthermore, several groups of potential recombinant
340 genomes with identical or highly similar mutational signatures were identified by Bolotie
341 (Figure 2). Such lineages are especially important for phylogenetic analysis, as they may affect
342 the topology of the trees more significantly. Even minor perturbations to the topology of the
343 tree in the presence of misclassified outliers may have adverse effects on the studies of
344 dynamics and transmission of the pathogen lineages in the population (Awadalla 2003). This
345 once again illustrates the importance of properly handling anomalous sequences in
346 phylogenetic analysis. It must also be noted that different clade assignments for the SARS-CoV-
347 2 genomes currently exist (Shu and McCauley 2017; Hadfield *et al.* 2018), at least in part due to
348 differences in tree-building strategies (Supplementary Table 1). However, even though
349 discrepancies in clade assignment may present a challenge to Bolotie, the results will still be of
350 use for the refinement of phylogenetic trees and ultimately clade assignments of the genomes.
351 Grouping together conflicting clades and smaller clades should increase the specificity of the
352 results.

353 To the best of our knowledge, there currently exists one other recent evaluation of
354 recombination in SARS-CoV-2 genomes (VanInsberghe *et al.* 2020). Several studies hinted at the
355 possibility of recombination occurring, but the results were inconclusive due to very small
356 sample size at the time (Yi 2020). Bolotie adds 221 candidate recombinants to the 5 proposed
357 by VanInsberghe *et al.* Upon comparison of our results we noted that isolate EPI_ISL_464547
358 was reported by VanInsberghe *et al.* but was absent in the output of Bolotie. Of the 5 genomes
359 reported in that study, EPI_ISL_464547 had the shortest length of the second segment. Upon
360 closer investigation of conditional probabilities computed by Bolotie for that genome, we found
361 that none of the variants are clade-defining, nor does the dominant clade 0 ever drop in
362 probability below the baseline of 0.25 at which all genomes are equally likely. The only
363 exception was a variant at the tail end of the sequence, which was equally likely for clades 0, 1
364 and 3 at 0.33 probability (Supplementary Figure 1). Thus, we concluded that the assignment
365 done by Bolotie was likely correct for that genome.

366 However, it is possible that differences in partitioning of clades 0, 1 and 3 between our method
367 and that used by VanInsberghe *et al.* would result in different conditional probabilities at that
368 variant position. As shown in Supplementary Table 1, one such discrepancy does exist between
369 clades 1 and 3, where NextStrain classification adds 334 sequences from clade GH to clade G.
370 However, clade 0 would still have a conditional probability at that position greater than the
371 baseline. As a result, we find it difficult to provide a conclusive assessment of the
372 EPI_ISL_464547 isolate given currently available data.

373 Our work focused primarily on the development of a highly efficient, scalable general-purpose
374 method for detecting recombination events in viral genomes irrespective of the divergence
375 rates in the pool of collected isolates. An additional purpose was to search for convincing
376 evidence that recombination does indeed happen in SARS-CoV-2. Hence, we imposed multiple
377 conservative criteria, such as assessment of only 4 clades. However, it is important to note that
378 recombination likely happens within lineages of the same clade. While not targeted in our
379 analysis, future studies may choose to evaluate these intra-clade events, possibly yielding much
380 higher numbers of recombination events.

381 **Conclusion**

382 Given how recently this novel strain of coronavirus appeared, much remains to be learned
383 about SARS-CoV-2 and how it may change over time. Recombination events, which may have
384 been responsible for the initial emergence of SARS-CoV-2 (Zhang *et al.* 2020), may have
385 significant impact on future transmission and virulence of the virus. As such, our ability to
386 detect recombination events in a timely manner is crucial in the ongoing efforts to find a
387 solution to the pandemic and prevent additional casualties.

388 Utilizing an enormous collection of SARS-CoV-2 isolates sequenced by thousands of researchers
389 around the globe (Shu and McCauley 2017), we were able to develop a method that can reliably
390 detect sequences with anomalous mutation patterns which are indicative of recombination
391 events. Using the proposed method, we identified 225 high likelihood recombinant sequences.
392 Our findings suggest that recombination in SARS-CoV-2 is much more common than previously
393 reported and that several recombinant lineages may have become established in the
394 population.

395 We hope that the software presented here along with provided pre-built indices will help to
396 detect future recombination events quickly and reliably, and aid in efforts to track the spread of
397 the SARS-CoV-2 virus.

398 **Data Availability Statement**

399 The core method is implemented in C++ and based on the SeqAn (Reinert *et al.* 2017) and KSW2
400 (Li 2018; Suzuki and Kasahara 2018) libraries, the tree building is performed by IQ-TREE (Minh
401 *et al.* 2020). The code and test data are available for download on GitHub:

402 <https://github.com/salzberg-lab/bolotie>. The SARS-CoV-2 index built using the genomes in our
403 analyses is also available for download at ftp://ftp.ccb.jhu.edu/pub/data/bolotie_sars_cov_2/.

404 A wrapper script is provided in the GitHub repository to run all steps of the protocol. This script,
405 while convenient, is intended for replicability and testing and lacks some of the available
406 features of Bolotie.

407 Supplemental Material available at figshare: <https://doi.org/10.25386/genetics.14553696>

408 **Acknowledgements**

409 We acknowledge the work of all the authors, originating and submitting laboratories who
410 contributed nucleotide sequences to the GISAID EpiCov Database (87,695 isolates, 2 September
411 2020). We would also like to thank Dr. Martin Steinegger for the helpful discussions on
412 algorithms for clustering large collections of viral genomes.

413 **Funding**

414 This work was supported in part by Fast Grants (part of Emergent Ventures at George Mason
415 University) and by the US National Institutes of Health [grants R01-AI141009 and R35-
416 GM130151].

417 **Author Contributions**

418 A.V. and C.P. developed and implemented the methods, ran experiments. A.V., C.P., S.L.S and
419 M.P. conceptualized the study, methods and wrote the manuscript.

420 **Competing Interests**

421 The authors have no conflicts of interest to declare.

422 **References**

- 423 Awadalla P., 2003 The evolutionary genomics of pathogen recombination. *Nat. Rev. Genet.* 4:
424 50–60.
- 425 Bruen T. C., H. Philippe, and D. Bryant, 2006 A simple and robust statistical test for detecting
426 the presence of recombination. *Genetics* 172: 2665–2681.
- 427 Demir A. B., D. Benvenuto, Y. H. ABACIOĞLU, S. Angeletti, and M. Ciccozzi, 2020 Identification
428 of the nucleotide substitutions in 62 SARS-CoV-2 sequences from Turkey. *Turkish J. Biol.*
429 44: 178–184.
- 430 Dong E., H. Du, and L. Gardner, 2020 An interactive web-based dashboard to track COVID-19 in
431 real time. *Lancet Infect. Dis.* 20: 533–534.
- 432 Dorp L. van, D. Richard, C. C. S. Tan, L. P. Shaw, M. Acman, *et al.*, 2020 No evidence for

433 increased transmissibility from recurrent mutations in SARS-CoV-2. *Nat. Commun.* 11: 1–8.
434 Drake J. W., and J. J. Holland, 1999 Mutation rates among RNA viruses. *Proc. Natl. Acad. Sci.* 96:
435 13910–13913.
436 Forney G. D., 1973 The viterbi algorithm. *Proc. IEEE* 61: 268–278.
437 Hadfield J., C. Megill, S. M. Bell, J. Huddleston, B. Potter, *et al.*, 2018 Nextstrain: real-time
438 tracking of pathogen evolution. *Bioinformatics* 34: 4121–4123.
439 Kalyanamoorthy S., B. Q. Minh, T. K. F. Wong, A. Von Haeseler, and L. S. Jermini, 2017
440 ModelFinder: fast model selection for accurate phylogenetic estimates. *Nat. Methods* 14:
441 587–589.
442 Kimura M., 1980 A simple method for estimating evolutionary rates of base substitutions
443 through comparative studies of nucleotide sequences. *J. Mol. Evol.* 16: 111–120.
444 Lam H. M., O. Ratmann, and M. F. Boni, 2018 Improved algorithmic complexity for the 3SEQ
445 recombination detection algorithm. *Mol. Biol. Evol.* 35: 247–251.
446 Langmead B., and S. L. Salzberg, 2012 Fast gapped-read alignment with Bowtie 2. *Nat. Methods*
447 9: 357.
448 Li H., B. Handsaker, A. Wysoker, T. Fennell, J. Ruan, *et al.*, 2009 The sequence alignment/map
449 format and SAMtools. *Bioinformatics* 25: 2078–2079.
450 Li H., 2018 Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* 34: 3094–
451 3100.
452 Minh B. Q., H. A. Schmidt, O. Chernomor, D. Schrempf, M. D. Woodhams, *et al.*, 2020 IQ-TREE 2:
453 New models and efficient methods for phylogenetic inference in the genomic era. *Mol.*
454 *Biol. Evol.* 37: 1530–1534.
455 Posada D., 2002 Evaluation of methods for detecting recombination from DNA sequences:
456 empirical data. *Mol. Biol. Evol.* 19: 708–717.
457 Reinert K., T. H. Dadi, M. Ehrhardt, H. Hauswedell, S. Mehringer, *et al.*, 2017 The SeqAn C++
458 template library for efficient sequence analysis: A resource for programmers. *J. Biotechnol.*
459 261: 157–168.
460 Shu Y., and J. McCauley, 2017 GISAID: Global initiative on sharing all influenza data—from vision
461 to reality. *Eurosurveillance* 22: 30494.

462 Su S., G. Wong, W. Shi, J. Liu, A. C. K. Lai, *et al.*, 2016 Epidemiology, genetic recombination, and
463 pathogenesis of coronaviruses. Trends Microbiol. 24: 490–502.
464 Suzuki H., and M. Kasahara, 2018 Introducing difference recurrence relations for faster semi-
465 global alignment of long sequences. BMC Bioinformatics 19: 33–47.
466 VanInsberghe D., A. S. Neish, A. C. Lowen, and K. Koelle, 2020 Identification of SARS-CoV-2
467 recombinant genomes. bioRxiv.
468 Yi H., 2020 2019 novel coronavirus is undergoing active recombination. Clin. Infect. Dis.
469 Zhang T., Q. Wu, and Z. Zhang, 2020 Probable pangolin origin of SARS-CoV-2 associated with the
470 COVID-19 outbreak. Curr. Biol.

471
472

473 **Figure 1.** *The maximum conditional probability for each nucleotide is highlighted in gray, while*
474 *the path with the maximum likelihood is highlighted in bold. By penalizing switching of clades,*
475 *insignificant differences in probabilities between clades as well as short windows representing a*
476 *switch to a different clade are avoided. For clarity transitions between nodes on non-optimal*
477 *paths are indicated in gray without labeled probabilities.*

478 **Figure 2.** *An unrooted topological cladogram of 4,249 SARS-CoV-2 genomes including 225*
479 *recombinants labeled as red bars. Arcs link each recombinant to both inferred parental*
480 *genomes. The color of the arc corresponds to the color of the clade to which a recombinant was*
481 *clustered within the tree. Clades correspond to the GISAID clades GR (0), GH (1), G (2) and all*
482 *minor lineages combined (3).*

483 **Figure 3.** *Four examples of inferred recombinant sequences: A. EPI_ISL_439137; B.*
484 *EPI_ISL_468407; C. EPI_ISL_509874; D. EPI_ISL_417420. The top section of each plot shows*
485 *conditional probabilities of a clade given a nucleotide at each position. Bars are plotted for the*
486 *two parent clades and the other clades are shown as dots of the corresponding color. Each peak*
487 *>0.1 above the baseline (0.25) is labeled with the number of genomes it appears in. An average*
488 *is reported whenever there are multiple variants in close proximity on the plot, listing the*
489 *number of averaged variants in parentheses. The three lower panels of each plot show the*

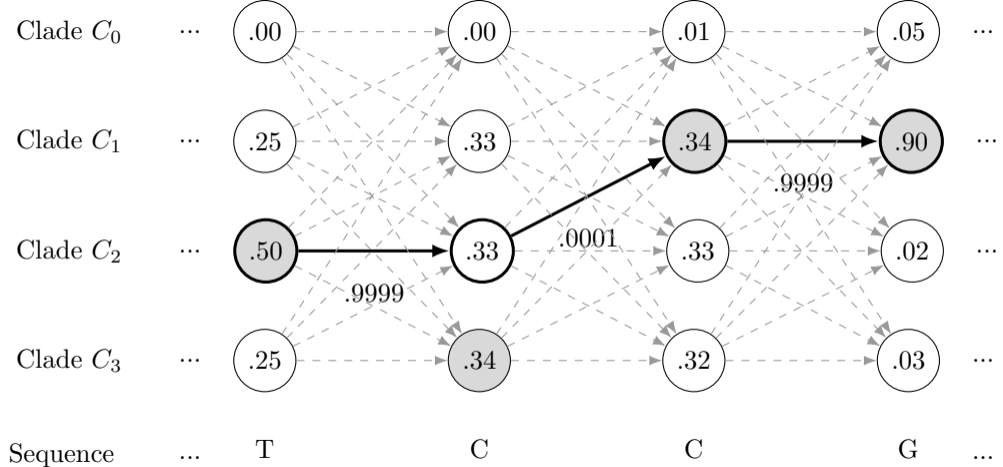
490 frequency of variants at each position for parental clades (top and bottom rows) and variants
 491 observed on the recombinant genome (middle row).

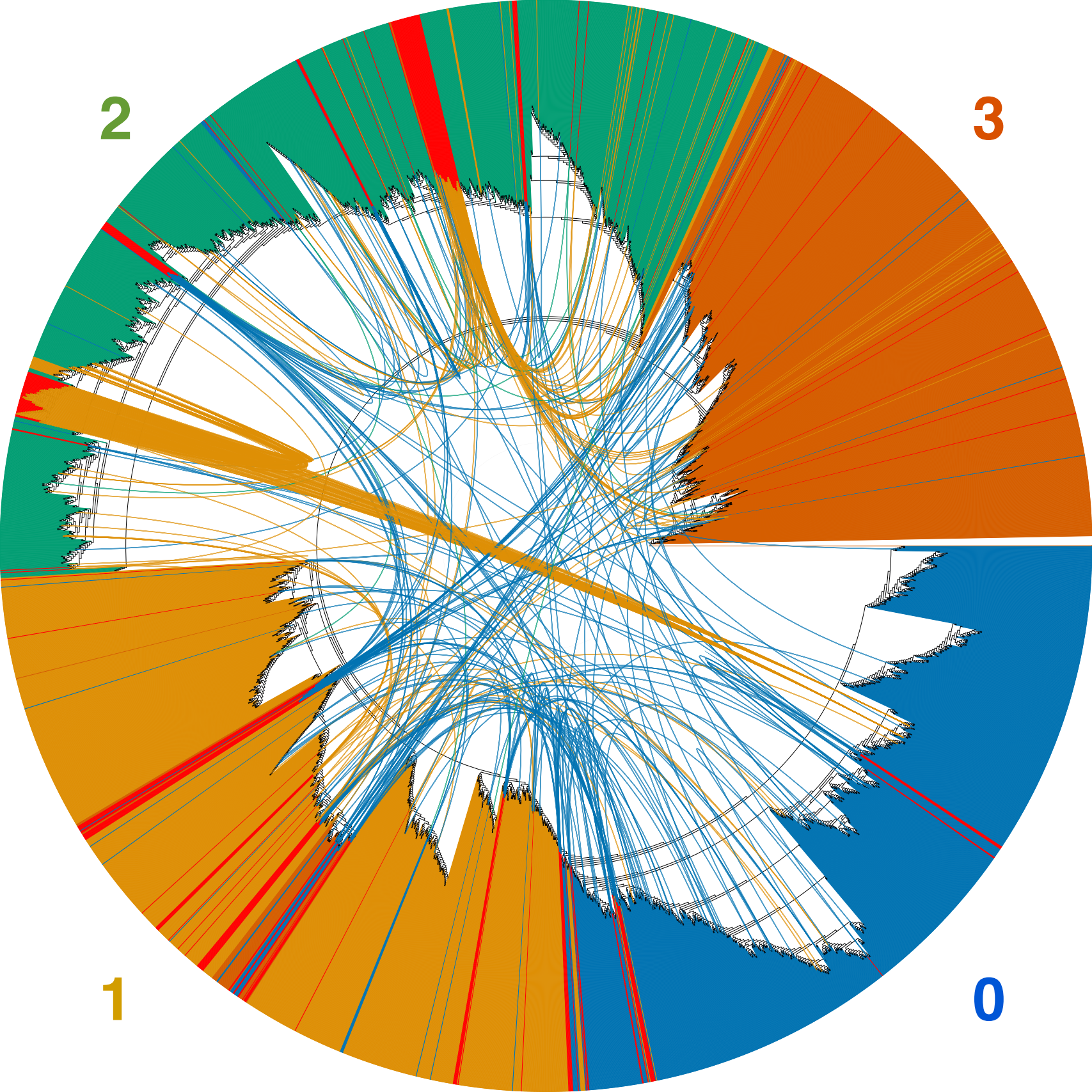
492 **Figure 4.** Effects of sequence composition on the topology of the phylogenetic tree for SARS-
 493 CoV-2. A tree obtained directly from NextStrain (**A**) is first compared to (**B**) the tree computed
 494 using Bolotie consensus sequences for the same set of isolates. (**C**) Shows a tree computed for
 495 the same set of isolates with 210 additional recombinant sequences as identified by Bolotie.
 496 Leaf nodes that correspond to recombinant genomes are labeled with red dots.

497

Clade	Position	Reference	Observed	P(0 Base)	P(1 Base)	A	C	G	T
0	240	C	C	0.3975	0.2358	4	6492	5	1
	3036	C	C	0.9998	0.0001	1	1422	1	9
	8781	C	T	0.9929	0.0012	0	0	0	140
1	14407	C	T	0.0092	0.3297	1	3	2	7359
	17125	T	C	0.0217	0.9783	1	7991	0	6
	20267	A	G	0.0061	0.9935	0	0	13	0
	23402	A	G	0.0066	0.3311	0	0	196	0
0	28143	T	C	0.9988	0.0005	2	8004	0	4

498 **Table 1.** Mutational signature of the EPI_ISL_439137 recombinant isolate (Figure 3A). The table
 499 shows all positions with defining conditional probabilities for each of the parental clades. Read
 500 counts extracted from the data deposited in EBI are provided to illustrate the likely single-isolate
 501 origin of the genome.



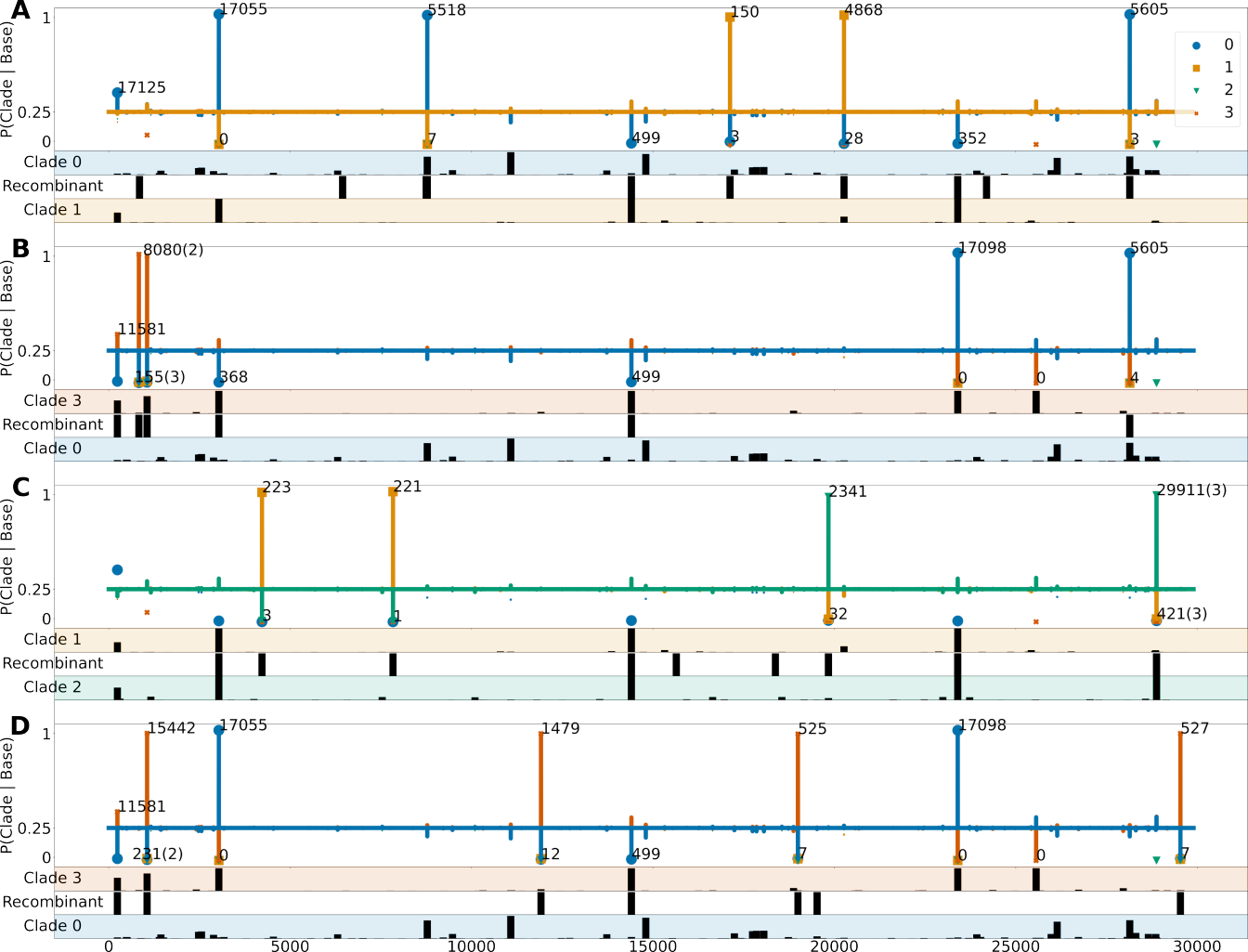


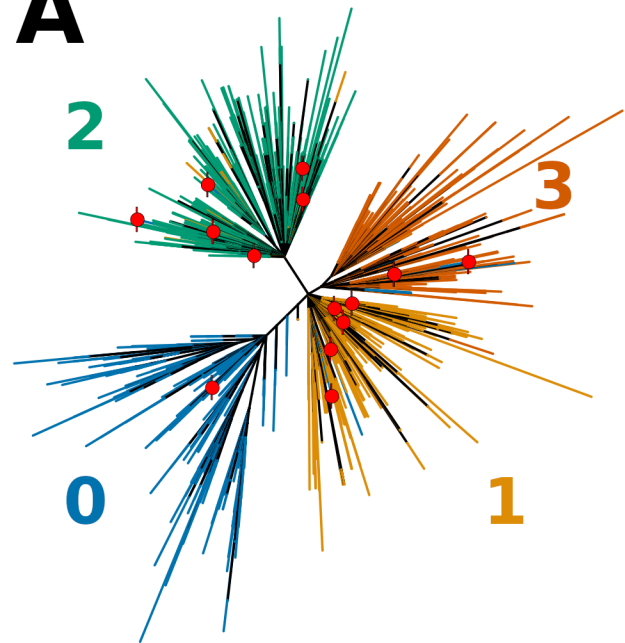
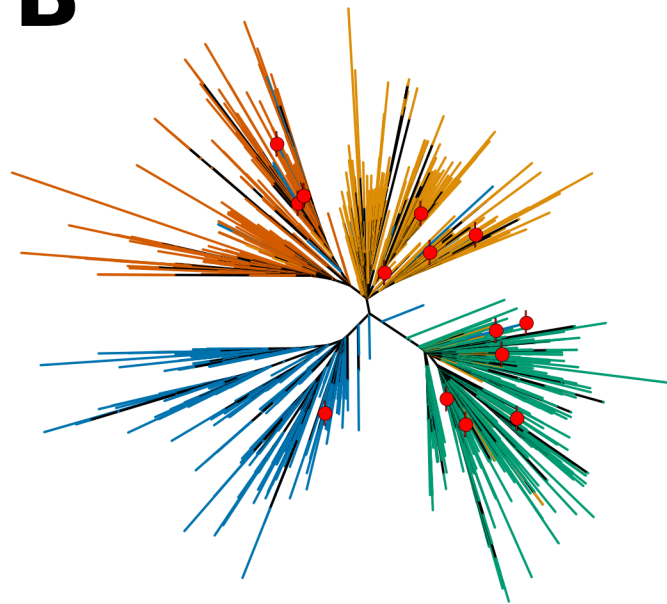
2

3

1

0



A**B****C**