1 **Robust expansion of phylogeny for fast-growing genome sequence data**

2

3 **Authors**

4

5 Yongtao Ye[1,2*], Marcus H. Shum[1,2*], Joseph L. Tsui[1,2*], Guangchuang Yu[3*], David K. Smith[1],

6 Huachen Zhu[1,2,4,5], Joseph T. Wu[1,2], Yi Guan[1,2,4,5], Tommy T. Lam[1,2,4,5,6]

7

8 **Affiliations**

9

10 [1] State Key Laboratory of Emerging Infectious Diseases, School of Public Health, The University

11 of Hong Kong, Hong Kong SAR, P. R. China

12 [2] Laboratory of Data Discovery for Health Limited, 19W Hong Kong Science & Technology

13 Parks, Hong Kong SAR, P. R. China

14 [3] Department of Bioinformatics, School of Basic Medical Sciences, Southern Medical University,

15 Guangzhou, Guangdong, China

16 [4] Guangdong-Hongkong Joint Laboratory of Emerging Infectious Diseases, Joint Institute of

17 Virology (Shantou University/The University of Hong Kong), Shantou, Guangdong, 515063, P.

18 R. China.

19 [5] EKIH (Gewuzhikang) Pathogen Research Institute, Futian District, Shenzhen City, Guangdong,

20 518045, P. R. China.

21 [6] Centre for Immunology & Infection Limited, 17W Hong Kong Science & Technology Parks,

22 Hong Kong SAR, P. R. China

23

24 * These authors contributed equally to the research.

25 # Correspondence: ttylam@hku.hk

26    **Abstract**

27

28    Massive sequencing of SARS-CoV-2 genomes has led to a great demand for adding new samples

29    to a reference phylogeny instead of building the tree from scratch. To address such challenge, we

30    proposed an algorithm 'TIPars' by integrating parsimony analysis with pre-computed ancestral

31    sequences. Compared to four state-of-the-art methods on four benchmark datasets (SARS-CoV-2,

32    Influenza virus, Newcastle disease virus and 16S rRNA genes), TIPars achieved the best

33    performance in most tests. It took only 21 seconds to insert 100 SARS-CoV-2 genomes to a 100k-

34    taxa reference tree using near 1.4 gigabytes of memory. Its efficient and accurate phylogenetic

35    placements and incrementation for phylogenies with highly similar and divergent sequences

36    suggest that it will be useful in a wide range of studies including pathogen molecular

37    epidemiology, microbiome diversity and systematics.

38

39   **Introduction**

40

41   Next generation sequencing (NGS) technologies enable large-scale exploration of diversity and

42   monitoring temporal evolution of organisms, which often involve generating and analyzing large

43   numbers of sequences from new organisms on an ongoing basis. For instance, over 5 million of

44   SARS-CoV-2 genomes have been sequenced within two years of the pandemic (Shu &

45   McCauley, 2017), largely facilitating transmission tracking and disease control. Conventional

46   methods of phylogeny inference from scratch such as those implemented in  IQ-TREE2 (Minh et

47   al., 2020) and FastTree2 (Price, Dehal, & Arkin, 2010) could hardly cope with such rapidly

48   growing huge sequence datasets. Therefore, determining the evolutionary position of new

49   sequences as they become available by placing or inserting them into the reference tree becomes a

50   more efficient alternative. Such 'phylogenetic placement' has been useful for taxonomic

51   classification, while accumulative addition of sequences (incrementing the phylogeny as a result)

52   allow efficient update of the growing phylogeny in a global context.

53

54   Previously published methods such as PhyClass (Filipski, Tamura, Billing-Ross, Murillo, &

55   Kumar, 2015), EPA-ng (Barbera et al., 2019) and pplacer (Matsen, Kodner, & Armbrust, 2010)

56   utilize minimum evolution or maximum likelihood criteria to infer the evolutionary position of

57   the query sequence and place it directly onto a pre-built phylogeny. These algorithms require

58   relatively large computer memory or long runtime which makes massive sequence insertion

59   difficult. Recently, in respect of tracking the diversity of the large amount of SARS-CoV-2 virus

60   genomes, UShER (Yatish Turakhia et al., 2021) was developed to tackle this problem by

61   calculating the 'branch parsimony score' to search for positions of taxa placement only based on

62   sequence mutations to a particular reference. It is extremely fast as compared to the other existing

63   programs. Although the performance of UShER on the SARS-CoV-2 genomes is promising, the

64    placement performance for genome sequences with greater divergence is not well studied.

65

66    We hereby introduce a new approach TIPars, which inserts sequences into a reference phylogeny

67    based on parsimony criterion with the aids of a full multiple sequence alignment of taxa and pre-

68    computed ancestral sequences. The ancestral sequences are useful and efficient in assisting the

69    search of the best placed position because these ancestral sequences often contain rich

70    information in the evolution context of a phylogenetic tree (Loytynoja, Vilella, & Goldman,

71    2012). Recent ancestral sequence reconstruction methods such as PastML (Ishikawa, Zhukova,

72    Iwasaki, & Gascuel, 2019) and RASP4 (Y. Yu, Blair, & He, 2020) have improved speed and

73    accuracy to become feasible in the huge SARS-CoV-2 phylogeny. TIPars searches the position

74    for insertion by calculating the triplet-based minimal substitution score for the query sequence on

75    all branches (Fig. 1A). To compare the performances of different phylogenetic

76    placement/insertion methods including TIPars, UShER, EPA-ng, IQ-TREE2 and PAGAN2

77    (Loytynoja et al., 2012), we applied them on four benchmark datasets (SARS-CoV-2, Influenza

78    virus, Newcastle disease virus and 16S rRNA genes). The first test is single taxon placement. We

79    pruned one taxon from a given phylogenetic tree and applied the methods to place it back. The

80    second is multiple taxa insertion in which a set of taxa was removed and sequentially inserted

81    back. We compared the topology and log likelihood for the trees before pruning and after

82    reinsertion. Our evaluation tests aimed to assess the robustness of the methods on both highly

83    similar sequences and divergent sequences, and whether the phylogenetic tree could be efficiently

84    updated with new sequences that are continuously generated.

85

86    **Results**

87

88    **Computational performance of TIPars and other methods**

89    A number of approaches have been proposed for phylogenetic placement or insertion, but dealing

90    with the vast number of SARS-CoV-2 genome sequences has rendered most of these methods

91    impractical or computationally prohibitive. Based on a reference SARS-CoV-2 phylogenetic tree

92    (SARS2-100k) generated from 96,020 unmasked SARS-CoV-2 sequences of high quality (details

93    in Methodology), we evaluated our proposed program TIPars with UShER, EPA-ng, IQ-TREE2

94    and PAGAN2 by sequentially inserting 100 new sequence samples. Only TIPars and UShER

95    were practicable in terms of running time and memory usage. PAGAN2 were not able to

96    complete the insertion within 96 hours and hence no data was available. Although IQ-TREE2

97    used a lower peak memory than EPA-ng, the running time was the highest among all programs. In

98    contrast, EPA-ng achieved a faster running time than IQ-TREE2 but the peak memory usage was

99    around 1 terabyte (TB) which would not be practicable for general users. As for TIPars, it took

100   only 21 seconds (excluding the input loading time) on a 64-cores server and required about 1.4

101   gigabytes (GB) peak memory usage (Table 1). Another computational performance comparison

102   on smaller dataset with 800 bacterial 16S rRNA sequences (16S) can be checked in table S1 in

103   which PAGAN2 was runnable. Overall, in the SARS2-100k phylogenetic tree, TIPars ran 10-300

104   folds faster than EPA-ng and PAGAN2 with 98.5% to 99.9% less memory used, an efficiency

105   that is comparable to that of the leading program UShER.

106

107   **Single taxon placement**

108

109   Adding a single sequence sample (query) into a reference tree is useful to obtain the phylogenetic

110   placement of the new data, and can be the basic step for expanding the phylogeny with new

111   sequences. We tested TIPars, UShER and EPA-ng on four datasets, including the SARS-CoV-2

112   genomes (SARS2-100k), 16S ribosomal RNA genes (16S), hemagglutinin genes of human

113   seasonal influenza A viruses (H3N2), and Newcastle disease virus genomes (NDV) where the

114    average pairwise genetic distances (substitutions per site) of SARS2-100k and H3N2 are less than

115    0.04 (similar sequences) while those of 16S and NDV are greater than 0.12 (divergent sequences)

116    (details in Methodology; table S2). For the SARS2-100k dataset, EPA-ng was not applied due to

117    impractically large memory requirement and long runtime.

118

119    Based on the postorder traversal, between every 10 taxa we selected one sequence from the

120    SARS2-100k sequence alignment resulting in 9,602 sequences, i.e., 10% of the total taxa in the

121    tree. These selected sequences were individually removed from the reference tree and multiple

122    sequence alignment (MSA) one at a time and used as the query sample for single taxon

123    placement. In datasets of 16S, H3N2 and NDV, all taxa were removed individually and used for

124    the placement test.

125

126    To evaluate the accuracy of each single taxon placement, we calculated the Robinson-Foulds (RF)

127    distance (Robinson & Foulds, 1981) between the reference tree before the taxon removal and the

128    resulting tree after the placement using corresponding programs. An RF distance measures the

129    topological clustering difference between two trees. A zero RF distance indicates that the two

130    trees are topologically identical, and hence the single taxon placement position is exactly the same

131    as the original position, i.e. a true positive.

132

133    With the aid of ancestral information and MSA of full sequences, TIPars performed accurately on

134    phylogenies made of highly similar (SARS2-100k and H3N2) and divergent (16S and NDV)

135    sequences (Fig. 1B). However, a drop in accuracy on more divergent sequences was observed

136    from UShER, perhaps because UShER was only based on the sequence mutations to a particular

137    reference sequence as input, which may lose the insertion information (Yatish Turakhia et al.,

138    2021). In addition, we noted that due to the massive sequencing of SARS-CoV-2 by different

139 research groups, sequencing quality varies and ambiguity bases often occur in the consensus

140 genome sequence data, which could affect the placement accuracy. To account for ambiguity data

141 in sequencing, we used a specific substitution scoring table based on the IUPAC nucleotide

142 ambiguity codes (table S3) for the taxon placement and insertion process (details in

143 Methodology), which achieved a robust performance in sequences of different qualities.

144

145 Notably, when searching through the whole phylogeny for the best position to place a taxon, there

146 may be cases where multiple branches achieve equal minimum substitution scores, and thus the

147 placement will be uncertain. As demonstrated in Fig. 1C, TIPars produced the least number of

148 multiple ambiguously optimal placements in all testing datasets. For example, TIPars generated

149 23% fewer multiple placements than UShER in the SARS2-100k dataset.

150

151 A possible reason for the relatively poor performance of EPA-ng could be that RF distance may

152 not be a reliable metric to compare binary trees derived from the phylogeny with polytomy

153 because there is a very skewed distribution of RF distance when comparing two random binary

154 trees (Bryant & Steel, 2009; Lin, Rajan, & Moret, 2012; Moon & Eulenstein, 2019). It is notable

155 that EPA-ng only processes binary trees. To address this issue, a relaxed criterion for true positive

156 was applied based on whether there are common sister taxa for the removed and re-placed single

157 taxon, as previously used (Yatish Turakhia et al., 2021). With the adjusted true positive

158 measurement, TIPars achieved the highest accuracy in all datasets (fig. S1). While the accuracy of

159 EPA-ng was substantially improved, it was still the lowest among the three tested programs.

160

161 To assess the practicability for extremely large phylogenies, we applied TIPars and UShER in

162 single taxon placement test over the global SARS-CoV-2 phylogenetic tree with 659,885 masked

163 genome sequences (SARS2-660k) downloaded from the Global Initiative on Sharing All

164 Influenza Database (GISAID) (Shu & McCauley, 2017) on the 6th September 2021. A total of

165 65,989 sequences (10% of the total taxa in the tree) were removed and re-inserted individually.

166 Cumulative proportion of single taxon placement result with different RF distance cutoff was

167 shown in Fig. 1D. TIPars produced trees with significantly higher topological similarity to the

168 reference tree with a median RF distance of 0.5 and mean of 5.8 (99% confidence interval (CI) =

169 [5.5-6.1]) as compared to UShER (median RF distance is 3.0 and mean is 31.2 (99% CI = [30.0-

170 32.4])) at 99% significance level (p-value $< 10^{-10}$).

171

172 **Multiple taxa insertion**

173

174 Multiple taxa insertion was an alternative method in determining the phylogenetic position of new

175 sequences over conventional complete phylogeny construction from scratch. TIPars and other

176 three programs (IQ-TREE2, PAGAN2 and UShER) were applied on the four datasets to conduct a

177 comprehensive evaluation of performance.

178

179 In the SARS2-100k dataset, we performed multiple taxa insertion for 100 sets of $10^2$ and $10^3$

180 randomly selected sequences (an example is shown in Fig. 2A) (random100 and random1000)

181 and 100 sets of $10^2$ and $10^3$ successively selected sequences (i.e., a set of successive taxa

182 following the tree postorder traversal; an example is shown in Fig. 2B) (successive100 and

183 successive1000). In the 16S, H3N2 and NDV datasets, 100 sets of 50 sequences were randomly

184 selected. The selected sequences are pruned from the corresponding reference tree and become

185 multiple taxa to be reinserted for each testing set.

186

187 RF distance and tree log-likelihood (LL) were used to evaluate the performance of the multiple

188 sequence insertion. To evaluate the topology accuracy, the resulting tree produced by the four

189     programs were compared to the original reference tree (leaf taxa unpruned) to obtain the RF

190     distance. At the same time, Gamma20 log-likelihoods of the reference tree and the resulting tree

191     after optimizing the branch length were also computed using FastTree2 (double-precision version)

192     and their differences were used for evaluation.

193

194     For the random100 and random1000 datasets, only analyses using TIPars and UShER were able

195     to complete within a reasonable computation time, hence no result from IQ-TREE2 and PAGAN2

196     was present. The resulting trees from multiple taxa insertion using TIPars had a significantly

197     smaller RF distance than those generated using UShER (Fig. 3A). In addition, the log-likelihood

198     of the resulting trees from TIPars was significantly higher than that of UShER (Fig. 3B).

199     Moreover, TIPars resulting trees tended to be very close to the reference tree with smaller log-

200     likelihood differences (fig. S2, A and B). A demonstration of the taxa-insertion was illustrated in

201     Fig. 2A by adding 1000 samples. We observed there were more crossing lines from reference tree

202     to UShER resulting tree indicating more misplaced insertions.

203

204     As to 16S, H3N2 and NDV datasets, TIPars mostly outperformed IQ-TREE2, PAGAN2 and

205     UShER with a significantly lower RF distance and a higher log-likelihood of resulting trees (Fig.

206     3, E to H; fig. S3). In the H3N2 dataset, there was no significant tree log-likelihood difference

207     between TIPars and UShER (Fig. 3G), and in NDV dataset, TIPars performed better than IQ-

208     TREE2 with higher mean log-likelihood but without statistical significance (Fig. 3H). The

209     demonstrations of the taxa-insertion result were visualized in Fig. 2C where UShER, IQ-TREE2

210     and PAGAN2 were less accurate than TIPars.

211

212     For the successive100 and successive1000 datasets, TIPars resulting trees had a significantly

213     larger RF distance than those of UShER (Fig. 3C). However, the log-likelihood of the TIPars

214 resulting trees was significantly higher than that of UShER (Fig. 3D; fig. S2, C and D). By

215 comparing the trees generated from TIPars and UShER (Fig. 2B), the difference is that TIPars

216 inserted some of query taxa (green lines in Fig. 2B; successive taxa pruned from the reference

217 tree) into two subtrees where one of them (the one containing over half the queries) had the same

218 topology as the one in the reference tree. Whereas UShER inserted those queries mostly within a

219 monophyletic clade but it was different from the reference tree. As a result, UShER retained the

220 local topology (better RF distance) (Lin et al., 2012; Smith, 2021) but missed the global topology

221 (worse log-likelihood). Through a RF distance comparison specifically to each query taxon

222 instead of all query taxa, we found that the RF distance resulted from UShER was not

223 significantly higher than that of TIPars (table S4).

224

225 On the other hand, we may suppose that in the situation of random100 and random1000 tests, RF

226 distance would be a suitable metric for comparing the performance of taxa insertions as they are

227 similar to the case of single taxon placements, where most removed taxa are within different

228 monophyletic clades due to randomness (Bryant & Steel, 2009).

229

230 To make the log-likelihood of the resulting trees comparable, we applied FastTree2 to reoptimize

231 the branch lengths with fixed topology (Price et al., 2010). However, compared to the efficiency

232 of taxa insertion (Table 1), the re-optimization is time-consuming. For example, the optimization

233 for a SARS2-100k tree took 10 to 12 hours and required around 125 GB memory (table S5).

234 Therefore, we also computed the log-likelihoods with fixed branch lengths (FLL) using IQ-

235 TREE2, and TIPars still outperformed UShER significantly (fig. S4) by achieving a higher log-

236 likelihood in the resulting tree output directly from the program.

237

238 **Inserting novel sequences**

239    To verify practicability of TIPars in adding novel sequences into a given phylogeny, we further

240    performed an experiment to insert novel real-world SARS-CoV-2 samples into the SARS2-100k

241    reference tree. We randomly selected SARS-CoV-2 samples from GISAID which were not

242    included in the SARS2-100k dataset. Twenty sets of 100, 1000, 5000 and 10000 genome

243    sequences were generated as the queries for taxa insertion using TIPars and UShER.

244

245    Log-likelihoods of the resulting trees from each program were calculated and their pairwise

246    differences between TIPars and UShER were used to evaluate the performance. RF distance was

247    not a suitable metric in this experiment as a comparable reference tree was not available. TIPars

248    provided a resulting tree with a significantly better log-likelihood than UShER in all situations (p-

249    values <0.05; Fig. 4A).

250

251    In addition to tree log-likelihood, we also compared the PANGO lineages (PANGOlins)

252    assignment of the added samples (Rambaut et al., 2020) to validate the accuracy. Only

253    PANGOlins that existed in the reference tree were considered. We assigned each newly inserted

254    sequence with the lineage name of the subtree under the parental node of the inserted position.

255    The subtree was annotated by its descendant reference taxa if all of them were monophyletic

256    (McBroome et al., 2021). A true positive was defined as when the assigned lineage of a query

257    sequence was identical to its original PANGOlins label. In case of queries within unannotated

258    subtrees, we ignored them in the calculation. TIPars outperformed UShER by achieving higher

259    true positive samples on the 100, 1000, 5000 and 10000 insertion datasets with an average of 92%

260    PANGOlins accuracy. The superiority of TIPars was statistically significant under a right-tailed

261    paired t-test (p-values < 0.001) on the 1000, 5000 and 10000 datasets (Fig. 4B and table S6).

262

263    **Discussion**

264    TIPars showed promising taxa placement and insertion accuracy in the phylogenies with

265    homogenous (H3N2 and SARS2-100k) and divergent (16S and NDV) sequences, and in

266    extremely large phylogeny (SARS2-660k) with reasonable runtime and memory usage. Although

267    UShER has a lower accuracy in the divergent sequence datasets (16S and NDV), it ran faster than

268    TIPars (Table 1).

269

270    Reconstruction of ancestral sequences are associated with all taxa across the phylogenetic tree,

271    which could be done using maximum likelihood statistical models or other advanced techniques

272    (Ishikawa et al., 2019; Kosakovsky Pond et al., 2020; Pupko, Pe'er, Shamir, & Graur, 2000; Y.

273    Yu et al., 2020). So ancestral sequences may reveal more accurate (especially intermediate)

274    evolutionary information than the consensus mutation lists along each individual lineage as

275    UShER does. The evolutionary information can be used to distinguish insertion, deletion and

276    substitution events in the searching of taxon placement (Löytynoja & Goldman, 2005), which

277    may help TIPars to be robust on more divergent phylogenies (Loytynoja et al., 2012). Overall,

278    compared to existing phylogenetic placement programs, TIPars is a robust method for a variety of

279    datasets with densely sampled and highly similar sequences of a single species which are

280    common in tracking pathogen epidemiology and transmission, as well as the sequences with

281    greater intraspecific divergence such as the genome datasets at genus, families or higher

282    taxonomic levels for systematics studies.

283

284    Although we showed that TIPars resulting trees with higher tree log-likelihood compared to other

285    programs, a general limitation of the phylogenetic placement method is that errors from incorrect

286    placements accumulate as multiple sequences are inserted sequentially. In order to minimize the

287    error due to large numbers of sequence insertions, it is suggested to conduct tree refinements on

288    not only branch length but also tree topology using different techniques such as nearest-neighbor

289    interchanges (NNIs) and subtree-pruning-regrafting (SPRs) (Price et al., 2010). Furthermore,

290    starting such optimization process with an initial tree of higher log-likelihood may achieve a final

291    tree with better log-likelihood using certain of time (Price et al., 2010). As demonstrated in table

292    S7, for the resulting trees of equal RF distance from both TIPars and UShER (n=28), the branch

293    length optimized trees for TIPars had higher (n=14) or equal (n=12) tree log-likelihoods than the

294    ones resulted from UShER.

295

296    TIPars could facilitate the future development of sequence analysis methods that make use of the

297    phylogenetic placement information. For instance, genome assembly of NGS read data from the

298    metagenome can use phylogenetic positions of the short-read sequences to distinguish between

299    related microbial strains or lineages. With the aid of TIPars, NGS sequences could be inserted to

300    the branches of specific strains or lineages in a reference phylogeny. This can be used in

301    calculating the proportion of strains in mixed infection even when one of the strains is at low

302    abundance in which *de novo* assembly may generate incomplete contigs.

303

304    Since the start of the COVID-19 pandemic, over 5 million SARS-CoV-2 genome sequences have

305    been made publicly available (Shu & McCauley, 2017). With the reduction in cost, the rate of

306    genome sequencing is expected to skyrocket in the future. By providing rapid and memory

307    efficient taxa insertions at high accuracies, TIPars may improve real-time tracing and monitoring

308    of SARS-CoV-2 transmission through the large-scale global phylogenetic analysis of the ever-

309    increasing SARS-CoV-2 genome sequences.

310

311    **Materials and Methods**

312

313    **Implementation of TIPars**

314    After assigning the ancestral sequences at every internal node and taxa sequences at external

315    nodes, TIPars inserts a set of new samples into the reference phylogenetic tree sequentially based

316    on parsimony criteria.

317

318    For a query sequence Q, TIPars computes the minimal substitution score against every branch in

319    the tree. While inserting query Q into to the branch A-B (parent node - child node) at a potential

320    newly added node P (Fig. 1A), the minimal substitution score is the sum of substitution scores

321    that sequence Q differs from both sequence A and sequence B based on a specific substitution

322    scoring table based on the IUPAC nucleotide ambiguity codes (table S3). The single branch with

323    the minimum substitution score $\sigma$ is reported as the best placement.

324

325    However, in terms of multiple placements where more than one branch have the same minimum

326    substitution score, TIPars applies simple but practical rules to filter them to a single best

327    placement such that multiple queries would be inserted sequentially based on one resulting tree.

328    The first priority is to select the branch with node A containing the most numbers of child nodes.

329    The second priority is to select the branch with node A of the lowest node height, that is the total

330    branch length on the longest path from the node to a leaf (Suchard et al., 2018). Finally, in the

331    case where the ambiguity cannot be resolved by the first two priorities, TIPars just turns to a pick

332    up randomly. Even though TIPars will filter out multiple placements, these potential placements

333    will also be printed out for user notice.

334

335    We proposed a local estimation model to calculate the pendant length of the newly introduced

336    branch P-Q ($l_{P-Q}$) which is considering the branch lengths of the local triplet subtree (A,(B,Q))

337    (Fig. 1A). Pendant length is defined as $l_{P-Q} = \sigma / (\delta_A + \delta_B) * l_{A-B}$, where $\delta_A$ and $\delta_B$ are the unique

338    mismatch substitution scores of Q to A and B, and $l_{A-B}$ is the original length of branch A-B. The

339    location of P on branch A-B is determined by the ration of $\delta_A$ and $\delta_B$, i.e., Distal length:

340    $l_{A-P} = \delta_A / (\delta_A + \delta_B) * l_{A-B}$, and Sibling length: $l_{B-P} = \delta_B / (\delta_A + \delta_B) * l_{A-B}$. The ancestral sequence

341    of node P is estimated by majority vote of the nucleotide bases of sequence A, B and Q. To retain

342    the topology of reference tree, a potential nucleotide base of Q will be only derived from A or B.

343    For a special case of $l_{A-B}$ is zero but $\sigma$ is not, TIPars will consider upper branch of A's parent to

344    A for scaling.

345

346    We implemented TIPars using Java with BEAST library (Suchard et al., 2018). Both FASTA and

347    VCF formats are acceptable for loading sequences while NEWICK format is for the tree file.

348    FASTA file is the default setting, but VCF file is more memory efficient for large dataset of high

349    similar sequences, e.g. SARS-CoV-2 virus. To convert a FASTA file to VCF file with all

350    sequence mutations, i.e. insertion, deletion and substitution, we used a Python package

351    PoMo/FastaToVCF.py (Schrempf, Minh, De Maio, von Haeseler, & Kosiol, 2016).

352

353    **Benchmark datasets preparation**

354

355    Unmasked SARS-CoV-2 MSA from GISAID was downloaded on 6th July 2021. Then all SARS-

356    CoV-2 viral genome sequences collected before 1st January 2021 were extracted from the MSA.

357    In order to ensure the sequences used for downstream analysis were complete, SARS-CoV-2

358    genomes with sequence length < 29,000 bp and > 0.5% Ns were removed (namely

359    GISAID202101). To ensure that the global phylogenetic diversity is well represented in the sub-

360    sampled dataset, sequences from all lineages as designated by the PANGO nomenclature system

361    (Rambaut et al., 2020) were sub-sampled. Where fewer than 50 sequences of a given lineage were

362    found in the global dataset, all sequences of the lineage were included. This resulted in a final

363    sub-sampled dataset of 96,020 sequences from 1,249 PANGO lineages, with hCoV-

364    19/Wuhan/WIV04/2019/EPI_ISL_402124 included as the reference genome (namely SARS2-

365    100k). The SARS2-100k reference tree was then built using IQ-TREE2 with GTR model using

366    the EPI_ISL_402124 as root. Ancestral sequences of each internal node were estimated using

367    PastML with the MSA and the IQ-TREE2 generated tree as input.

368

369    Three small but representative nucleotide sequence datasets namely, bacterial 16S rRNA (16S),

370    hemagglutinin genes of human seasonal influenza A viruses (H3N2), and Newcastle disease virus

371    genomes (NDV), were prepared for programs performance comparison. The 16S dataset was

372    downloaded from Genomes OnLine Database (Mukherjee et al., 2019) and randomly down-

373    sampled to 800 sequences. HA sequences of 800 H3N2 viruses were randomly extracted from

374    Influenza Research Database (Zhang et al., 2017). The 235 NDV sequences were downloaded

375    from GenBank. Alignments were constructed using MUSCLE (Edgar, 2004). Reference trees of

376    these datasets were built using RAxML (Stamatakis, 2014) standard hill-climbing heuristic search

377    with 100 multiple inferences and GTRGAMMA model. Ancestral sequences were estimated

378    using ML joint method (Pupko et al., 2000).

379

380    **Novel SARS-CoV-2 query sequence dataset**

381

382    To generate novel query sequences for the 20 sets of 100, 1000, 5000 and 10000 sequences,

383    SARS-CoV-2 genomes that were not included in the SARS2-100k dataset were randomly

384    selected from the GISAID202101 dataset. Selected sequences were then aligned to the SARS2-

385    100k sequences alignment by opening necessary gaps to obtain the full-length MSA. The newly

386    selected sequences were extracted to obtain the final query sample sets. Corresponding new gaps

387    were also added back to the ancestral sequence alignment for each dataset generated. PANGO

388  lineages data for the novel SARS-CoV-2 query sequences and the taxa of reference tree was

389  downloaded from GISAID on 6th July 2021.

390

391  **Benchmark programs**

392

393  We compared TIPars to four state-of-the-art phylogenetic placement tools, namely UShER, EPA-

394  ng, IQ-TREE2 and PAGAN2 while EPA-ng only works for single taxon placement and IQ-

395  TREE2 and PAGAN2 were only used for multiple taxa insertion.

396

397  For the SARS2-100k dataset, only TIPars and UShER were considered as the other programs

398  were not able to complete the computation within a reasonable runtime (Table 1). For the three

399  smaller datasets, we compared all of them comprehensively. Details of the commands used for

400  different programs could be found in table S8.

401

402  TIPars, UShER and EPA-ng would report multiple placements for single taxon insertion. The

403  marked best placements of TIPars and UShER by themselves were used for our accurac

404  evaluation. EPA-ng reports its results sorted by log-likelihood, so the placement with the highest

405  log-likelihood was applied for assessment.

406

407  For any tools that accept only binary tree, i.e., EPA-ng and PAGAN2, we first converted the

408  original polytomous tree to a binary tree using the Ape R package (Paradis & Schliep, 2019).

409

410  When adding unaligned query samples, it is suggested to align them to the MSA of taxa and

411  ancestral sequences in the reference tree using MAFFT ('--add' option) (Katoh & Standley,

412  2013).

413 **Evaluation metrics**

414

415 For single taxon placement evaluation, we first pruned one taxon from the reference tree and re-

416 inserted it back. To assess the consistency between placement algorithms and the typical tree-

417 constructing approach, we proposed using Robinson–Foulds (RF) Distance as a measure of the

418 tree topology accuracy, as calculated by TreeCmp (Bogdanowicz, Giaro, & Wróbel, 2012). When

419 the RF distance between a hypothetical tree and the reference tree is zero, the topology of the

420 hypothetical tree is the same as the reference tree which means the algorithm inserts the query

421 sample into the reference tree topological correctly. Another performance comparison with

422 different true positive definition was conducted for binary trees derived from trees with polytomy

423 using the measurement of whether sister node sets are identical to reference (Y. Turakhia et al.,

424 2020).

425

426 For multiple taxa insertion evaluation, we randomly pruned a set of taxa from the reference tree

427 and re-inserted them back. In addition to using RF distance to compare the hypothetical tree

428 against the reference tree, we also calculated the log-likelihood of the hypothetical tree as a

429 measurement of the accuracy of the taxa insertions. We applied two methods to compute log-

430 likelihoods including FastTree2 (double-precision version) (Gamma20 Log-Likelihood) (Price et

431 al., 2010) for optimized branch length, and IQ-TREE2 (Log-Likelihood (Fixed Br)) for fixed

432 branch length.

433

434 EPA-ng outputs the placement information (placed branch, distal length, and pendant length) for a

435 query without the construction of the final tree. In order to compute the RF distance, we assisted

436 EPA-ng in inserting the query into the reference tree to generate the hypothetical tree.

437

438    IQ-TREE2 and PAGAN2 support initial tree, but they are not exactly based on the input tree

439    topology for construction, so RF distance to original reference tree is not suitable for them.

440    Note that UShER outputs the final constructed tree using the number of mutations as branch

441    length (otherwise no branch length would be specified at branches modified), so we modified its

442    branch length as number of mutations divided by alignment length in calculation of log-likelihood

443    with fixed branch length model.

444

445    **Statistics**

446

447    99% t-test confident intervals and 99% paired t-test p-value (right tail) for the results of TIPars

448    against other programs were computed by Matlab R2013b. All violin graphs were generated by R

449    4.1.1 using the package *ggstatsplot* (Patil, 2021). Illustration and annotation of phylogenetic trees

450    were done using the R package *ggtree* (G. Yu, Smith, Zhu, Guan, & Lam, 2017).

451

452    **Data and materials availability**

453

454    SARS2-CoV-2 data used in this work were all downloaded from GISAID

455    (https://www.gisaid.org/). TIPars is available at https://github.com/id-bioinfo/TIPars.

456

457    **Acknowledgments**

458

459    We gratefully acknowledge the following Authors from the Originating laboratories responsible

460    for obtaining the specimens and the Submitting laboratories where genetic sequence data were

461    generated and shared via GISAID Initiative, on which this research is based. A full

462    acknowledgement table can be found with two EPI_SET-IDs, i.e., EPI_SET_20211201vz and

463     EPI_SET_20211206tc, in Data Acknowledgement Locator under GISAID resources

464     (https://www.gisaid.org/).

465

472

473     **Competing interests**

474

475     Authors declare that they have no competing interests.

476

477     **References**

478

479     Barbera, P., Kozlov, A. M., Czech, L., Morel, B., Darriba, D., Flouri, T., & Stamatakis, A.
480        (2019). EPA-ng: Massively Parallel Evolutionary Placement of Genetic Sequences. *Syst*
481        *Biol, 68*(2), 365-369. doi:10.1093/sysbio/syy054
482     Bogdanowicz, D., Giaro, K., & Wróbel, B. (2012). TreeCmp: Comparison of Trees in Polynomial
483        Time. *Evolutionary Bioinformatics Online, 8*, 475-487. doi:10.4137/EBO.S9657
484     Bryant, D., & Steel, M. A. (2009). Computing the Distribution of a Tree Metric. *IEEE ACM*
485        *Trans. Comput. Biol. Bioinform., 6*, 420-426.
486     Edgar, R. C. (2004). MUSCLE: a multiple sequence alignment method with reduced time and
487        space complexity. *BMC Bioinformatics, 5*, 113. doi:10.1186/1471-2105-5-113
488     Filipski, A., Tamura, K., Billing-Ross, P., Murillo, O., & Kumar, S. (2015). Phylogenetic
489        placement of metagenomic reads using the minimum evolution principle. *BMC Genomics,*
490        *16*(1), S13. doi:10.1186/1471-2164-16-S1-S13
491     Ishikawa, S. A., Zhukova, A., Iwasaki, W., & Gascuel, O. (2019). A Fast Likelihood Method to
492        Reconstruct and Visualize Ancestral Scenarios. *Molecular Biology and Evolution, 36*(9),
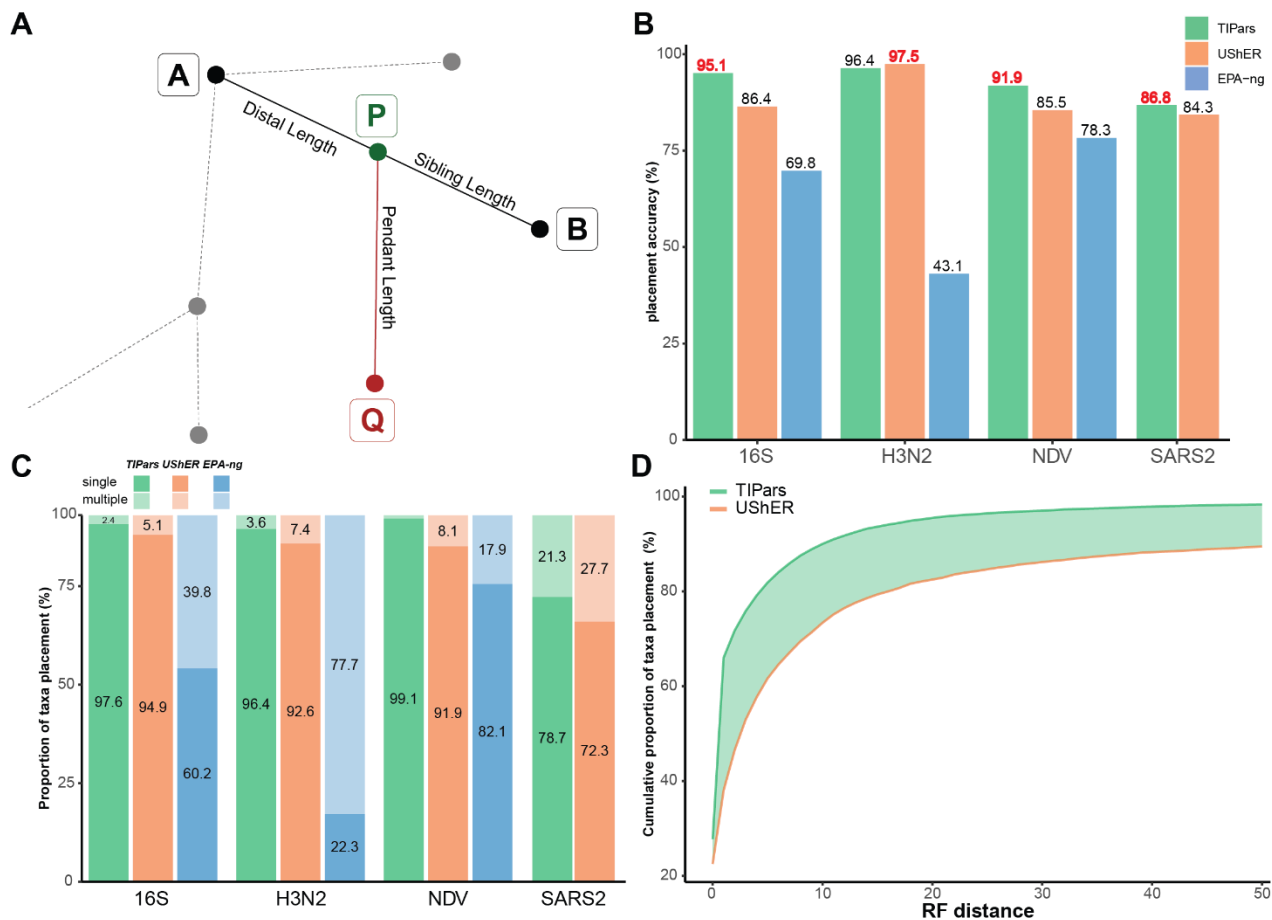493        2069-2085. doi:10.1093/molbev/msz131
494     Katoh, K., & Standley, D. M. (2013). MAFFT Multiple Sequence Alignment Software Version 7:
495        Improvements in Performance and Usability. *Molecular Biology and Evolution, 30*(4),
496        772-780. doi:10.1093/molbev/mst010

497 Kosakovsky Pond, S. L., Poon, A. F. Y., Velazquez, R., Weaver, S., Hepler, N. L., Murrell,
498         B., . . . Muse, S. V. (2020). HyPhy 2.5—A Customizable Platform for Evolutionary
499         Hypothesis Testing Using Phylogenies. *Molecular Biology and Evolution, 37*(1), 295-299.
500         doi:10.1093/molbev/msz197
501 Lin, Y., Rajan, V., & Moret, B. M. E. (2012). A Metric for Phylogenetic Trees Based on
502         Matching. *IEEE/ACM Transactions on Computational Biology and Bioinformatics, 9*(4),
503         1014-1022. doi:10.1109/TCBB.2011.157
504 Löytynoja, A., & Goldman, N. (2005). An algorithm for progressive multiple alignment of
505         sequences with insertions. *Proceedings of the National Academy of Sciences of the United
506         States of America, 102*(30), 10557. doi:10.1073/pnas.0409137102
507 Loytynoja, A., Vilella, A. J., & Goldman, N. (2012). Accurate extension of multiple sequence
508         alignments using a phylogeny-aware graph algorithm. *Bioinformatics, 28*(13), 1684-1691.
509         doi:10.1093/bioinformatics/bts198
510 Matsen, F. A., Kodner, R. B., & Armbrust, E. V. (2010). pplacer: linear time maximum-
511         likelihood and Bayesian phylogenetic placement of sequences onto a fixed reference tree.
512         *BMC Bioinformatics, 11*(1), 538. doi:10.1186/1471-2105-11-538
513 McBroome, J., Thornlow, B., Hinrichs, A. S., Kramer, A., De Maio, N., Goldman, N., . . .
514         Turakhia, Y. (2021). A Daily-Updated Database and Tools for Comprehensive SARS-
515         CoV-2 Mutation-Annotated Trees. *Molecular Biology and Evolution*.
516         doi:10.1093/molbev/msab264
517 Minh, B. Q., Schmidt, H. A., Chernomor, O., Schrempf, D., Woodhams, M. D., von Haeseler, A.,
518         & Lanfear, R. (2020). IQ-TREE 2: New Models and Efficient Methods for Phylogenetic
519         Inference in the Genomic Era. *Mol Biol Evol, 37*(5), 1530-1534.
520         doi:10.1093/molbev/msaa015
521 Moon, J., & Eulenstein, O. (2019, 2019//). *The Cluster Affinity Distance for Phylogenies.* Paper
522         presented at the Bioinformatics Research and Applications, Cham.
523 Mukherjee, S., Stamatis, D., Bertsch, J., Ovchinnikova, G., Katta, H. Y., Mojica, A., . . . Reddy,
524         T. (2019). Genomes OnLine database (GOLD) v.7: updates and new features. *Nucleic
525         Acids Res, 47*(D1), D649-D659. doi:10.1093/nar/gky977
526 Paradis, E., & Schliep, K. (2019). ape 5.0: an environment for modern phylogenetics and
527         evolutionary analyses in R. *Bioinformatics, 35*(3), 526-528.
528         doi:10.1093/bioinformatics/bty633
529 Patil, I. (2021). Visualizations with statistical details: The 'ggstatsplot' approach. *PsyArXiv*.
530         doi:10.21105/joss.03167
531 Price, M. N., Dehal, P. S., & Arkin, A. P. (2010). FastTree 2--approximately maximum-likelihood
532         trees for large alignments. *PLoS One, 5*(3), e9490. doi:10.1371/journal.pone.0009490
533 Pupko, T., Pe'er, I., Shamir, R., & Graur, D. (2000). A fast algorithm for joint reconstruction of
534         ancestral amino acid sequences. *Mol Biol Evol, 17*(6), 890-896.
535         doi:10.1093/oxfordjournals.molbev.a026369
536 Rambaut, A., Holmes, E. C., O'Toole, Á., Hill, V., McCrone, J. T., Ruis, C., . . . Pybus, O. G.
537         (2020). A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic
538         epidemiology. *Nature Microbiology, 5*(11), 1403-1407. doi:10.1038/s41564-020-0770-5
539 Robinson, D. F., & Foulds, L. R. (1981). Comparison of phylogenetic trees. *Mathematical
540         Biosciences, 53*(1), 131-147. doi:https://doi.org/10.1016/0025-5564(81)90043-2
541 Schrempf, D., Minh, B. Q., De Maio, N., von Haeseler, A., & Kosiol, C. (2016). Reversible
542         polymorphism-aware phylogenetic models and their application to tree inference. *J Theor
543         Biol, 407*, 362-370. doi:10.1016/j.jtbi.2016.07.042
544 Shu, Y., & McCauley, J. (2017). GISAID: Global initiative on sharing all influenza data - from
545         vision to reality. *Euro surveillance : bulletin Europeen sur les maladies transmissibles =*

546      *European communicable disease bulletin, 22*(13), 30494. doi:10.2807/1560-
547          7917.ES.2017.22.13.30494
548  Smith, M. R. (2021). Information theoretic generalized Robinson–Foulds metrics for comparing
549          phylogenetic trees. *Bioinformatics, 37*(14), 2077-2078.
550          doi:10.1093/bioinformatics/btab200
551  Stamatakis, A. (2014). RAxML version 8: a tool for phylogenetic analysis and post-analysis of
552          large phylogenies. *Bioinformatics, 30*(9), 1312-1313. doi:10.1093/bioinformatics/btu033
553  Suchard, M. A., Lemey, P., Baele, G., Ayres, D. L., Drummond, A. J., & Rambaut, A. (2018).
554          Bayesian phylogenetic and phylodynamic data integration using BEAST 1.10. *Virus*
555          *evolution, 4*(1), vey016. doi:10.1093/ve/vey016. (Accession No. 29942656)
556  Turakhia, Y., Thornlow, B., Hinrichs, A. S., De Maio, N., Gozashti, L., Lanfear, R., . . . Corbett-
557          Detig, R. (2020). Ultrafast Sample Placement on Existing Trees (UShER) Empowers
558          Real-Time Phylogenetics for the SARS-CoV-2 Pandemic. *bioRxiv*.
559          doi:10.1101/2020.09.26.314971
560  Turakhia, Y., Thornlow, B., Hinrichs, A. S., De Maio, N., Gozashti, L., Lanfear, R., . . . Corbett-
561          Detig, R. (2021). Ultrafast Sample placement on Existing tRees (UShER) enables real-
562          time phylogenetics for the SARS-CoV-2 pandemic. *Nature Genetics, 53*(6), 809-816.
563          doi:10.1038/s41588-021-00862-7
564  Yu, G., Smith, D. K., Zhu, H., Guan, Y., & Lam, T. T.-Y. (2017). ggtree: an r package for
565          visualization and annotation of phylogenetic trees with their covariates and other
566          associated data. *Methods in Ecology and Evolution, 8*(1), 28-36.
567          doi:https://doi.org/10.1111/2041-210X.12628
568  Yu, Y., Blair, C., & He, X. (2020). RASP 4: Ancestral State Reconstruction Tool for Multiple
569          Genes and Characters. *Molecular Biology and Evolution, 37*(2), 604-606.
570          doi:10.1093/molbev/msz257
571  Zhang, Y., Aevermann, B. D., Anderson, T. K., Burke, D. F., Dauphin, G., Gu, Z., . . .
572          Scheuermann, R. H. (2017). Influenza Research Database: An integrated bioinformatics
573          resource for influenza virus research. *Nucleic Acids Res, 45*(D1), D466-D474.
574          doi:10.1093/nar/gkw857
575

576 **Figures and Tables**

577



580 **Fig. 1. Illustration of phylogenetic placement and single taxon placement performance.** (**A**)

581 Illustration of the placement for a query sequence. "Q" indicates the query sequence, "A" and "B"

582 represent the existing nodes in the reference tree. "P" represents the parental node of "Q"

583 generated by TIPars. Minimum substitution score is calculated based on the triplet formed by A-

584 B-Q. (**B**) Bar charts represent the accuracy of single taxon placement on 16S, H3N2, NDV and

585 SARS2-100k datasets using TIPars, UShER and EPA-ng respectively. Accuracy is indicated on

586 top of each bar and the highest accuracy in each dataset is highlighted in red. (**C**) Stacked bar

587 charts show the proportion of single and multiple taxon placement result for TIPars (Green),
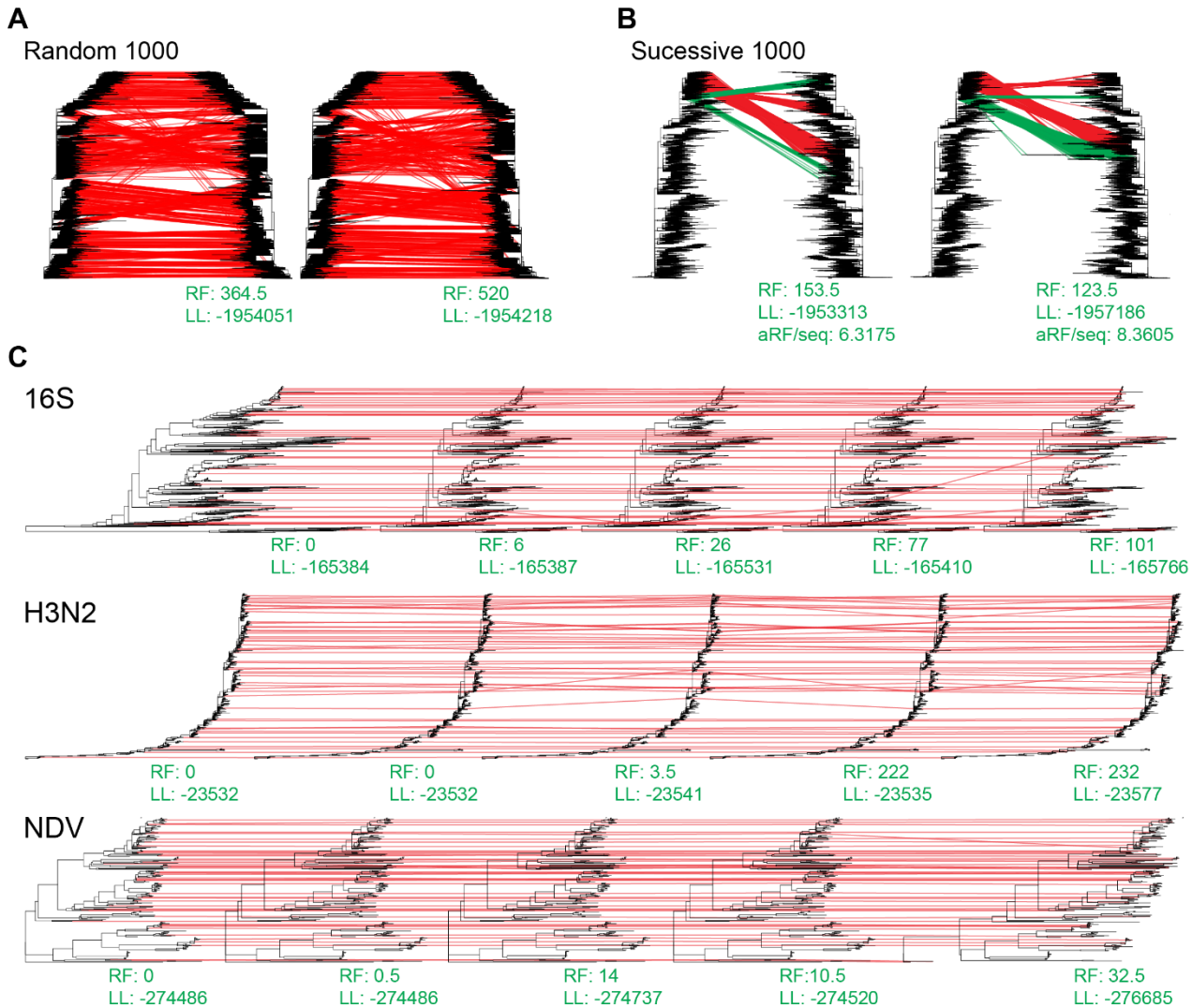
588 UShER (Orange) and EPA-ng (Blue) on the 16S, H3N2, NDV and SARS2-100k datasets.

589 Proportion with > 1% is indicated within the bar. (**D**) Cumulative proportion of single taxa

590    placement on the SARS2-660k dataset with different RF distance cutoff. Highlighted area

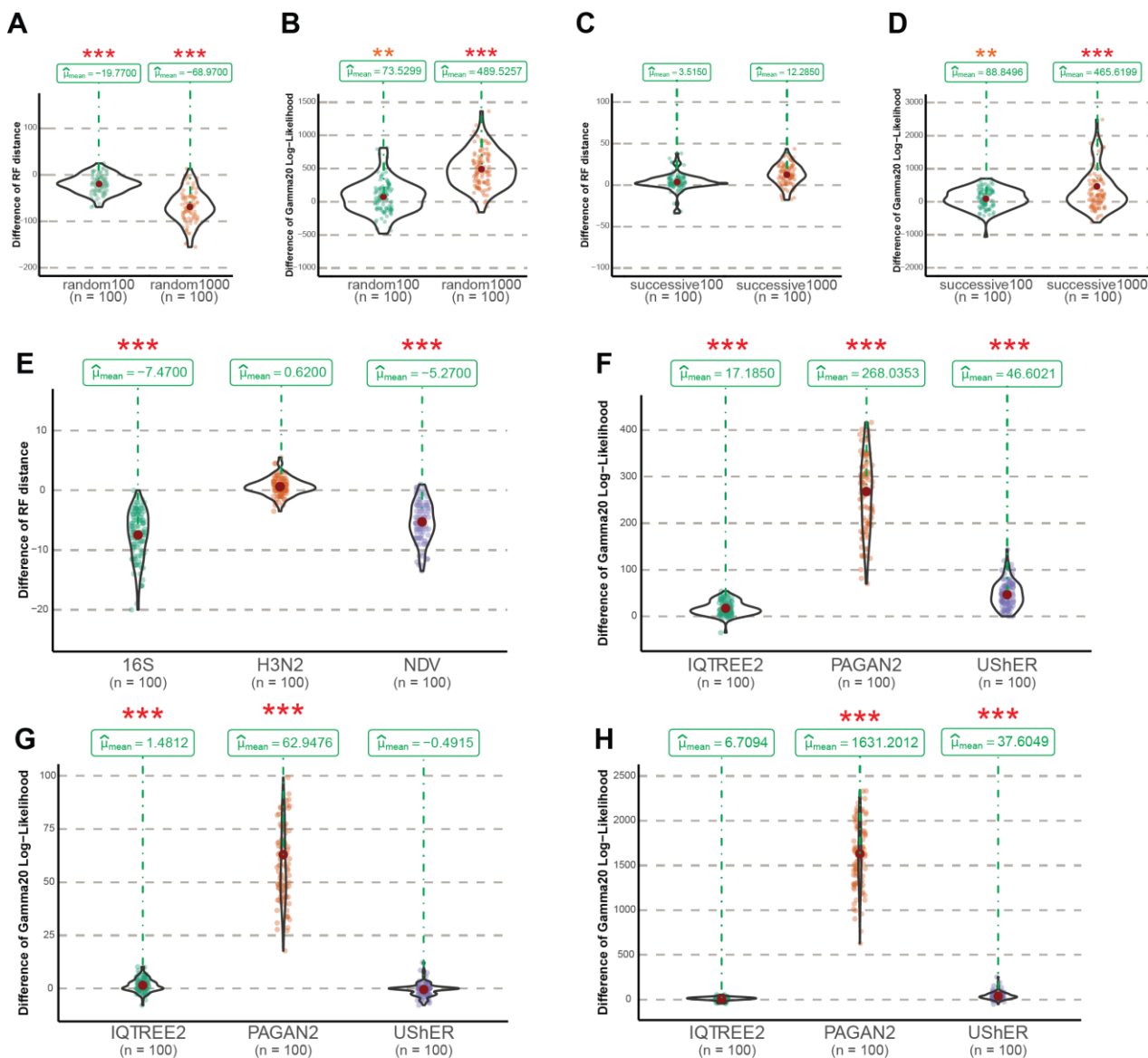591    represents the difference between TIPars and UShER.

592

**A** Random 1000

RF: 364.5
LL: -1954051

RF: 520
LL: -1954218

**B** Sucessive 1000

RF: 153.5
LL: -1953313
aRF/seq: 6.3175

RF: 123.5
LL: -1957186
aRF/seq: 8.3605

**C**

16S

RF: 0
LL: -165384

RF: 6
LL: -165387

RF: 26
LL: -165531

RF: 77
LL: -165410

RF: 101
LL: -165766

H3N2

RF: 0
LL: -23532

RF: 0
LL: -23532

RF: 3.5
LL: -23541

RF: 222
LL: -23535

RF: 232
LL: -23577

NDV

RF: 0
LL: -274486

RF: 0.5
LL: -274486

RF: 14
LL: -274737

RF:10.5
LL: -274520

RF: 32.5
LL: -276685

**Fig. 2. Taxa insertion visualization.** (**A**) A demonstration of TIPars resulting tree (Left) and UShER resulting tree (Right) paired with the reference SARS2-100k reference tree (Left tree in both figures) for the insertion of randomly selected 1000 taxa sequences. Red lines link the corresponding positions of inserted taxa between reference and resulting tree. (**B**) A demonstration of TIPars resulting tree (Left) and UShER resulting tree (Right) paired with the reference SARS2-100k reference tree (Left tree in both figures) for the insertion of successively selected 1000 taxa sequences. Green lines indicate different taxa insertion positions between TIPars and UShER. Averaged RF-distance per sequence (aRF/seq) comparing to the reference tree is shown at the bottom. (**C**) Demonstrations of the resulting trees for randomly selected 50

604 taxa in NDV, 16S (Midpoint rooted) and H3N2 datasets. From the left to the right are trees of

605 reference, TIPars, UShER, IQ-TREE2 and PAGAN2. RF distance (RF) compared to the reference

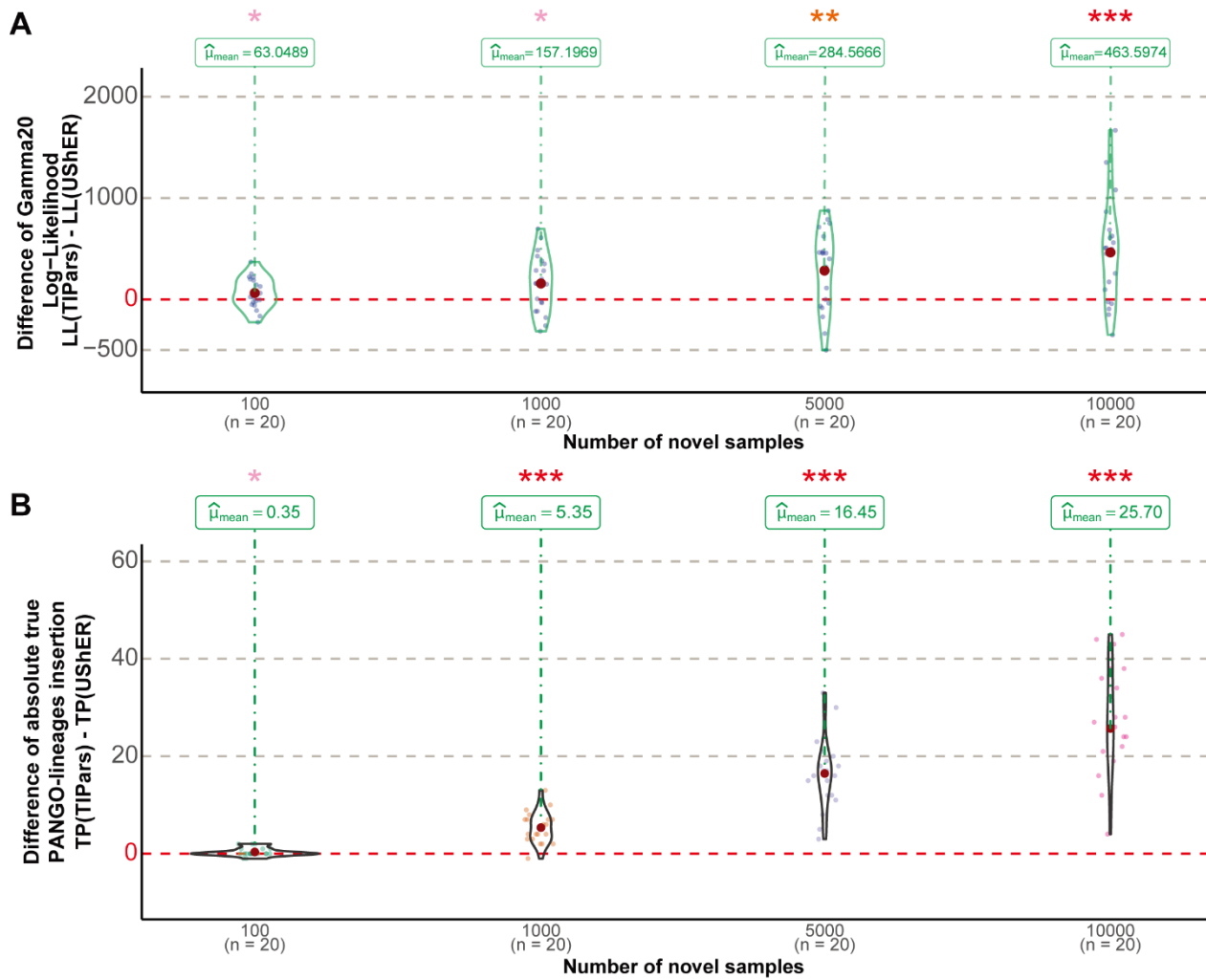606 tree and the Gamma20 log-likelihood (LL) are shown at the bottom of each resulting tree.

607

**Fig. 3. Multiple sequences insertion performance.** (**A-D**) Violin graphs show the distribution of paired differences of the RF distance and the Gamma20 log-likelihood between the optimized resulting trees generated by TIPars and UShER (TIPars - UShER) for the random 100, 1000 and successive 100 and 1000 multiple sequences insertions. (**E**) Distribution of the paired difference of the RF distance between the optimized resulting trees generated by TIPars and UShER (TIPars - UShER) on 16S, H3N2 and NDV random 50 multiple sequences insertions. (**F-H**) Distribution of the paired difference of the Gamma20 log-likelihood between the optimized resulting trees generated by TIPars and the three other programs (TIPars - Others) on 16S (**F**), H3N2 (**G**) and

618    NDV (**H**) random 50 multiple sequences insertions. P-value for the right-sided paired t-test is

619    indicated by the asterisk on top of each violin diagram, where p<0.05 is indicated by one pink

620    asterisk (*), p<0.01 by two orange asterisks (**) and p<0.001 by three red asterisks (***).

621

**Fig. 4. Performance of inserting actual novel sequences.** (**A, B**) Violin graph represents the distribution of the paired differences between the Gamma20 log-likelihood (LL) (**A**) and the absolute number of true PANGO-lineages insertion (TP) (**B**) of TIPars over UShER. p-value for the right-sided paired t-test was indicated by the asterisk on top of each violin diagram, where p<0.05 indicated by one pink asterisk (*), p<0.01 by two orange asterisks (**) and p<0.001 by three red asterisks (***).

631 **Table 1. Average running time and memory used through 10 repeated runs of**

632 **inserting/placing 100 genome samples into SARS2-100k reference tree.** Tests were running on

633 a server of 64 Intel Xeon Gold 6242 CPU cores and 1500 GB RAM. We also compared TIPars

634 with UShER on a general computer with 8 CPU cores. TIPars ran with a JAVA setting of

635 Xmx1G. The running time of UShER contains its necessary computation of 'mutation-annotated

636 tree'. PAGAN2 was not runnable for this dataset. N/A indicates that data are not applicable.

637

| Tools | CPU cores assigned | Mean insertion time (HH:MM:SS) | Mean running time (HH:MM:SS) | Mean peak memory (GB) |
|-------|--------------------|--------------------------------|------------------------------|-----------------------|
| **TIPars** | 64 | 0:00:21 | 0:00:52 | 1.39 |
| **TIPars** | 8 | 0:00:31 | 0:01:03 | 1.18 |
| **UShER** | 64 | 0:00:02 | 0:03:14 | 0.84 |
| **UShER** | 8 | 0:00:05 | 0:05:14 | 0.16 |
| **EPA-ng** | 64 | 0:04:45 | 0:10:25 | 1022.14 |
| **IQ-TREE2** | 64 | N/A | 5:49:10 | 101.10 |
| **PAGAN2** | 64 | N/A | N/A | N/A |

638

639