

1 **Supplementary Methods**

2

3 **1. Simulating recombining bacterial populations**

4 Multiple sequence alignments and their corresponding clonal frames were jointly
5 simulated 1000 times under a neutral coalescent model with bacterial
6 recombination and a Jukes-Cantor model of nucleotide substitution using
7 SimMLST (1). Each population comprised 100 genome sequences of 1 million
8 base pairs in length, which we partitioned into 1000 loci of equal length for
9 computational efficiency. The population-scaled mutation rate, $\theta = 2N_e u$ (where
10 N_e is the effective population size and u the mutation rate per site per generation)
11 was fixed to 1%, a typical value for many bacterial species (2) and the average
12 recombination tract length to 500 base pairs, similar to estimates from several
13 species (Fearnhead et al. 2005; Jolley et al. 2005; Kennemann et al. 2011;
14 Everitt et al. 2014). The population-scaled recombination rate, $\rho = 2N_e r$ (where r
15 is the rate of initiation of recombination per site per generation), was fixed to 0%,
16 0.1% or 1%. At $\rho = 1\%$, recombination events are initiated as often as mutation
17 events, but the overall effect of recombination on the substitution process, known
18 as r/m , is greater than that of mutation ($r/m = 5$) because each recombination
19 event affects many sites. Therefore the range of recombination rates investigated
20 encompasses those seen in the majority of bacteria, with the notable exception
21 of extremely promiscuous species such as *Helicobacter pylori*, *Streptococcus*
22 *pneumoniae* and *Salmonella enterica* (7).

23

24 Since recombination has been shown to produce spurious signals of exponential
25 growth in phylogenies (8), we also studied the effect of recombination on branch
26 accuracy in populations with different rates of growth. The growth rate parameter,
27 $g = N_e m$ (where m is the exponential growth rate per generation), was
28 investigated across a range of values of $g = 0, 1$ and 10 , covering both low and
29 high growth rates. Since exponential growth reduces the total number of
30 substitutions across the tree, data simulated with the same value of θ under
31 higher growth rates will comprise fewer substitutions. To make a fairer
32 comparison we scaled θ in order to maintain the same expected number of
33 mutations per tree across all demographic models. This was achieved by
34 simulating the ratio of the average tree length constructed under a model of
35 constant population size and one of exponential growth ($g = 1$ and $g = 10$).

36

37 **2. Phylogenetic tree construction**

38 Phylogenetic trees were constructed for each simulated dataset using the
39 distance-based neighbor joining (NJ) and UPGMA methods, maximum likelihood
40 (ML), and BEAST, which is a Bayesian inference method. In each analysis, a
41 Jukes-Cantor (JC) model of nucleotide substitution was used (9). ML trees were
42 constructed using PhyML with the following command line arguments: `-m HKY85`
43 `-v 0 -t 1 -f 0.25,0.25,0.25,0.25 -c 1 -s BEST -b 100` (10). Bayesian phylogenetic
44 trees were constructed in BEAST v.1.7.5 (Bayesian Evolutionary Analysis by
45 Sampling Trees), using a strict molecular clock (uniform prior) and exponential
46 growth model (populations size fixed to 1.0) (XML available on request) (11, 12).

47 Since ML was shown to reconstruct the most accurate tree topology, this tree
48 was used as a starting tree in each BEAST analysis. This required midpoint
49 rooting of the tree and rescaling of terminal branches, such that tip heights were
50 zero (a requirement of starting trees for exponential growth models in BEAST).
51 Two independent Markov chain Monte Carlo (MCMC) chains were run for 10
52 million steps each, which provided sufficient mixing and convergence to the
53 stationary distribution. Parameters and trees from both runs were sampled every
54 1000 steps and combined using LogCombiner. Model parameters were
55 summarized using LogAnalyser. NJ and UPGMA trees were constructed using
56 the APE and phangorn libraries in R respectively (13, 14).

57

58 Bootstrapping of ML trees was performed in PhyML and of NJ and UPGMA trees
59 in R, using 100 replicates in each case. Posterior probabilities of branches in
60 BEAST trees were calculated by constructing the maximum clade credibility
61 (MCC) tree for each distribution in TreeAnnotator (11, 12) .

62

63 **3. Calculating tree topology accuracy of estimated trees**

64 The accuracy of tree topology was calculated using the Robinson-Foulds
65 Symmetric Difference metric (15) between the clonal frame and reconstructed
66 tree. This was used to obtain the proportion of branches in the clonal frame
67 correctly reconstructed in the estimated tree, i.e. accuracy = (total number of
68 branches – (Symmetric Distance/2))/total number of branches. Unrooted trees
69 were used for each comparison and the accuracy for each model was averaged

70 over 1000 simulated datasets. The accuracy of each posterior distribution of
71 BEAST trees was quantified as the average accuracy across 1000 trees in the
72 distribution, which did not differ from those estimated using the MCC tree.

73

74 In order to investigate how accuracy of branches varied by branch length, each
75 branch in a tree was assigned to one of three intervals based on its length.
76 Intervals were defined so that the number of branches per interval was
77 approximately equal (mean: 65.3). The average accuracy of branches within
78 each interval is plotted in Figure S2.

79

80 The accuracy of bootstrap values (ML, NJ, UPGMA) and posterior probabilities
81 (BEAST) was measured as the mean proportion of correctly estimated branches
82 within each of ten intervals of branch support value.

83

84 **4. Removal of homoplasious sites from sequence alignment**

85 Recombination events within a population can give rise to homoplasies across
86 the phylogenetic tree. In order to remove all substitutions arising from
87 recombination, all sites at which a homoplasy had occurred were removed from
88 the alignment. Homoplasies were identified using maximum likelihood ancestral
89 state reconstruction (16) to reconstruct the sequences at internal nodes in the ML
90 tree and then counting the number of times each substitution arose on the tree. A
91 homoplasious site is defined as one at which the minimum number of
92 substitutions needed to explain the observed number of alleles is exceeded.

93 Supplementary References

- 94 1. **Didelot X, Lawson D, Falush D.** 2009. SimMLST: Simulation of multi-
95 locus sequence typing data under a neutral model. *Bioinformatics*
96 **25**:1442–1444.
- 97 2. **Pérez-Losada M, Browne EB, Madsen A, Wirth T, Viscidi RP, Crandall**
98 **KA.** 2006. Population genetics of microbial pathogens estimated from
99 multilocus sequence typing (MLST) data. *Infect. Genet. Evol.* **6**:97–112.
- 100 3. **Jolley KA, Wilson DJ, Kriz P, McVean G, Maiden MCJ.** 2005. The
101 influence of mutation, recombination, population history, and selection on
102 patterns of genetic diversity in *Neisseria meningitidis*. *Mol. Biol. Evol.*
103 **22**:562–569.
- 104 4. **Fearnhead P, Smith NGC, Barrigas M, Fox A, French N.** 2005. Analysis
105 of recombination in *Campylobacter jejuni* from MLST population data. *J.*
106 *Mol. Evol.* **61**:333–340.
- 107 5. **Kennemann L, Didelot X, Aebischer T, Kuhn S, Drescher B, Droege M,**
108 **Reinhardt R, Correa P, Meyer TF, Josenhans C, Falush D, Suerbaum**
109 **S.** 2011. *Helicobacter pylori* genome evolution during human infection.
110 *Proc. Natl. Acad. Sci. U. S. A.* **108**:5033–5038.
- 111 6. **Everitt RG, Didelot X, Batty EM, Miller RR, Knox K, Young BC,**
112 **Bowden R, Auton A, Votintseva A, Larner-Svensson H, Charlesworth**
113 **J, Golubchik T, Ip CLC, Godwin H, Fung R, Peto TE a., Walker a. S,**
114 **Crook DW, Wilson DJ.** 2014. Mobile elements drive recombination
115 hotspots in the core genome of *Staphylococcus aureus*. *Nat. Commun.*
116 **5**:3956.
- 117 7. **Vos M, Didelot X.** 2009. A comparison of homologous recombination rates
118 in bacteria and archaea. *ISME J.* **3**:199–208.
- 119 8. **Schierup MH, Hein J.** 2000. Consequences of recombination on
120 traditional phylogenetic analysis. *Genetics* **156**:879–891.
- 121 9. **Jukes TH, Cantor CR.** 1969. Evolution of protein molecules, p. 21–132. *In*
122 *In Mammalian protein metabolism*, Vol. III (1969), pp. 21-132.
- 123 10. **Guindon S, Gascuel O.** 2003. A simple, fast, and accurate algorithm to
124 estimate large phylogenies by maximum likelihood. *Syst. Biol.* **52**:696–704.
- 125 11. **Drummond AJ, Rambaut A.** 2007. BEAST: Bayesian evolutionary
126 analysis by sampling trees. *BMC Evol. Biol.* **7**:214.

- 127 12. **Drummond AJ, Suchard MA, Xie D, Rambaut A.** 2012. Bayesian
128 phylogenetics with BEAUti and the BEAST 1.7. *Mol. Biol. Evol.* **29**:1969–
129 1973.
- 130 13. **Paradis E, Claude J, Strimmer K.** 2004. APE: Analyses of phylogenetics
131 and evolution in R language. *Bioinformatics* **20**:289–290.
- 132 14. **Schliep KP.** 2011. phangorn: phylogenetic analysis in R. *Bioinformatics*
133 **27**:592–593.
- 134 15. **Robinson DF, Foulds LR.** 1981. Comparison of phylogenetic trees. *Math.*
135 *Biosci.*
- 136 16. **Pupko T, Pe'er I, Shamir R, Graur D.** 2000. A fast algorithm for joint
137 reconstruction of ancestral amino acid sequences. *Mol. Biol. Evol.* **17**:890–
138 896.
- 139