

A coarse-graining, ultrametric approach to resolve the phylogeny of prokaryotic strains with frequent homologous recombination

Tin Yau Pang (✉ pang@hhu.de)

Research article

Keywords: phylogenetics, homologous recombination, ultrametric tree

Posted Date: December 10th, 2019

DOI: <https://doi.org/10.21203/rs.2.18054/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.
[Read Full License](#)

Version of Record: A version of this preprint was published at BMC Evolutionary Biology on May 7th, 2020. See the published version at <https://doi.org/10.1186/s12862-020-01616-5>.

1

2 **A coarse-graining, ultrametric approach to resolve the phylogeny of prokaryotic strains**
3 **with frequent homologous recombination**

4

5 Tin Yau Pang

6

7 Institute for Computer Science, Heinrich Heine University, Düsseldorf, 40225, Germany

8 Tel: +49-211-81-11651; Fax: +49-211-81-15767;

9

10 Email: pang@hhu.de

11

12

13

14

15

16 **ABSTRACT**

17 **Background**

18 A frequent event in the evolution of prokaryotic genomes is homologous recombination,
19 where a foreign DNA stretch replaces a genomic region similar in sequence. Recombination
20 can affect the relative position of two genomes in a phylogenetic reconstruction in two
21 different ways: (i) one genome can recombine with a DNA stretch that is similar to the other
22 genome, thereby reducing their pairwise sequence divergence; (ii) one genome can
23 recombine with a DNA stretch from an outgroup genome, increasing the pairwise divergence.
24 While several recombination-aware phylogenetic algorithms exist, many of these cannot
25 account for both types of recombination; some algorithms can, but do so inefficiently.
26 Moreover, many of them reconstruct the ancestral recombination graph (ARG) to help infer
27 the genome tree, and require that a substantial portion of each genome has not been affected
28 by recombination, a sometimes unrealistic assumption.

29 **Results**

30 Here, we propose a coarse-graining approach for phylogenetic reconstruction (CGP), which
31 is recombination-aware but forgoes ARG reconstruction, applicable even if all genomic
32 regions have experienced substantial amounts of recombination, and can be used on both
33 nucleotide and amino acid sequences. CGP considers the local density of substitutions along
34 pairwise genome alignments, fitting a model to the empirical distribution of substitution
35 density to infer the pairwise coalescent time. Given all pairwise coalescent times, CGP
36 reconstructs an ultrametric tree representing vertical inheritance. Based on simulations, we
37 show that the proposed approach can reconstruct ultrametric trees with accurate topology,
38 branch lengths, and root positioning. Applied to a set of *E. coli* strains, the reconstructed trees
39 are most consistent with gene distributions when inferred from amino acid sequences, a data
40 type that cannot be utilized by many alternative approaches.

41 **Conclusions**

42 The CGP algorithm is more accurate than alternative recombination-aware methods for
43 ultrametric phylogenetic reconstructions.

44 **KEYWORDS**

45 phylogenetics, homologous recombination, ultrametric tree

46 **BACKGROUND**

47 The transfer of DNA stretches from one prokaryotic genome to another—also called
48 horizontal gene transfer (HGT) or lateral gene transfer (LGT)—is a major driver of
49 prokaryotic evolution [1]. It is caused by a variety of mechanisms, including transformation,
50 transduction, conjugation, and gene transfer agents [2, 3]. Many prokaryotic genomes encode
51 defense systems against foreign DNA, such as the restriction modification system [4]. A
52 foreign DNA stretch that enters the prokaryotic cell and survives these host defenses may be
53 incorporated into the host genome. If the incoming DNA stretch is highly similar to a stretch
54 on the host genome, homologous recombination may occur, where the incoming DNA stretch
55 homologously recombines with the host stretch and overwrites it [5, 6]. Alternatively, the
56 incoming stretch may be inserted directly into the host genome through non-homologous
57 recombination.

58 HGT allows the fast spread of genes in prokaryotic pangenomes, and facilitates rapid
59 adaptation to environmental changes. A point in case is the spread of antibiotic resistance
60 genes in pathogenic bacteria via HGT [7]. But recombination is also crucial for the long-term
61 maintenance of prokaryotic populations, as it helps to repair DNA damaged by deleterious
62 mutations, thereby avoiding the mutational meltdown of Muller’s ratchet [8] in that sense,
63 prokaryotic recombination may fulfill the same function as does sex in eukaryotes.

64 Computational modelling also suggests that recombination may help prokaryotes to purge
65 selfish mobile genetic elements [9].

66 Recombination can severely affect phylogeny reconstructions. Its effects on genome
67 divergence are complex. Depending on the circumstances, recombination can speed up the
68 divergence of a genome pair or slow it down [6]; its effects may severely affect the accuracy
69 of estimated branch lengths of phylogenetic tree. For example, (i) when a stretch of genome
70 X is replaced by DNA from genome Y, some of the single nucleotide polymorphisms that
71 previously differentiated X and Y will be erased, shortening the apparent evolutionary
72 distance between the two genomes. Conversely, (ii) when X recombines with a DNA stretch
73 of an outgroup genome (a genome that diverged before the split of the X and Y lineages),
74 then it introduces further nucleotide polymorphisms into X, thereby increasing the apparent
75 X-Y distance.

76 Multilocus sequence typing (MLST) aims to extract sequences of housekeeping genes
77 from prokaryotic genomes, which can then be utilized to resolve evolutionary relationships
78 [10]. However, MLST genes may also experience frequent recombination, and phylogenetic
79 reconstruction without accounting for recombination can compromise the resulting trees [11].
80 In fact, the frequency with which recombination affects a gene can be of the same order of
81 magnitude as the corresponding mutation rate [5]. Thus, if two lineages recombine with each
82 other, application of conventional phylogenetic algorithms without accounting for
83 recombination will generally lead to an underestimation of the age of the common ancestors
84 [12]. When there are more than two strains, recombination affects the reliability of inference
85 of relative divergence times between strains and may hence compromise both tree topology
86 and branch length estimates.

87 There are several popular recombination-aware algorithms, including
88 ClonalFrameML (and its predecessor ClonalFrame) [13, 14], the Bacter package in BEAST2

89 (which implements the ClonalOrigin model) [15, 16], and Gubbins [17]; there are also non-
90 phylogenetic algorithms that detect recombination, such as BratNextGen and fastGEAR [18,
91 19]. These recombination-aware algorithms may reconstruct the ancestral recombination
92 graph (ARG), which describes the history of transfer and homologous recombination of local
93 genomic stretches across the genomes, to help infer the tree of phylogenetic inheritance of the
94 genomes. While these algorithms can identify genomic stretches with high numbers of
95 substitutions due to recombination with distant strains and thus account for type (ii)
96 recombination effects, many do not take type (i) recombination effects into account; Bacter
97 can account for type (i) effects, but is not computationally efficient for long genome
98 sequences. Some of these algorithms rely on the assumption of low recombination rates, such
99 that a substantial part of a genome remains clonal and has not been affected by
100 recombination. This is unrealistic at least for some bacteria: *e.g.*, *E. coli* strains whose DNA
101 sequences have diverged by more than 1.3% share very few stretches larger than a few kb
102 that have not been affected by recombination in at least one of the two lineages [5].
103 Moreover, ARG may only reveal the latest recombination events on a genomic stretch, but its
104 ability to recover the earlier events on the same stretch is limited, since each recombination
105 erases the history of previous recombinations; this uncertainty on earlier recombination
106 events may introduce error in branch length prediction.

107 In this paper, we propose a novel approach to phylogenetic reconstruction that neither
108 assumes low recombination rates nor relies on ARG reconstruction. Our approach follows a
109 coarse-graining model, which considers the local density of substitutions on a sequence
110 alignment instead of site-specific substitutions [5]; the model describes how different
111 parameters, such as mutation rate, recombination rate or coalescent time between a pair of
112 genomes, affect the shape of their distribution of substitution density. It fits the empirical
113 distribution of substitution density of genome pairs to the model, which allows the inference

114 of the matrix of coalescent time between the genome pairs and thereafter their ultrametric
115 phylogenetic tree. In short, it forgoes the reconstruction of ARG and infers the branch
116 lengths—coalescent times—from the relative abundance of genomic segments with different
117 number of substitutions, and is also applicable to both nucleotide and amino acid sequences.
118 The source code implementing this model is available at [https://github.com/TinPang/coarse-](https://github.com/TinPang/coarse-graining-phylogenetics)
119 [graining-phylogenetics](https://github.com/TinPang/coarse-graining-phylogenetics).

120 **METHODS**

121 **Overview of the CGP algorithm that reconstructs ultrametric phylogenetic tree**

122 The proposed coarse-graining phylogenetic (CGP) algorithm takes n aligned genome
123 sequences as input (Figure 1a), which can be either nucleotide or amino acid sequences. For
124 each pair of sequences, it divides them into L_{seg} equal-sized segments, each segment has l_s
125 sites, and enumerates the single sites polymorphisms (SSPs)—sites with substitution—on
126 each segment to obtain the distribution of local SSP density of the genome sequence pair
127 (Figure 1b). This algorithm considers segments instead of nucleotide / amino-acid sites as the
128 basic unit of a genome, because the local SSP density can be defined conveniently on
129 segments; an SSP can be a single nucleotide polymorphism (SNP) on a nucleotide sequence,
130 or a single amino acid polymorphism (SAP) on an amino acid sequence. It then fits the
131 empirical SSP distribution of all pairs to a model to infer the matrix of coalescent time of the
132 genome pairs (Figure 1d). While searching for the best fit model parameters and pairwise
133 coalescent times, it constrains these $n(n-1)/2$ coalescent times (Figure 1e), so that the matrix
134 can be bijectively mapped to a UPGMA tree that describes the phylogenetic inheritance of
135 the n genome sequences (Figure 1f).

136 **Model describing the evolution of single site polymorphism distributions of genome**
137 **pairs in a Fisher-Wright population**

138 The model behind CGP, which is used to fit the empirical distribution of local SSP (Figure
139 1d), is based on a Fisher-Wright haploid population with non-overlapping generations,
140 constant population size, and homologous recombination [20, 21]. In this framework, a node
141 in one generation inherits the genome of a random node in the previous generation, followed
142 by mutation and homologous recombination. A genome sequence is divided into L_{seg}
143 consecutive and non-overlapping segments, where every segment has length l_s (*i.e.*, consists
144 of l_s sites). The rate of mutation is μ per site per generation. The rate for a site to be covered
145 by a foreign DNA stretch attempting to recombine with the host genome is ρ per generation;
146 the rate for a segment to be covered by a recombination-attempting DNA stretch is also
147 approximately ρ , assuming that the segment is much shorter than the DNA stretch. Here,
148 $\rho = \rho_{ini} L$, where ρ_{ini} is the probability for a recombination-attempting foreign DNA stretch to
149 start at any given site, and L is the average length of the foreign DNA stretch. When a
150 recombination attempt happens on a segment, it either succeeds, and the foreign DNA
151 replaces the host DNA at the segment, or it fails. The success rate of an attempt is
152 approximately $\exp(-\delta/\delta_{TE})$, where δ is the divergence between the incoming DNA and the
153 host DNA, and δ_{TE} is the transfer efficiency, a constant that governs the success rate [5]. The
154 average site divergence in the population is denoted as θ , with $\theta = 2\mu N_e$ and population size N_e .

155 CGP's model [5, 6] considers the evolution of a SSP distribution between a pair of
156 genomes, X and Y. As the alignment of genome X and Y is divided into L_{seg} consecutive and
157 non-overlapping segments with l_s sites, let $f(x|t)$ be the distribution of segment divergence,
158 where $x=0, 1, \dots, l_s$ represents the number of SSPs on a segment of the XY alignment, $t \geq 0$ is
159 the (continuous) XY coalescent time, and $f(x|t)$ is normalized to unity by summing over x . To
160 save computational resources, we assume an upper bound $l_s^{cutoff} \leq l_s$ to x . At $t=0$, the most

161 recent common ancestor (MRCA) of XY splits into two lineages; initially, the two have
 162 identical genomes, and thus $f(x|0)=\delta_{x,0}$ (where $\delta_{x,0}$ is the Kronecker delta, *i.e.*, $f(x|0)$ is non-
 163 zero only at $x=0$). At $t>0$, mutations and recombinations occur, and the evolution of $f(x|t)$ is
 164 described by the following equation:

$$\frac{df(x|t)}{dt} = 2l_s\mu \sum_{y=0}^{l_s^{cutoff}} (M(x|y) - I(x|y))f(y|t) + 2\rho \sum_{y=0}^{l_s^{cutoff}} (P(x|y, \theta, \delta_{TE}, l_s) - I(x|y))f(y|t) \quad (1)$$

165 The first term of Eq (1) accounts for mutations on a segment— $M(x|y)$ models a
 166 mutation event, where a segment in the pair XY with y SSPs increases to $x=y+1$ SSPs during
 167 a mutation (*i.e.*, $M(x|y)=0$ for $x\neq y+1$); $I(x|y)$ is the identity matrix. For simplicity, we ignore
 168 back mutations.

169 The second term accounts for recombination— $P(x|y, \theta, \delta_{TE}, l_s)$ models a recombination
 170 event (see Eq. (S4) of *Dixit et al.* [5] or Eq. (3) of *Dixit et al.* [6] for a detailed derivation).
 171 Since a segment can recombine with its counterpart on another genome, it assumes that each
 172 segment of a genome, along with its counterparts in different genomes of the population,
 173 have their own phylogeny, which is detached from the genomes' phylogeny, and the segment
 174 population structure is approximated by the coalescent model. For an attempted
 175 recombination between Y and an external donor D, we can use the coalescent model to
 176 calculate the probability distribution for the segment divergence δ between D and X, and then
 177 obtain x from $x=l_s\delta$. As mutation and recombination can equally occur on either X or Y,
 178 there is a factor 2 attached to both terms. See Supplementary Text for the exact form of
 179 $P(x|y, \theta, \delta_{TE}, l_s)$. We solved Eq. (1) with boundary condition $f(x|0)=\delta_{x,0}$ to obtain the theoretical
 180 SSP distribution $f(x|t)$ at different coalescent times t .

181 We fit the theoretical distribution obtained with the CGP model to an empirical SSP
 182 distribution to infer the coalescent time of the pair of genomes. Let us consider an alignment

183 for a genome pair XY that is divided into L_{seg} segments, with empirical SSP distribution
 184 $g_{XY}(x)$ following the normalization condition:

$$\sum_{x=0}^{l_s^{cutoff}} g_{XY}(x) = L_{seg}$$

185 Let us denote the theoretical distribution as $f_{\mu,\rho,\theta,\delta_{TE}}(x|t)$, which is normalized to unity.
 186 The probability to observe the empirical distribution $g_{XY}(x)$ given the theoretical distribution
 187 $f_{\mu,\rho,\theta,\delta_{TE}}(x|t)$ is

$$\prod_x [f_{\mu,\rho,\theta,\delta_{TE}}(x|t)]^{g_{XY}(x)} \quad (2)$$

188 If we take the logarithm of this expression, it becomes the (negative) cross entropy between
 189 $g(x)$ and $f_{\mu,\rho,\theta,\delta_{TE}}(x|t)$ [22, 23]. The higher their similarity, the higher is this negative cross
 190 entropy; it attains its maximum when $f_{\mu,\rho,\theta,\delta_{TE}}(x|t)$ is equal to $g(x)$.

191 Suppose that we have n genomes ($X_i, i=1,\dots,n$), where their phylogeny is described by
 192 an ultrametric tree T ; the $n(n-1)/2$ pairwise SSP distributions have evolved according to the
 193 model with parameters $\mu, \rho, \theta, \delta_{TE}$. Let $t_T(X_a, X_b)$ be the coalescent time of X_a and X_b in the
 194 tree T . We use a score function, $S(X_1, X_2, \dots, X_n | \mu, \rho, \theta, \delta_{TE}, T)$, which is defined as the
 195 logarithm of the probability to observe the $n(n-1)/2$ empirical SSP distributions given the
 196 model and the tree, to quantify the model fit to the empirical SSP distributions. This score is
 197 the summation of the $n(n-1)/2$ negative mutual entropy terms:

$$S(X_1, \dots, X_n | \mu, \rho, \theta, \delta_{TE}, T) = \sum_{\text{all } (X_a, X_b) \text{ pairs}} \log \left\{ \prod_x [f_{\mu,\rho,\theta,\delta_{TE}}(x|t_T(X_a, X_b))]^{g_{X_a X_b}(x)} \right\} \quad (3)$$

198 Since the $n(n-1)/2$ SSP distributions are not completely independent of each other, Eq. (3) is
 199 not exactly a probability and so we call it a score. We developed an algorithm that samples
 200 the tree+model space and searches for the configuration with the maximum score using

201 Monte Carlo simulation with annealing and Metropolis acceptance (See Supplementary Text
202 for details).

203 **RESULTS**

204 **A coarse-graining approach to phylogenetic reconstruction**

205 Figure 1 gives a brief illustration on how the proposed coarse-graining phylogenetics (CGP)
206 algorithms fits the distribution of local single site polymorphisms (SSPs) density of the
207 genome pairs to infer their phylogenetic tree, forgoing the reconstruction of ancestral
208 recombination graph. In short, CGP is based on a mathematical model [5, 6] that
209 quantitatively describes the evolution of genomic sequence divergence; this model is
210 applicable to both nucleotide and amino acid sequences, and does not assume low
211 recombination rate. Recombination can introduce DNA stretches characterized by a high
212 density of substitutions, and the model considers substitution densities defined on the
213 genomic segments. A nucleotide genome sequence alignment (or a corresponding
214 concatenation of amino acid sequence alignments) is divided into a chain of consecutive,
215 non-overlapping segments, each with l_s sites; for a pair of genomes, the single site
216 polymorphisms (SSPs) on each segment are counted, resulting in an SSP distribution. CGP
217 takes the SSP distribution of every pair of considered genomes as input. The coalescent time
218 of two genomes can be inferred by fitting the CGP model to the empirical SSP distribution.
219 The ultrametric tree describing the vertical component of inheritance among n genomes can
220 be inferred from the coalescent times resulting from the fits to the $n(n-1)/2$ empirical SSP
221 distributions, implemented by the score function of Eq. (3). We developed the CGP
222 algorithm, which employs Monte Carlo simulation to sample the model+tree space,
223 identifying the tree and parameters that result in the highest score.

224 **Testing the CGP algorithm on simulated genomes**

225 We performed Fisher-Wright simulations with recombination to generate genome sequences,
226 allowing us to test different phylogenetic reconstruction algorithms. In the simulation, each
227 recombination-attempting DNA stretch starts at a random site of a genome, with equal chance
228 to be either a micro (geometric distribution, mean 100bp) or a macro stretch (geometric
229 distribution, mean 1kb). We used three sets of parameters that correspond to prokaryotic
230 populations with $r/m=2, 40, 80$ (r/m is the ratio between substitutions contributed by
231 mutations and by recombinations; these three settings are denoted as low, intermediate, and
232 high recombination level, respectively), and prepared the test groups, each with 4-10 genome
233 sequences. For comparison, r/m values observed in nature range from 0.02 to 63.6 [11].

234 The MRCA of a group of random genomes in a simulated population has an average
235 age close to the age of the population root node t_{root} . We would like to mimic the condition
236 where a single lineage diverges from the rest of the population and forms its own
237 subpopulation, so that the genomes in its subpopulation continue exchanging DNA among
238 themselves and with the rest outside. Hence, when picking genomes in the population to form
239 test groups, we constrained the age of the MRCA of the genomes in a test group, $t_{\text{test-group-root}}$,
240 to be $t_{\text{test-group-root}} \ll t_{\text{root}}$ (see Supplementary Text for the details and Supplementary File S3 for
241 the genome sequences in each test-group).

242 We applied CGP, as well as the previously published methods RaxML [24],
243 ClonalFrameML [14], and Gubbins [17] to the sequences of each test group. The RaxML
244 and Gubbins trees are midpoint-rooted. ClonalFrameML requires an initial tree as input and
245 we used the RaxML tree. CGP uses segment size $l_s=150$ and $l_s^{\text{cutoff}}=100$. We compared each
246 reconstructed tree with the true tree, measuring their unrooted symmetric distance (SD) [25],
247 as well as their rooted and unrooted branch score distance (BSD) [26] to quantify the
248 accuracy of the reconstructed phylogeny (see Supplementary File S1 for the these values); the

249 lower the unrooted SD / unrooted BSD / rooted BSD, the more accurate is the topology /
250 branch lengths / root positioning, respectively. We normalized the branch lengths of each tree
251 by its total branch length when calculating BSD.

252 CGP can predict the topology of a phylogeny of vertical inheritance as accurately as
253 the other algorithms. Figure 2 shows the histograms of unrooted SD; ClonalFrameML is
254 excluded as it uses the topology of RAxML trees. The distributions of SD of CGP are not
255 significantly different from the distributions of RAxML and Gubbins. Two-sided Wilcoxon
256 signed-rank tests (WSRT) at low, intermediate, and high recombination levels resulted in
257 $p=0.25$, 0.69 , 0.54 between CGP and RAxML, and $p=0.25$, 0.38 , 0.92 between CGP and
258 Gubbins.

259 Branch length predictions are more accurate with CGP than with alternative programs
260 at higher levels of recombination. Figure 3 plots the distributions of the unrooted BSD of
261 different algorithms and their z-scores; the unrooted BSD of trees reconstructed by different
262 algorithms on the same test group sequences are pooled together to calculate their z-scores to
263 help data visualization. The distribution of the unrooted BSD of CGP is significantly lower
264 than RAxML and Gubbins ($p < 10^{-10}$ at low, intermediate, and high recombination level
265 compared to both RAxML and Gubbins). The unrooted BSD of CGP is significantly lower
266 than ClonalFrameML except at low recombination levels ($p=0.76$, 2.2×10^{-7} , 10^{-7} at low,
267 intermediate, and high recombination levels).

268 CGP can perform accurate root positioning. Figure 4 plots the distribution of the
269 rooted BSD and their z-scores; the rooted BSD of trees reconstructed by different algorithms
270 on the same test group sequences are pooled together to calculate their z-scores. The
271 distribution of the rooted BSD of CGP is significantly lower than the other algorithms
272 ($p < 2 \times 10^{-17}$ at all recombination levels for CGP compared with the other three algorithms).

273 **Testing the CGP algorithm on real *E. coli* genomes**

274 We tested CGP, RAxML, ClonalFrameML, and Gubbins using *E. coli* and *Shigella* genome
275 sequences (see Supplementary Table S1 for their names); we refer to them as *E. coli*, as these
276 two species have intertwined phylogenies. We prepared test groups, each with 10 random
277 strains, where each strain is represented by a nucleotide and an amino acid sequence made
278 from its concatenated core genes (see Supplementary File S2 for the strains in each test
279 group, and also the 1,636 orthologous gene families of core genes; see Supplementary File S3
280 for their sequences). We applied CGP, RAxML, ClonalFrameML (with the topology from
281 RAxML trees), and Gubbins on the nucleotide sequences, and CGP and RAxML on amino
282 acid sequences; the RAxML and Gubbins trees are midpoint-rooted. CGP uses segment
283 length $l_s=150$ and $l_s^{cutoff}=100$ for nucleotide sequences, and $l_s=50$, $l_s^{cutoff}=50$ for amino acid
284 sequences.

285 To assess the accuracy of the phylogenetic trees reconstructed by the different
286 algorithms, we compared the reconstructed trees with the phylogenetic signal inferred from
287 the distribution of orthologous gene families in different genomes. We applied the GLOOME
288 algorithm [27], which considers the interior nodes of the tree as ancestral strains and
289 reconstructs their gene distribution; it takes a tree and the presence-and-absence of genes
290 across the extant strains as input, and performs a reconstruction of presence-and-absence of
291 genes in the ancestral strains based on the GLOOME posterior likelihood (GPL). We used
292 GPL as a score to quantify the accuracy of the tree fed into GLOOME; the more consistent
293 the phylogenetic signal from the gene distributions with a given tree, the higher the GPL (see
294 Supplementary File S2 for the GPL values of the reconstructed trees).

295 Figure 5 plots the distribution of the GPLs and the corresponding z-scores; the GPLs
296 of trees of the same test groups reconstructed by different methods are pooled together to
297 calculate the z-scores. Trees reconstructed from amino acid sequences have a higher GPL

298 than trees calculated from nucleotide sequences; moreover, CGP trees based on amino acids
299 are more accurate than trees calculated using RAxML ($p < 4 \times 10^{-14}$ when comparing CGP on
300 amino acid sequences with any other algorithm; other recombination-aware algorithms are
301 not applicable to amino acid sequences). Considering only trees reconstructed from
302 nucleotide sequences, the CGP trees generally have higher GPL than RAxML,
303 ClonalFrameML, and Gubbins trees ($p = 5.2 \times 10^{-4}$, 0.049, 1.4×10^{-15} , respectively).

304 **DISCUSSION**

305 We introduced a coarse-graining phylogenetic (CGP) model, which infers a phylogenetic tree
306 from the estimated pairwise coalescent times of genomes. We conducted extensive analyses
307 to compare the accuracy of the CGP algorithm with other state-of-the-art algorithms to
308 demonstrate its ability to reliably predict the topology, branch lengths, and root positioning of
309 phylogenetic trees. The CGP model does not rely on the assumption of low recombination
310 rates, which allows it to predict branch lengths accurately even if the vast majority of the
311 considered genome segments have experienced recombination on the timescale covered by
312 the phylogeny.

313 Analyses performed on the real *E. coli* genome sequences showed that trees
314 reconstructed from core genome amino acid sequences are more accurate, *i.e.*, more
315 consistent with the signal inferred from the distribution of genes in the extant genomes, than
316 trees calculated from nucleotide sequences. Amino acid sequences of core genes tend to
317 evolve more slowly than the corresponding DNA sequences, as these genes show dN/dS
318 values < 1 [28], and accordingly, the divergence of a pair of amino acid sequences is lower
319 than that of their nucleotide counterparts (see Supplementary File S3 for nucleotide and
320 amino acid sequence divergence between *E. coli* genome pairs). Thus, amino acid sequences

321 may be more “clonal” than nucleotide sequences and thus may provide more accurate
322 phylogenetic signals.

323 The major source of error of the CGP algorithm comes from the mismatch between
324 the genomic segments as basic unit of the algorithm and the genomic stretches affected by
325 homologous recombination, as the algorithm does not try to match the boundary of the
326 segments to the boundary of the actual recombination stretches. This mismatch gives rise to
327 segments that lie on the boundary and cover multiple recombination stretches, which
328 subsequently reduces the accuracy of the predictions of the algorithm. Hence, a possible
329 direction for further development is to find out the criteria to fine tune the segment size l_s so
330 as to minimize these boundary-overlapping segments; alternatively, we can improve the
331 theoretical model so that the segments do not have to be equal-sized and we can match the
332 segments to the recombination stretches.

333 The computational demand of the CGP algorithm is independent of sequence length,
334 as CGP considers only SSP distributions that are represented by a vector of $1+l_s^{cutoff}$ elements
335 in the computer code. Calculation of the CGP score (Eq. (3)) involves multiplication of
336 $(1+l_s^{cutoff}) \times (1+l_s^{cutoff})$ matrices; thus, the computational time scales as $O((l_s^{cutoff})^k)$, where $k \leq 3$
337 depends on the algorithm that carries out the matrix operations. When reconstructing a tree of
338 n genomes, the score calculation involves the summation over $n(n-1)/2$ pairs, making it scale
339 as $O(n^2)$. The segment size l_s affects the efficiency and accuracy of the algorithm. While a
340 smaller l_s leads to lower accuracy, increasing l_s leads to higher computational demand; a large
341 l_s combined with a small l_s^{cutoff} can also reduce the accuracy. Hence, one needs to set l_s and
342 l_s^{cutoff} carefully to balance the need for speed and accuracy.

343 The current algorithm that implements the CGP model is very simple; it should be
344 considered a proof of concept. While it makes use of Monte Carlo simulation to sample the
345 tree+parameter space, a hill-climbing method may be more efficient. Other possible

346 improvements involve better local search moves in the ultrametric tree space; one might even
347 drop the stringent ultrametricity constraint, and replace it with a more flexible matrix-tree
348 mapping method that allows a more efficient search in the tree space. The mutation matrix in
349 the current model can be improved to include back mutations and a more complex mutation
350 model. We leave these possible improvements to future studies.

351

352 **DECLARATIONS**

353 **Ethics approval and consent to participate**

354 Not applicable

355 **Consent for publication**

356 Not applicable

357 **Availability of data and materials**

358 All data generated or analysed during this study are included in this published article, its
359 supplementary information files, and GitHub repository ([https://github.com/TinPang/coarse-](https://github.com/TinPang/coarse-graining-phylogenetics)
360 [graining-phylogenetics](https://github.com/TinPang/coarse-graining-phylogenetics)).

361 **Competing interests**

362 The authors declare that they have no competing interests

363 **Funding**

364 This work was supported by the German Research Foundation (CRC 680 and CRC 1310).

365 **Authors' contributions**

366 Not applicable

367 **Acknowledgements**

368 We would like to thank Martin Lercher and Arndt von Haeseler for helpful comments and
369 advice.

370

371

372

373 **SUPPLEMENTARY MATERIALS**

374 1. Source code of the CGP algorithm: [https://github.com/TinPang/coarse-graining-](https://github.com/TinPang/coarse-graining-phylogenetics)
375 [phylogenetics](https://github.com/TinPang/coarse-graining-phylogenetics)

376 2. Supplementary Text, Figures and Tables.

377 3. Supplementary File S1: Analyses of the simulated genomes: the symmetric distance
378 (SD) and branch score distance (BSD) between the reconstructed trees and the true
379 trees.

380 4. Supplementary File S2: Analyses of the real genomes: b-number of the *E. coli* core
381 genes used to make the 'super-gene' sequences, strains in each test-group, and also
382 the GLOOME posterior likelihood (GPL) of the reconstructed trees.

383 5. Supplementary File S3: Sequences and their phylogenetic trees: sequences of the
384 simulated genomes in different test-groups, their true trees and also the phylogenetic
385 trees reconstructed by different algorithms; sequences of the *E. coli* genomes and their
386 trees reconstructed by different algorithms; genes to orthologous gene families map
387 provided by ProteinORTHO.

389 **REFERENCES**

- 390 1. Pál C, Papp B, Lercher MJ. Adaptive evolution of bacterial metabolic networks by
391 horizontal gene transfer. *Nat Genet.* 2005;37:1372–5.
- 392 2. Ochman H, Lawrence JG, Groisman EA. Lateral gene transfer and the nature of bacterial
393 innovation. *Nature.* 2000;405:299–304.
- 394 3. Lang AS, Zhaxybayeva O, Beatty JT. Gene transfer agents: phage-like elements of genetic
395 exchange. *Nat Rev Microbiol.* 2012;10:472–82.
- 396 4. Wilson GG, Murray NE. Restriction and Modification Systems. *Annu Rev Genet.*
397 1991;25:585–627.
- 398 5. Dixit PD, Pang TY, Studier FW, Maslov S. Recombinant transfer in the basic genome of
399 *Escherichia coli*. *Proc Natl Acad Sci U S A.* 2015;112:9070–5.
- 400 6. Dixit PD, Pang TY, Maslov S. Recombination-Driven Genome Evolution and Stability of
401 Bacterial Species. *Genetics.* 2017;;genetics.300061.2017.
- 402 7. Huddleston JR. Horizontal gene transfer in the human gastrointestinal tract: potential
403 spread of antibiotic resistance genes. *Infect Drug Resist.* 2014;7:167–76.
- 404 8. Takeuchi N, Kaneko K, Koonin E. Horizontal Gene Transfer Can Rescue Prokaryotes
405 from Muller’s Ratchet: Benefit of DNA from Dead Cells and Population Subdivision. *G3*
406 *GenesGenomesGenetics.* 2014;4:325–39.
- 407 9. Croucher NJ, Mostowy R, Wymant C, Turner P, Bentley SD, Fraser C. Horizontal DNA
408 Transfer Mechanisms of Bacteria as Weapons of Intragenomic Conflict. *PLOS Biol.*
409 2016;14:e1002394.
- 410 10. Spratt BG. Multilocus sequence typing: molecular typing of bacterial pathogens in an era
411 of rapid DNA sequencing and the internet. *Curr Opin Microbiol.* 1999;2:312–6.
- 412 11. Vos M, Didelot X. A comparison of homologous recombination rates in bacteria and
413 archaea. *ISME J.* 2009;3:199–208.
- 414 12. Schierup MH, Hein J. Consequences of recombination on traditional phylogenetic
415 analysis. *Genetics.* 2000;156:879–91.
- 416 13. Didelot X, Falush D. Inference of Bacterial Microevolution Using Multilocus Sequence
417 Data. *Genetics.* 2007;175:1251–66.
- 418 14. Didelot X, Wilson DJ. ClonalFrameML: Efficient Inference of Recombination in Whole
419 Bacterial Genomes. *PLOS Comput Biol.* 2015;11:e1004041.
- 420 15. Didelot X, Lawson D, Darling A, Falush D. Inference of Homologous Recombination in
421 Bacteria Using Whole-Genome Sequences. *Genetics.* 2010;186:1435–49.

- 422 16. Vaughan TG, Welch D, Drummond AJ, Biggs PJ, George T, French NP. Inferring
423 Ancestral Recombination Graphs from Bacterial Genomic Data. *Genetics*. 2017;205:857–70.
- 424 17. Croucher NJ, Page AJ, Connor TR, Delaney AJ, Keane JA, Bentley SD, et al. Rapid
425 phylogenetic analysis of large samples of recombinant bacterial whole genome sequences
426 using Gubbins. *Nucleic Acids Res*. 2015;43:e15–e15.
- 427 18. Marttinen P, Hanage WP, Croucher NJ, Connor TR, Harris SR, Bentley SD, et al.
428 Detection of recombination events in bacterial genomes from large population samples.
429 *Nucleic Acids Res*. 2012;40:e6.
- 430 19. Mostowy R, Croucher NJ, Andam CP, Corander J, Hanage WP, Marttinen P. Efficient
431 Inference of Recent and Ancestral Recombination within Bacterial Populations. *Mol Biol*
432 *Evol*. 2017;34:1167–82.
- 433 20. Kingman JFC. Origins of the Coalescent: 1974-1982. *Genetics*. 2000;156:1461–3.
- 434 21. Fraser C, Hanage WP, Spratt BG. Recombination and the nature of bacterial speciation.
435 *Science*. 2007;315:476–80.
- 436 22. Rubinstein RY, Kroese DP. *The Cross-Entropy Method: A Unified Approach to*
437 *Combinatorial Optimization, Monte-Carlo Simulation and Machine Learning*. Springer
438 Science & Business Media; 2004.
- 439 23. Boer P-T de, Kroese DP, Mannor S, Rubinstein RY. A Tutorial on the Cross-Entropy
440 Method. *Ann Oper Res*. 2005;134:19–67.
- 441 24. Stamatakis A. RAxML Version 8: A tool for Phylogenetic Analysis and Post-Analysis of
442 Large Phylogenies. *Bioinformatics*. 2014;:btu033.
- 443 25. Robinson DF, Foulds LR. Comparison of phylogenetic trees. *Math Biosci*. 1981;53:131–
444 47.
- 445 26. Kuhner MK, Felsenstein J. A simulation comparison of phylogeny algorithms under
446 equal and unequal evolutionary rates. *Mol Biol Evol*. 1994;11:459–68.
- 447 27. Cohen O, Ashkenazy H, Belinky F, Huchon D, Pupko T. GLOOME: gain loss mapping
448 engine. *Bioinformatics*. 2010;26:2914–5.
- 449 28. Lapierre M, Blin C, Lambert A, Achaz G, Rocha EPC. The Impact of Selection, Gene
450 Conversion, and Biased Sampling on the Assessment of Microbial Demography. *Mol Biol*
451 *Evol*. 2016;33:1711–25.

452

453

454

455

456 **FIGURES CAPTIONS**

457 **Figure 1.** Illustration of the procedure of the proposed CGP algorithm. (a) The algorithm
458 takes n aligned sequences as input, which can be nucleotide or amino acid sequences;
459 substitutions on the sequences are represented by coloured markers. (b) Each of the $n(n-1)/2$
460 genome pairs is divided into equal-sized segments, and the pairwise substitutions on each
461 segment is enumerated to obtain the distribution of local single site polymorphisms (SSPs)
462 density (denoted as $g(x)$). (c) The algorithm aims to infer the distance matrix of the genome
463 sequence pairs from the $n(n-1)/2$ SSP distributions. (d) In particular, the algorithm fits the
464 empirical SSP distributions with a model; the input of this model involves a matrix of $n(n-$
465 $1)/2$ coalescent times and other model parameters (mutation rate μ , recombination rate ρ ,
466 average population divergence θ and transfer efficiency δ_{TE}). (e) In the fitting process, the
467 $n(n-1)/2$ coalescent times are constrained (matrix cells with the same colour have the same
468 value), such that the matrix can be bijectively mapped to a UPGMA tree. (f) the algorithm
469 explores the model parameter space and tree space to obtain the best fit ultrametric tree.

470 **Figure 2.** Histograms showing the distributions of the unrooted symmetric distance (SD)
471 between true trees and trees reconstructed by CGP, RaxML, and Gubbins, from genome
472 sequences derived from Fisher-Wright simulations with low, intermediate, and high
473 recombination levels.

474 **Figure 3.** Boxplots showing the distributions of the unrooted branch score distance (BSD)
475 and the distributions of the z-score of unrooted BSD between the true trees and trees
476 reconstructed by CGP, RaxML, and Gubbins, from genome sequences from Fisher-Wright
477 simulations with low, intermediate, and high recombination levels.

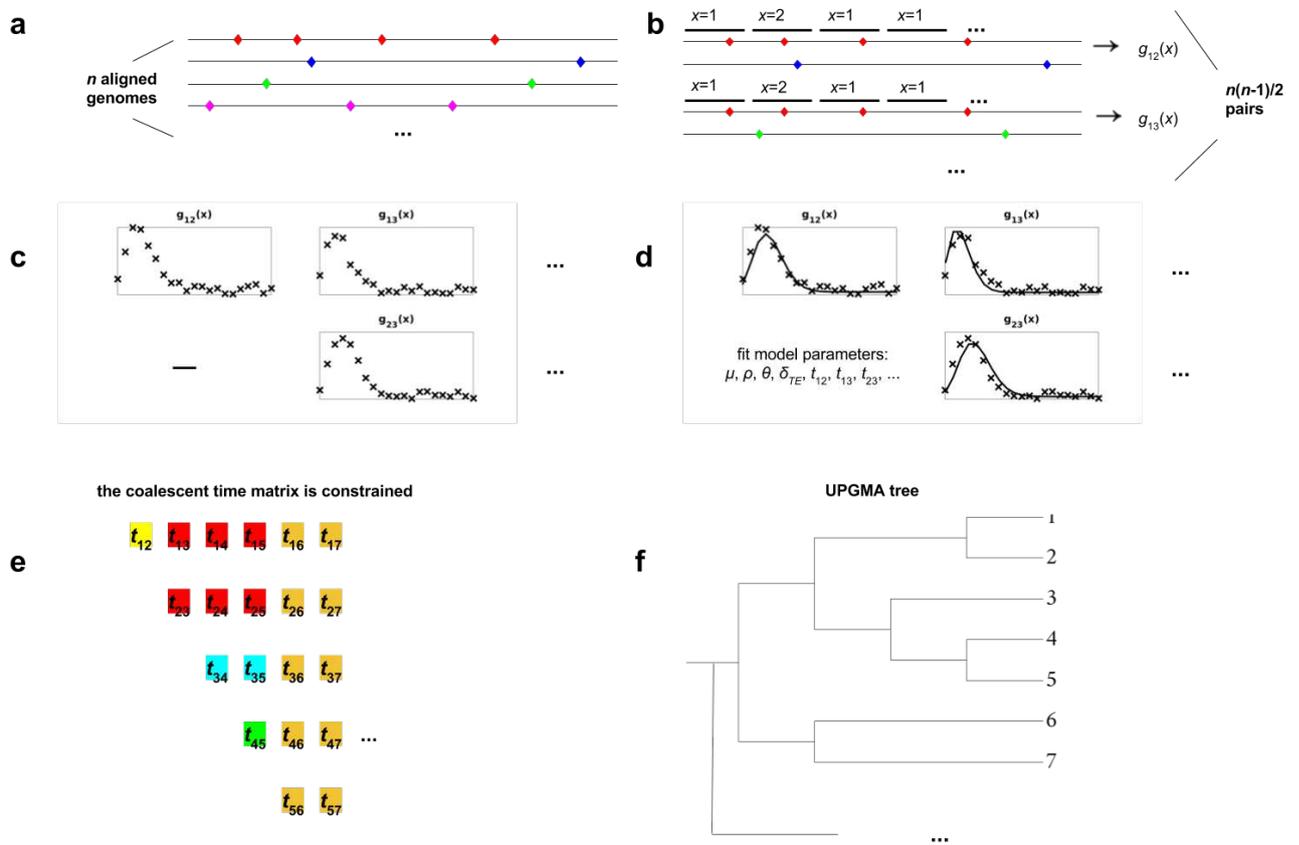
478 **Figure 4.** Boxplots showing the distributions of the rooted branch score distance (BSD) and
479 the distributions of the z-score of rooted BSD between the true trees and trees reconstructed

480 by CGP, RaxML, and Gubbins, from genome sequences from Fisher-Wright simulations with
 481 low, intermediate, and high recombination levels.

482 **Figure 5.** Boxplots showing the distributions of the GLOOME posterior likelihood (GPL)
 483 and the distributions of the z-scores of GPL. The trees were reconstructed from observed *E.*
 484 *coli* genome sequences, applying CGP, RAxML, ClonalFrameML, and Gubbins to nucleotide
 485 sequences, and CGP and RAxML to amino acid sequences.

486 **FIGURES**

487

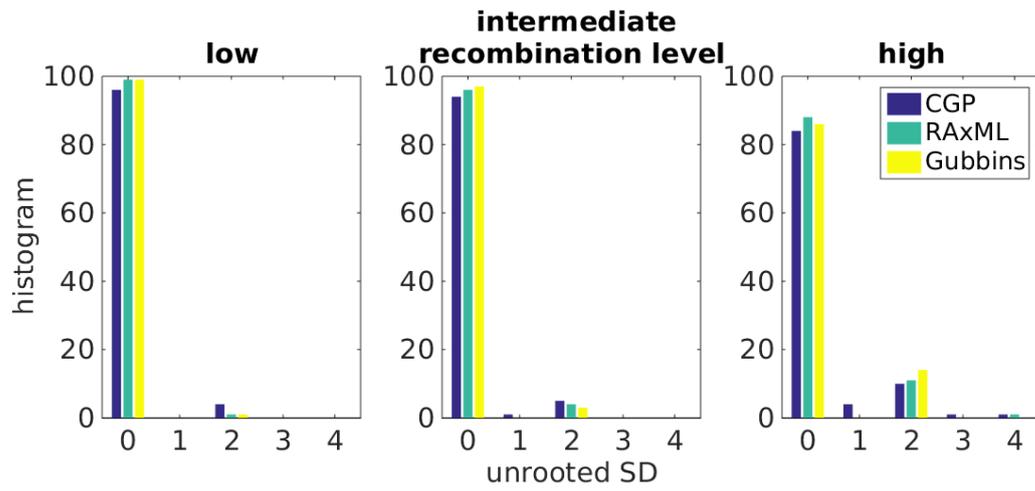


488

489 **Figure 3.**

490

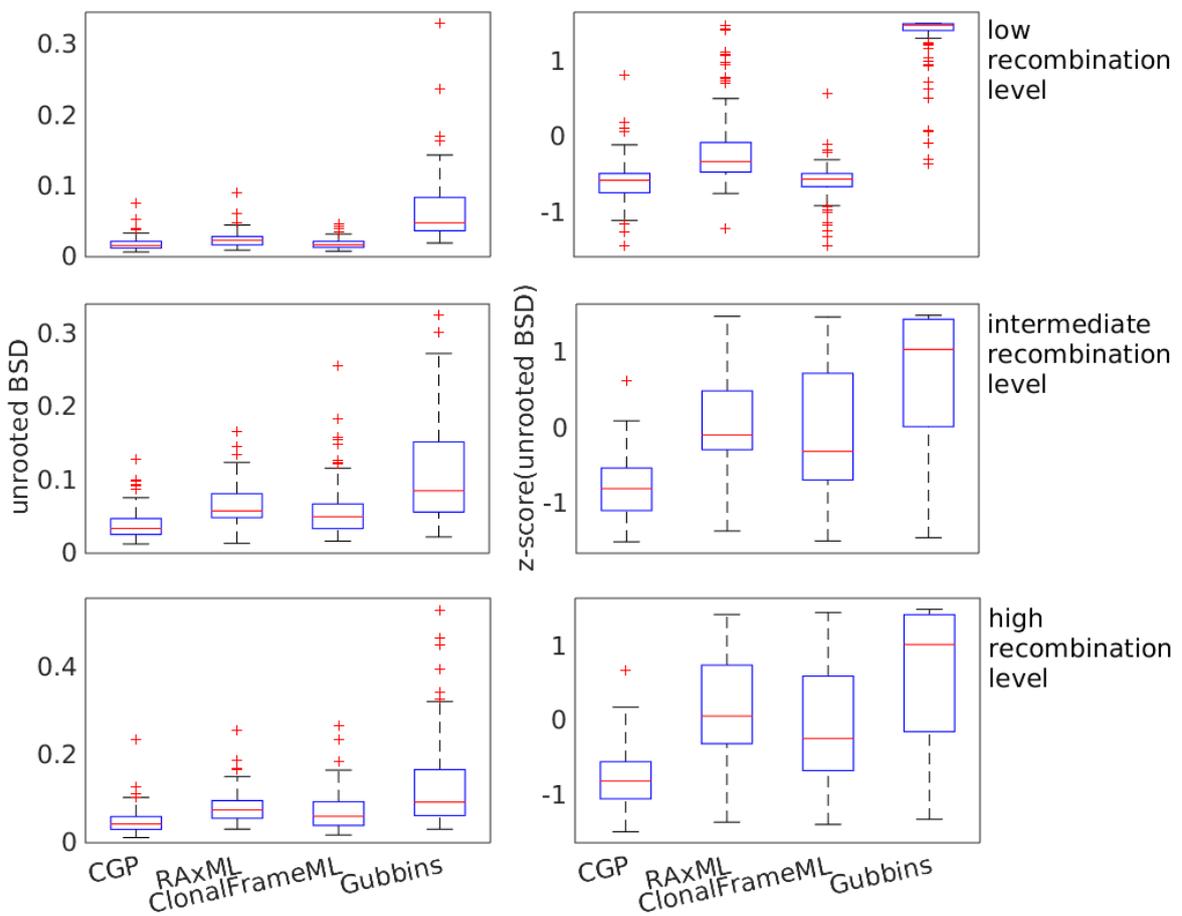
491



492

493 **Figure 2.**

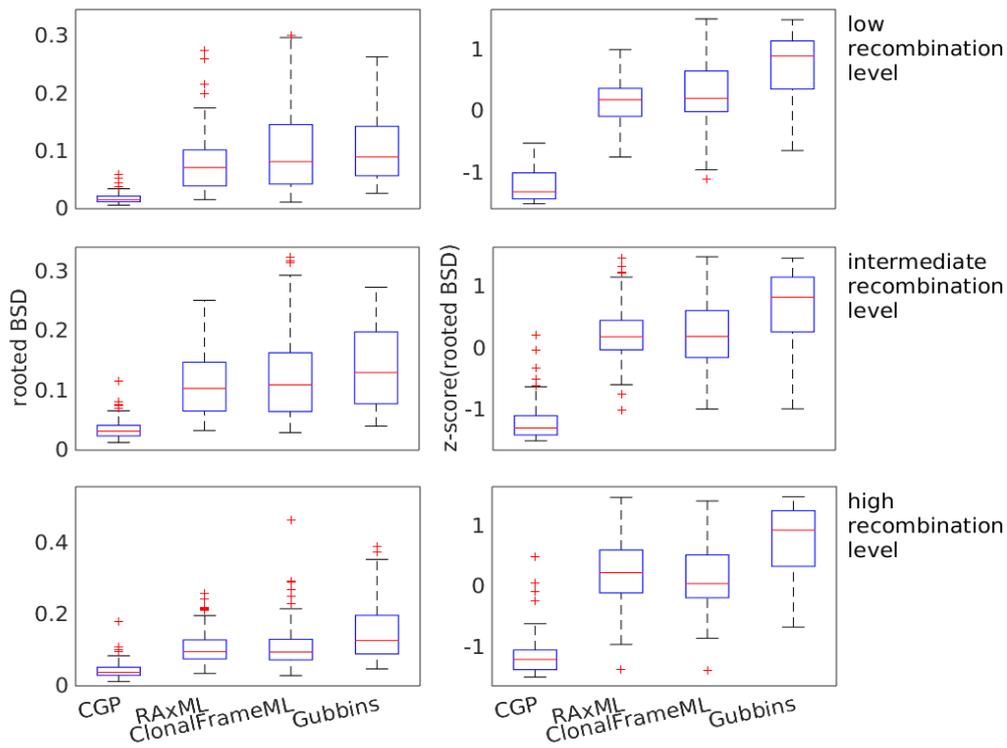
494



495

496 **Figure 3.**

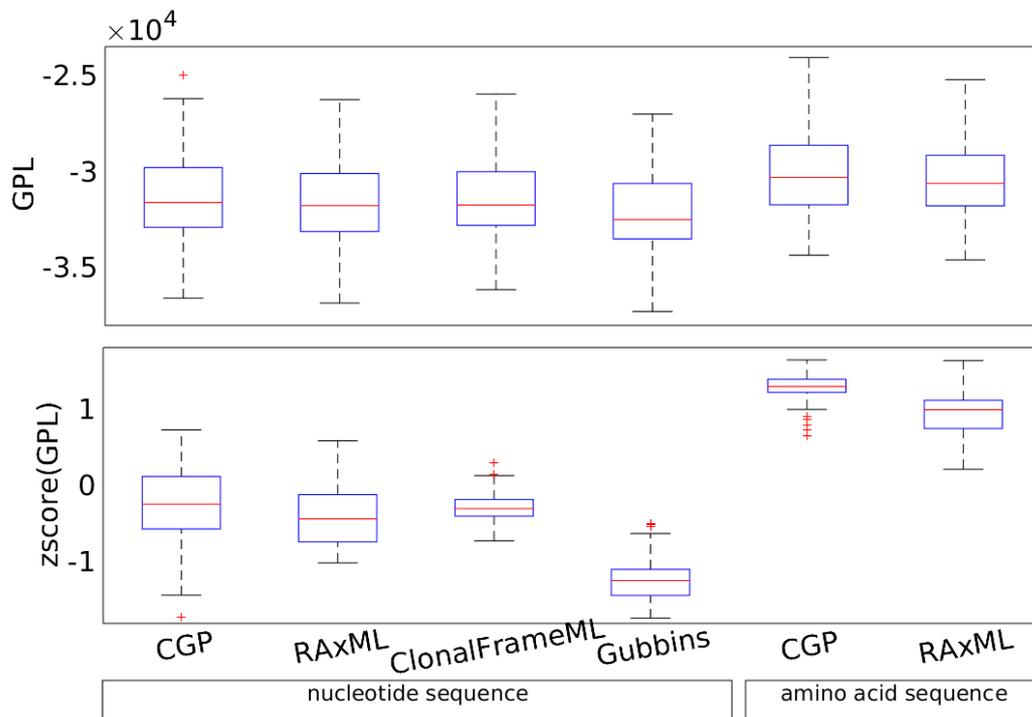
497



498

499 **Figure 4.**

500



501

502 **Figure 5.**

Figures

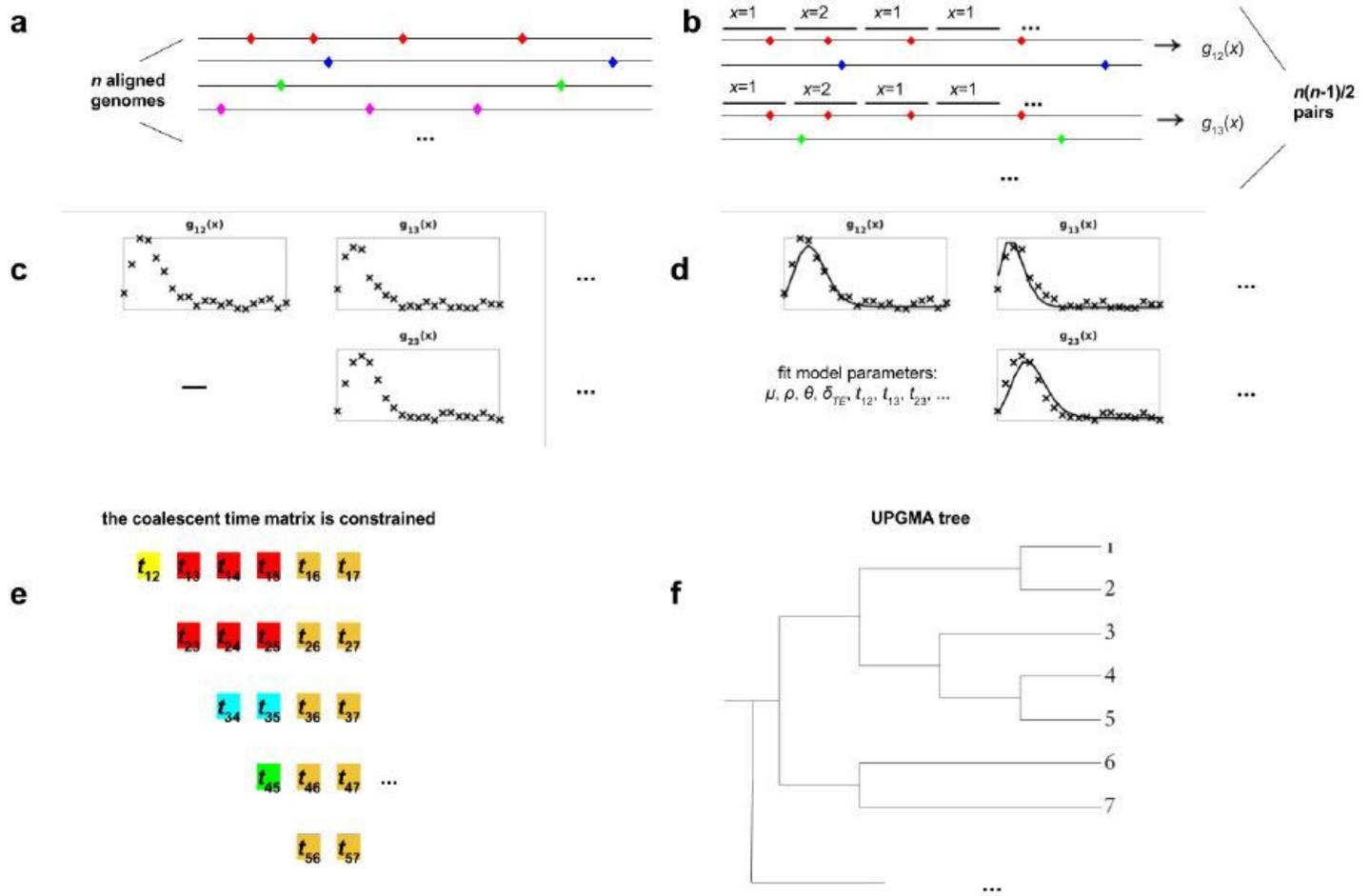


Figure 1

Illustration of the procedure of the proposed CGP algorithm. (a) The algorithm takes n aligned sequences as input, which can be nucleotide or amino acid sequences; substitutions on the sequences are represented by coloured markers. (b) Each of the $n(n-1)/2$ genome pairs is divided into equal-sized segments, and the pairwise substitutions on each segment is enumerated to obtain the distribution of local single site polymorphisms (SSPs) density (denoted as $g(x)$). (c) The algorithm aims to infer the distance matrix of the genome sequence pairs from the $n(n-1)/2$ SSP distributions. (d) In particular, the algorithm fits the empirical SSP distributions with a model; the input of this model involves a matrix of $n(n-1)/2$ coalescent times and other model parameters (mutation rate μ , recombination rate ρ , average population divergence θ and transfer efficiency δ_{TE}). (e) In the fitting process, the $n(n-1)/2$ coalescent times are constrained (matrix cells with the same colour have the same 467 value), such that the matrix can be bijectively mapped to a UPGMA tree. (f) the algorithm 468 explores the model parameter space and tree space to obtain the best fit ultrametric tree.

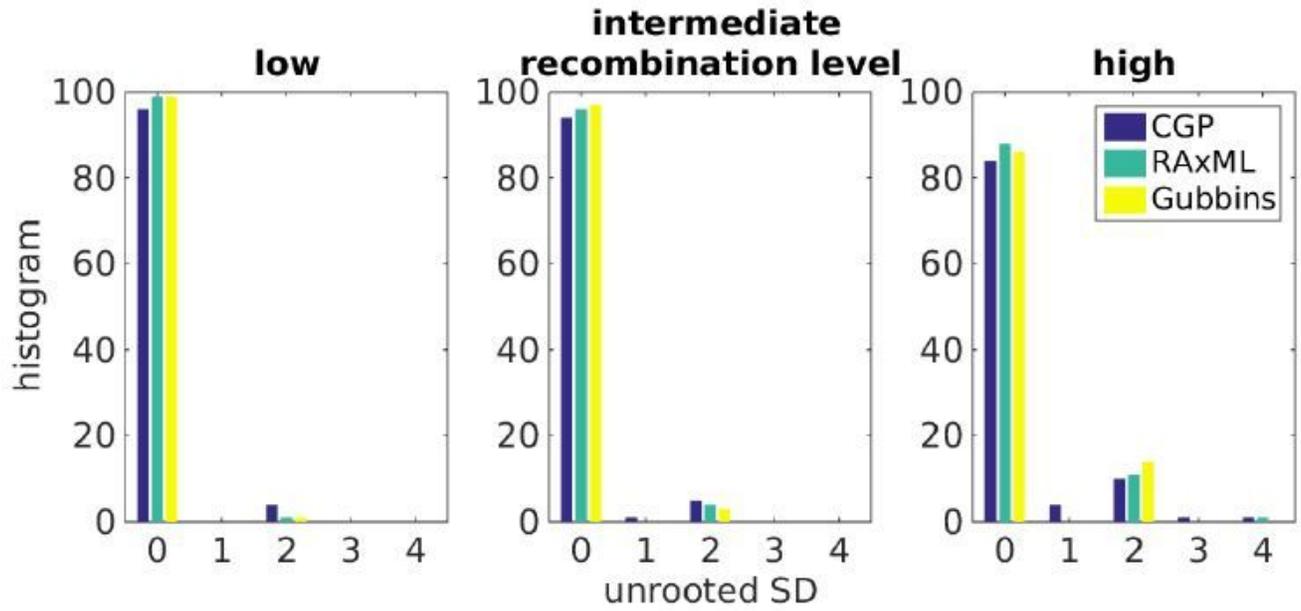


Figure 2

Histograms showing the distributions of the unrooted symmetric distance (SD) between true trees and trees reconstructed by CGP, RAxML, and Gubbins, from genome sequences derived from Fisher-Wright simulations with low, intermediate, and high recombination levels.

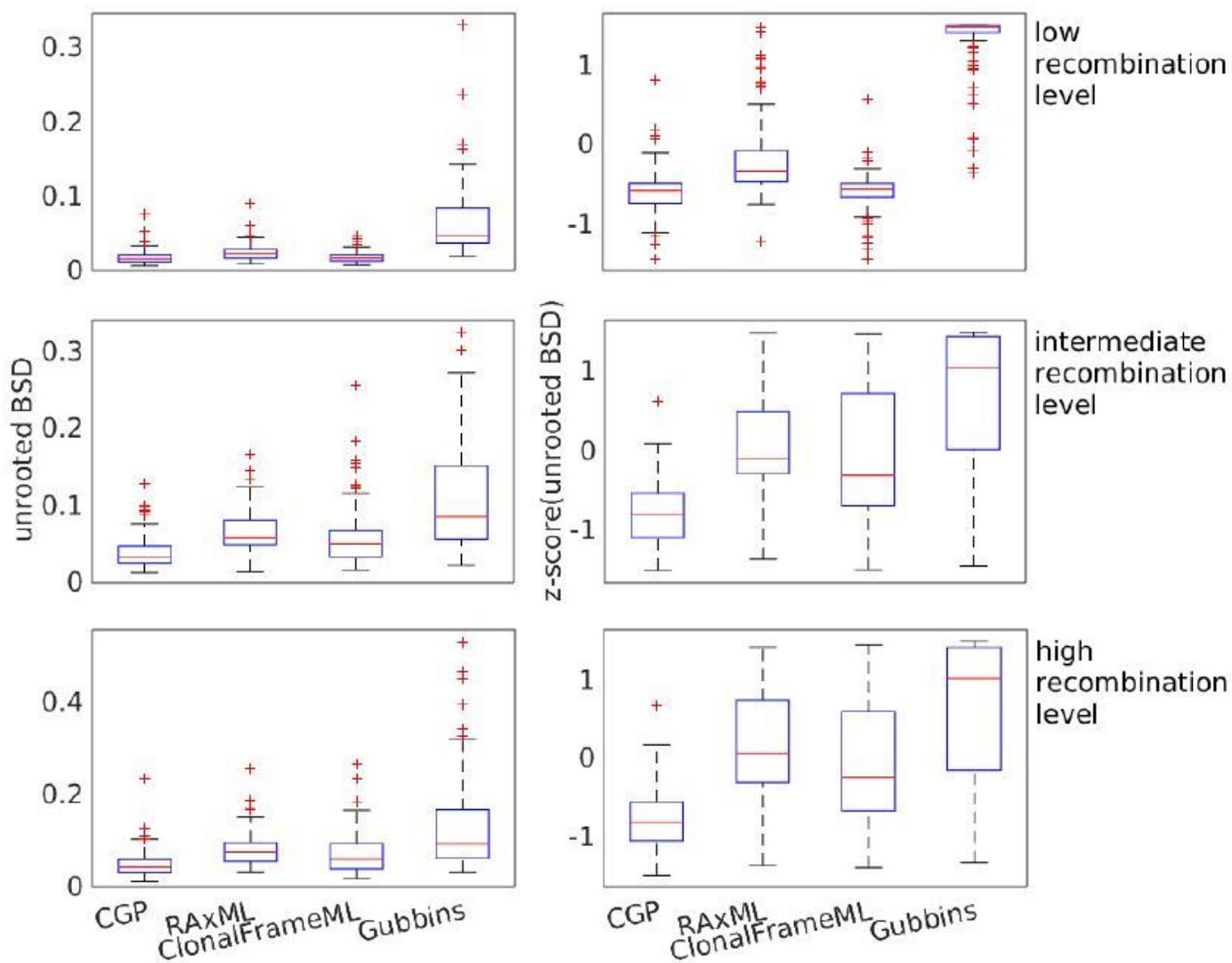


Figure 3

Boxplots showing the distributions of the unrooted branch score distance (BSD) and the distributions of the z-score of unrooted BSD between the true trees and trees reconstructed by CGP, RaxML, and Gubbins, from genome sequences from Fisher-Wright simulations with low, intermediate, and high recombination levels.

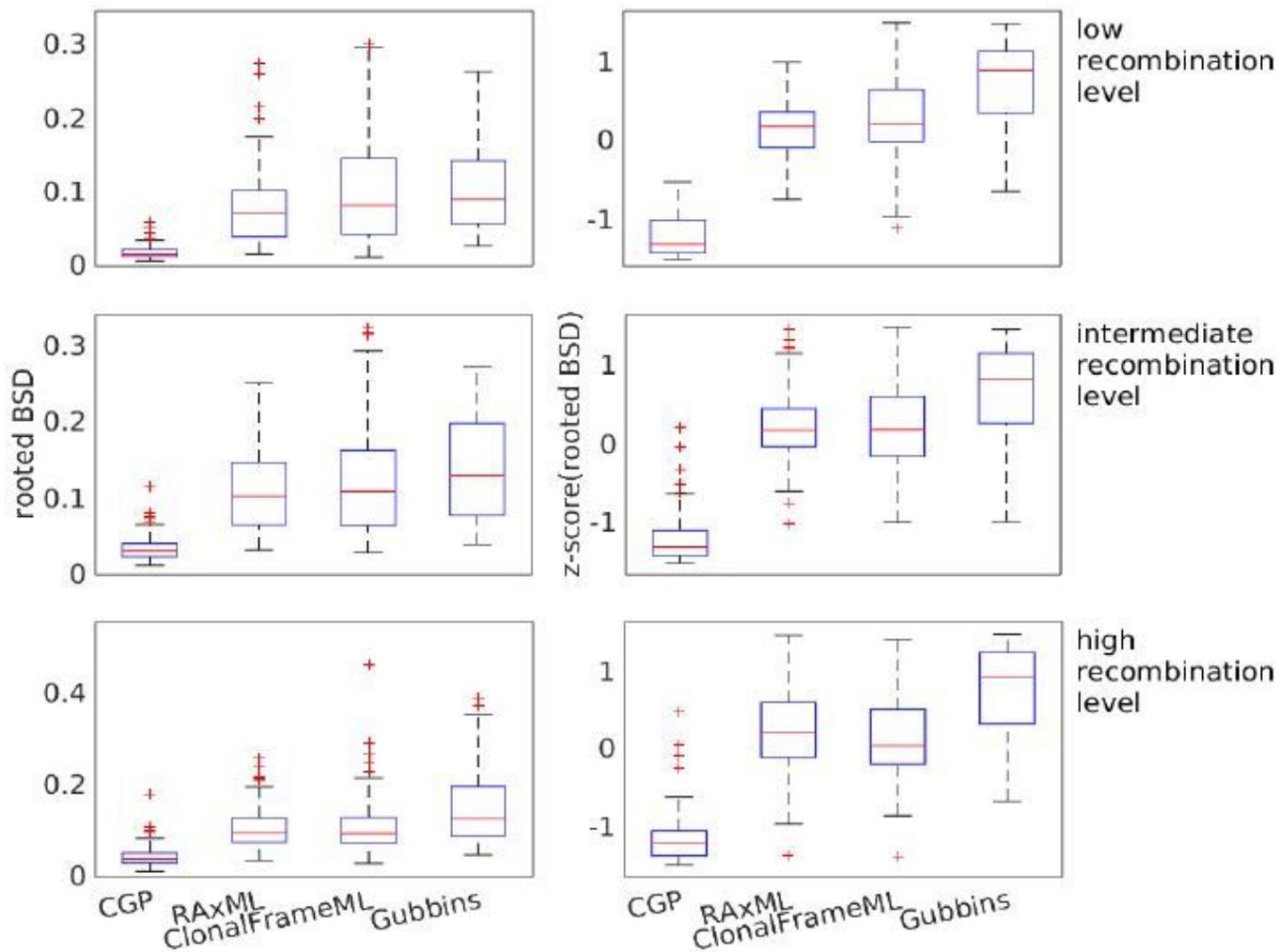


Figure 4

Boxplots showing the distributions of the rooted branch score distance (BSD) and the distributions of the z-score of rooted BSD between the true trees and trees reconstructed by CGP, RaxML, and Gubbins, from genome sequences from Fisher-Wright simulations with 480 low, intermediate, and high recombination levels.

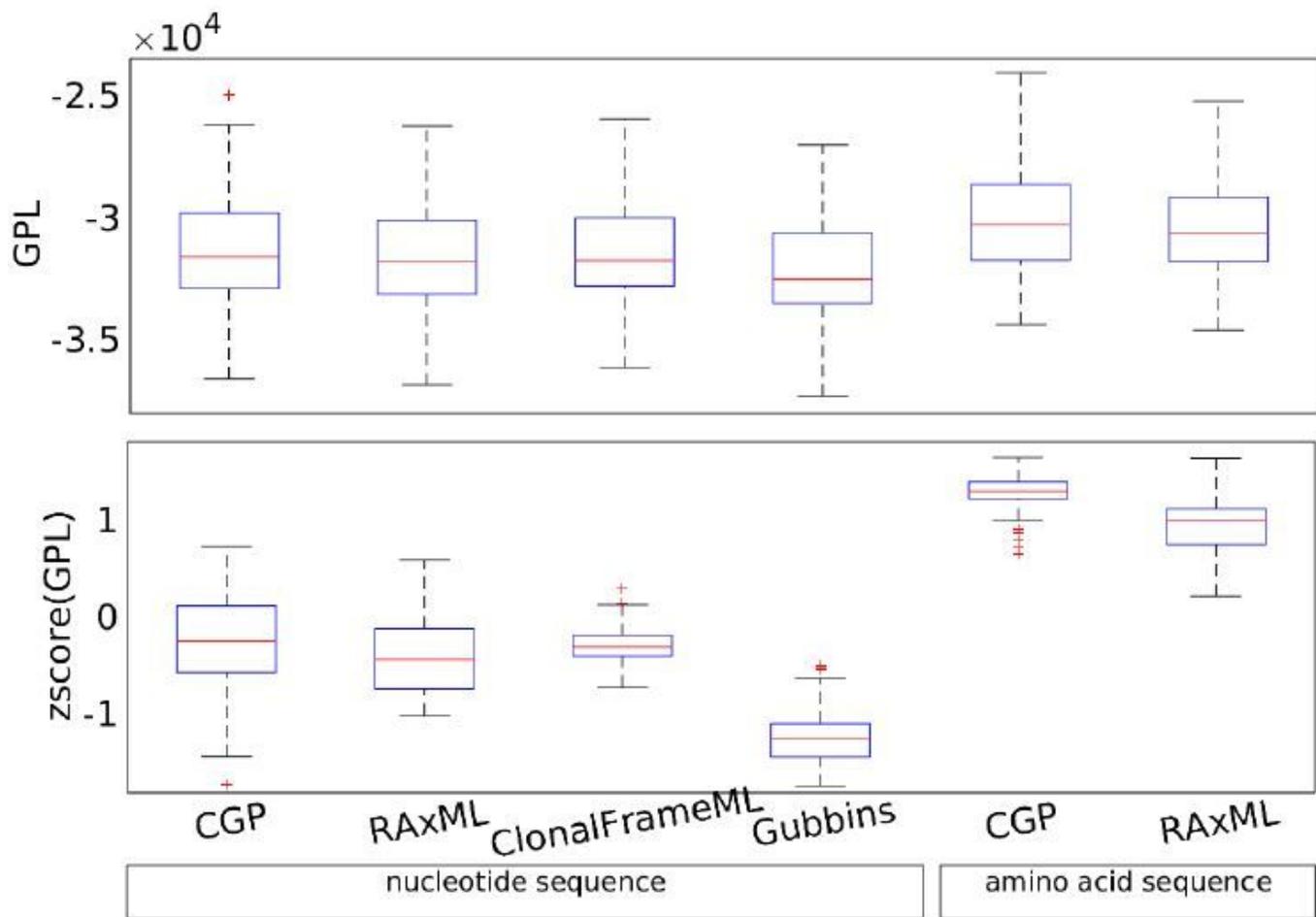


Figure 5

Boxplots showing the distributions of the GLOOME posterior likelihood (GPL) and the distributions of the z-scores of GPL. The trees were reconstructed from observed *E. coli* genome sequences, applying CGP, RAxML, ClonalFrameML, and Gubbins to nucleotide sequences, and CGP and RAxML to amino acid sequences.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplementaryFileS3.tar.xz](#)
- [supplementarytext.pdf](#)
- [SupplementaryFileS2resultsrealgenomes.xlsx](#)
- [SupplementaryFileS1resultssimulatedgenomes.xlsx](#)