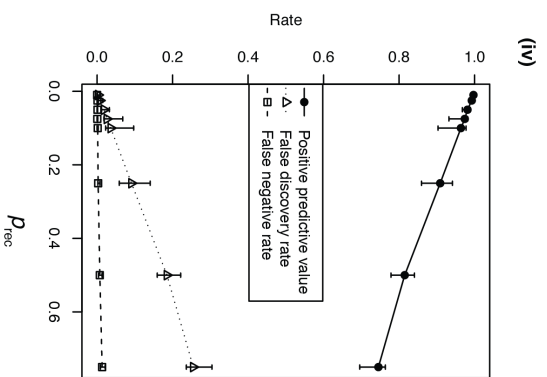**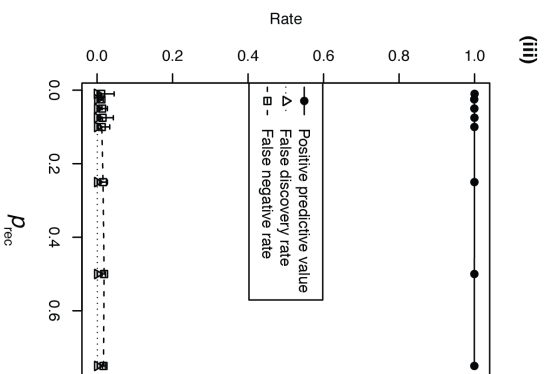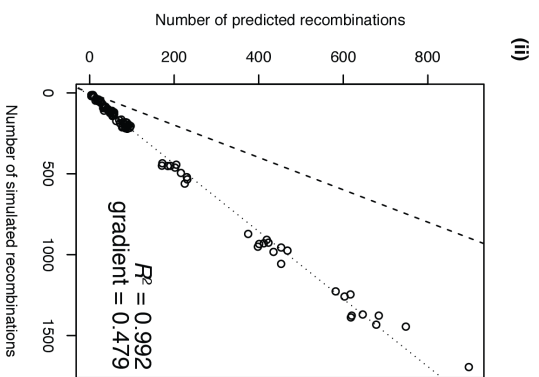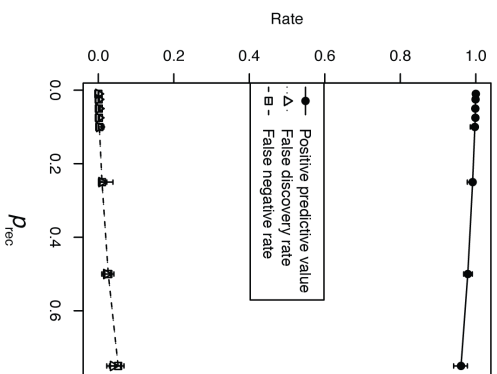Supplementary Materials - Rapid phylogenetic analysis of large samples of recombinant bacterial whole genome sequences using Gubbins**

**Figure S1**

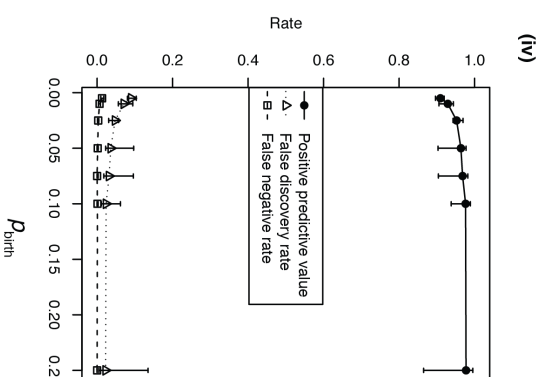**Figure S1:** Accuracy of Gubbins reconstructions from simulations in which recombinations may represent imports from divergent donors or exchanges between extant sequences within the simulation. Statistics representing the impact of varying (A) $p_{\text{rec}}$ and (B) $p_{\text{birth}}$ are shown as in Figure 1.

Figure S2

**Figure S2:** Accuracy of Gubbins reconstructions from simulations using recombination donors closely related to the recipient. Statistics representing the impact of v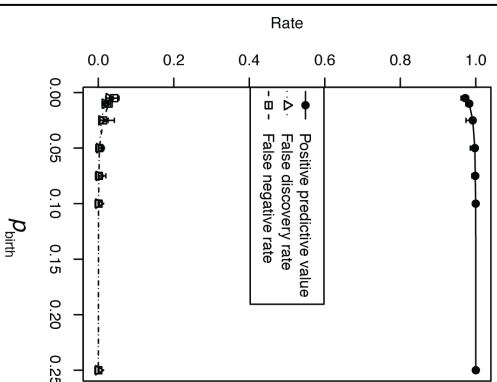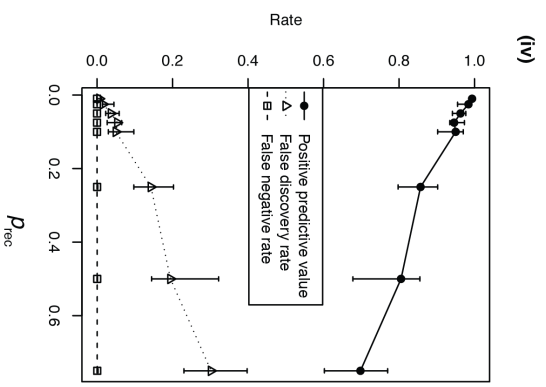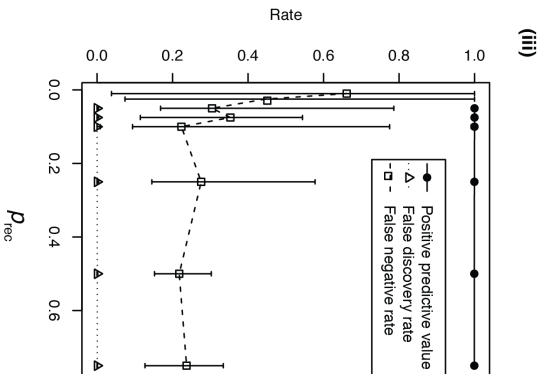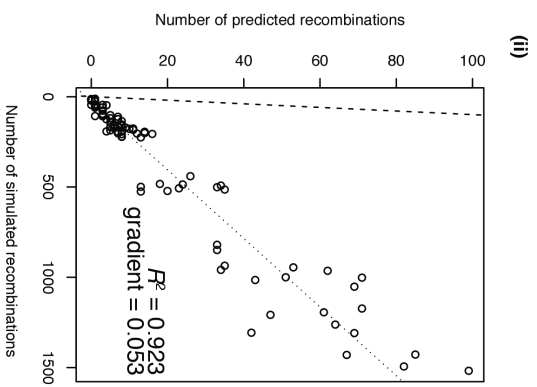arying (A) $p_{rec}$ and (B) $p_{birth}$ are shown as in Figure 1. At lower values of $p_{rec}$, and higher values of $p_{birth}$, some simulations did not result in any recombinations being predicted by Gubbins, and therefore the PPV and FDR for assigning true positive base substitutions to recombinations could not be evaluated; hence these points are omitted from panels (A)(iii) and (B)(iii).

**Figure S3**

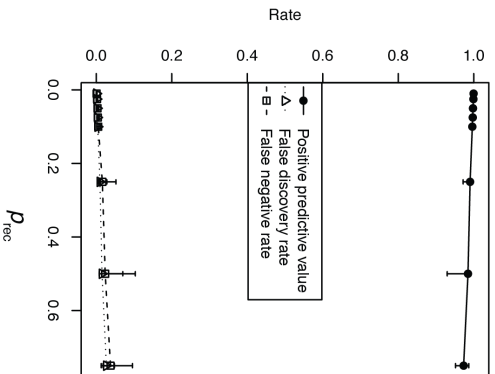**Figure S3:** Accuracy of phylogenetic reconstructions from Gubbins analyses of simulations in which recombinations may represent imports from divergent donors or exchanges between extant taxa within the simulation, displayed as in Figure 2.

**Figure S4**

**Figure S4:** Evidence for convergence of ClonalFrame in the analysis of *S. pneumoniae* PMEN1. The values of parameters estimated by the analysis are plotted over the duration of Monte Carlo Markov chains as blue lines, with the end of the burn-in iterations indicated by the vertical red dashed lines. The distribution of values is displayed 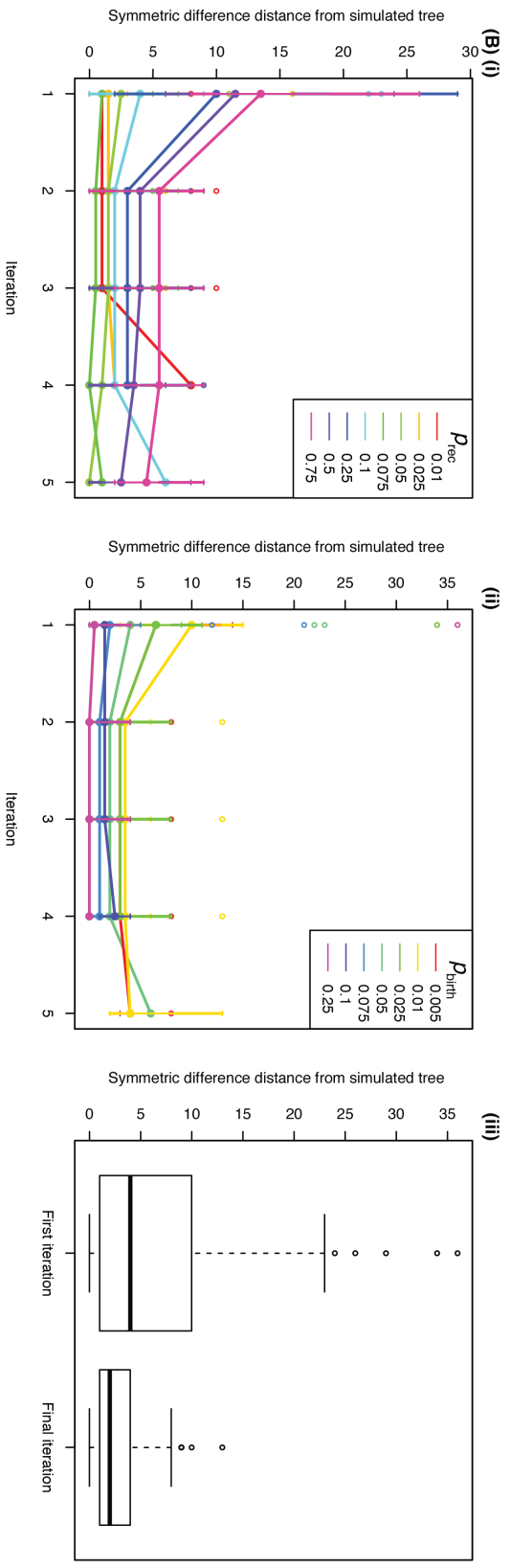by the histograms to the right of each plot. There is evidence of bistability in the estimates of some parameters (theta, R, TMRCA and total branch length); these all relate to the rate of mutation and recombination relative to branch length, and would not be expected to strongly impact on the comparison with the Gubbins analysis, which does not take into account the absolute length of the tree branches.

**Figure S5**

**Figure S5:** Predicted recombination size distributions from Gubbins and ClonalFrame. These histograms represent the lengths of the recombination events displayed in Figure 4 and Figure S6: (A) analysis performed using Gubbins on (i) *S. pneumoniae* PMEN1, (ii) *H. pylori* or (iii) *S. aureus* ST239; and (B) analysis performed using ClonalFrame on (i) *S. pneumoniae* PMEN1, (ii) *H. pylori* or (iii) *S. aureus* ST239. In each case, the red line represents the best-fitting exponential distribution.

(A)

0.5 Mb

Gambia94

26695

India7

Shi470

Sat464

F32

F30

F57

10000 substitutions

(B)

0.5 Mb

F30

F57

F32

Shi470

Sat464

26695

India7

Gambia94

0.17

**Figure S6:** Comparison of ClonalFrame and Gubbins analyses of *H. pylori* sequences. Eight complete publically available *H. pylori* genomes were analysed with both (A) Gubbins and (B) ClonalFrame. The results are displayed as described in Figure 5.

# Figure S7

**Figure S7:** Evidence for convergence of ClonalFrame in the analysis of *H. pylori*, as displayed in Figure S4. As node ages and theta were not estimated in this analysis, constraints that appeared necessary for the algorithm to converge in this case, total branch length, TMRCA and theta are invariant along the chain.

**Figure S8**

**(A)**



36 substitutions

S38
S7
DEN907
S2
S78
S26
TW20
S71
S93
S102
S40
S85
S87
S130

**(B)**

**(C)**

20 kb

φSa3(TW20)

Prophage from S102

**Figure S8:** Differences in prophage complement between *S. aureus* ST239 isolates. (A) The phylogeny produced by Gubbins, as displayed in Figure 5 (B)(ii). (B) Two different prophage sequences were concatenated together: that of φSa3 from TW20, and a similar prophage at the same insertion site from a *de novo* assembly of isolate S102. (C) A heatmap for each Illumina-sequenced isolate in the phylogeny represents the level of sequence read mapping across the two prophage sequences. Blue represents an absence of read mapping, red indicates high mapping, scaled to a maximum of 100-fold coverage. Read mapping was achieved using BWA as described in (25), except that only identical reads were mapped to the sequences in this instance to allow the similar viruses to be distinguished. Data were not plotted for isolate TW20, which contains φSa3, as this was the complete reference sequence against which all the other datasets were mapped.
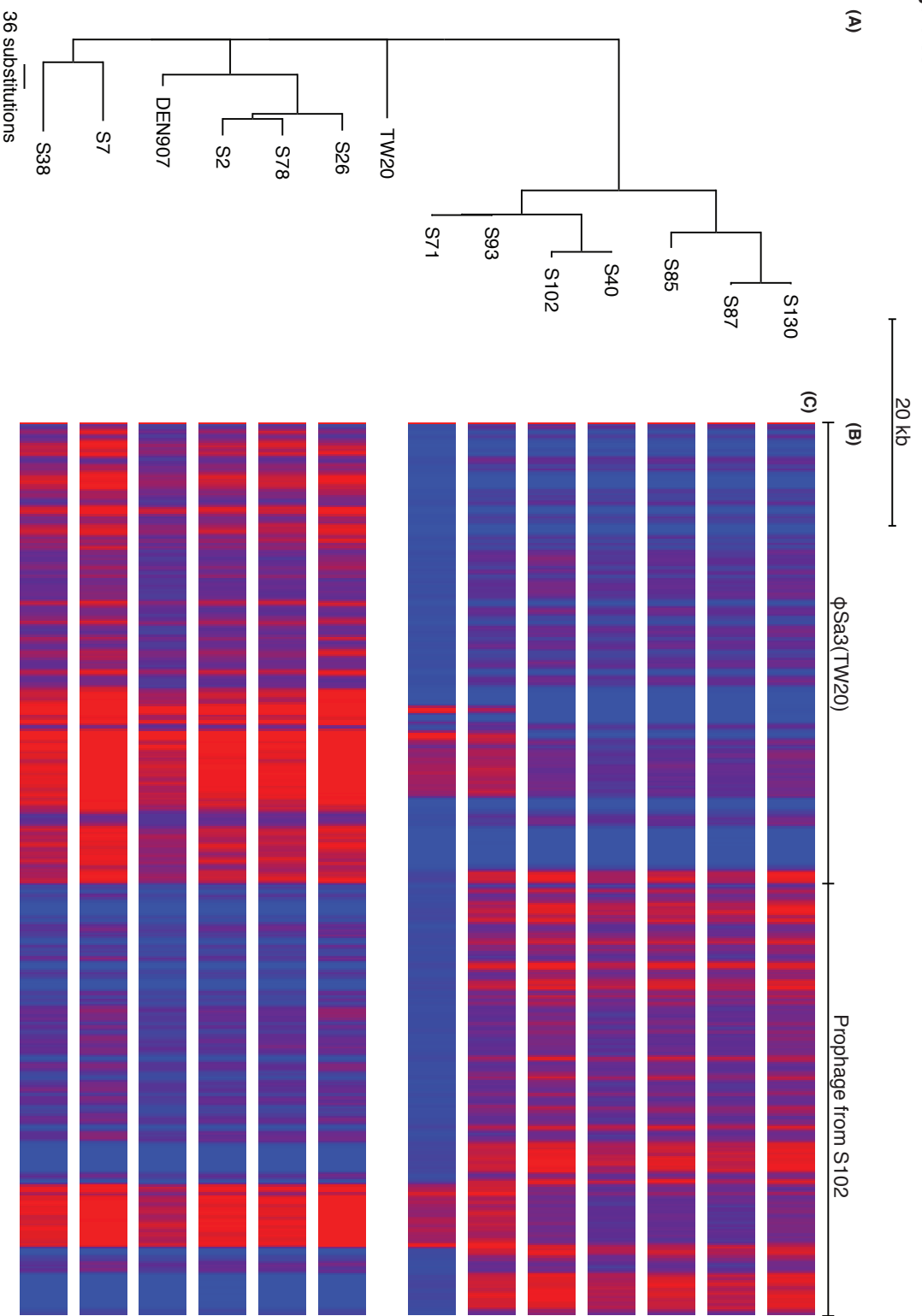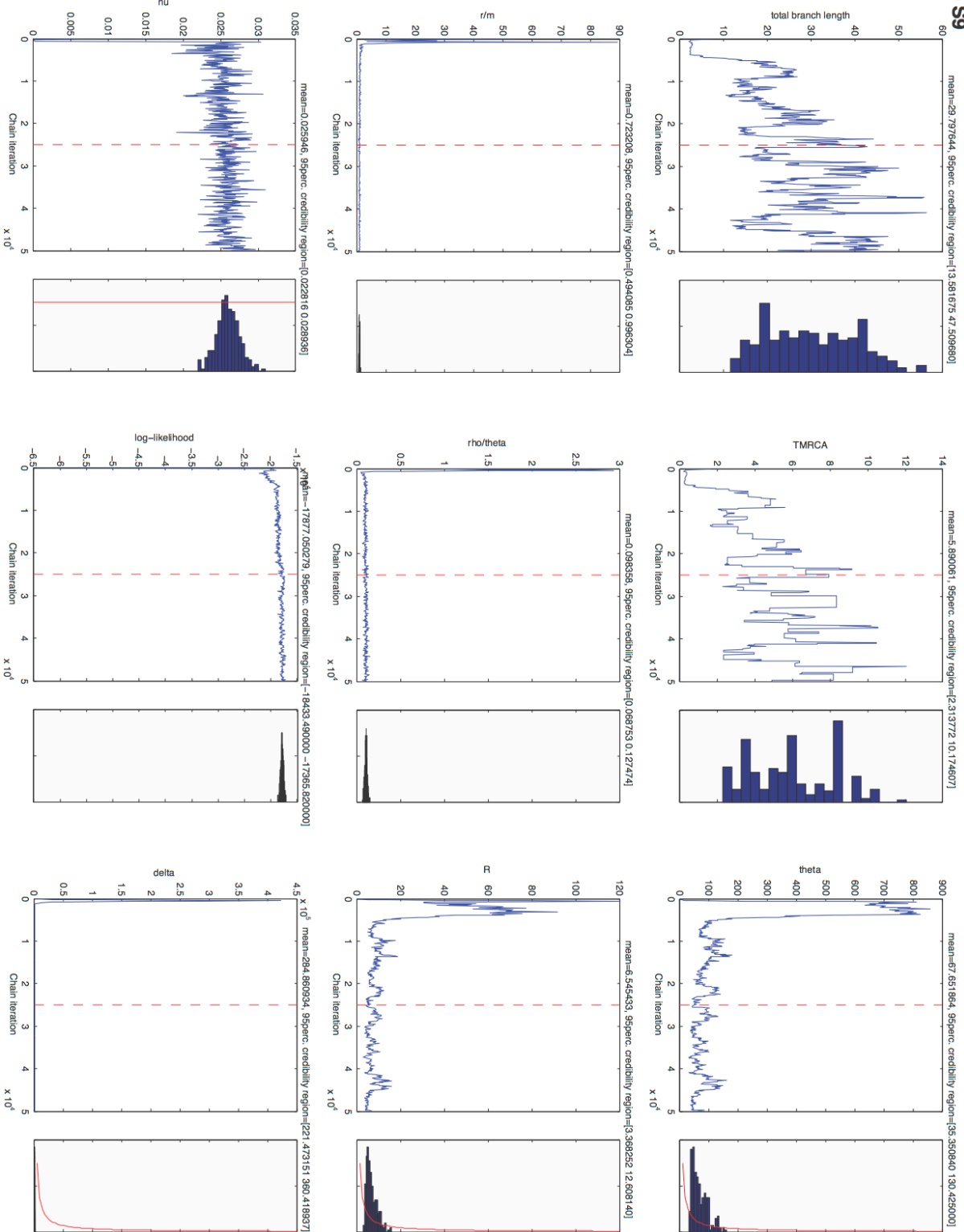
Figure S9

**Figure S9:** Evidence for convergence of ClonalFrame in the analysis of *S. aureus* ST239, as displayed in Figure S4.

**Table S1:** Details of genomes used as sequence donors in simulations analysed in Figure 1.

| Genome | Accession code |
|---|---|
| *Streptococcus pneumoniae* 670-6B | CP002176 |
| *Streptococcus pneumoniae* 70585 | CP000918 |
| *Streptococcus pneumoniae* AP200 | CP002121 |
| *Streptococcus pneumoniae* CGSP14 | CP001033 |
| *Streptococcus pneumoniae* D39 | CP000410 |
| *Streptococcus pneumoniae* G54 | CP001015 |
| *Streptococcus pneumoniae* gamPNI0373 | CP001845 |
| *Streptococcus pneumoniae* Hungary19A-6 | CP000936 |
| *Streptococcus pneumoniae* INV104 | FQ312030 |
| *Streptococcus pneumoniae* INV200 | FQ312029 |
| *Streptococcus pneumoniae* JJA | CP000919 |
| *Streptococcus pneumoniae* OXC141 | FQ312027 |
| *Streptococcus pneumoniae* P1031 | CP000920 |
| *Streptococcus pneumoniae* R6 | AE007317 |
| *Streptococcus pneumoniae* ST556 | CP003357 |
| *Streptococcus pneumoniae* Taiwan19F-14 | CP000921 |
| *Streptococcus pneumoniae* TCH8431/19A | CP001993 |
| *Streptococcus pneumoniae* TIGR4 | AE005672 |