

Supplementary Material

Genetic competence drives genome diversity in *Bacillus subtilis*

Patrícia H. Brito^{1,2*}, Bastien Chevreux³, Cláudia R. Serra⁴, Ghislain Schyns³, Adriano O. Henriques⁴,
José B. Pereira-Leal^{1,5*}

¹Instituto Gulbenkian de Ciência, Oeiras, Portugal; ²Nova Medical School, Faculdade de Ciências Médicas, Universidade Nova de Lisboa, Portugal; R&1D; ³DSM Nutritional Products, Ltd., Kaiseraugst, Switzerland; ⁴Instituto de Tecnologia Química e Biológica, Oeiras, Portugal; ⁵Ophiomics - Precision Medicine, Lisboa, Portugal.

* Corresponding authors

Patrícia H. Brito: pbrito@igc.gulbenkian.pt

José B. Pereira-Leal: jleal@igc.gulbenkian.pt

Instituto Gulbenkian de Ciência, Oeiras, Rua da Quinta Grande 6, 2780-156 Oeiras, Portugal; Phone: +351214407900; Fax: +351214407970

Supplementary methods

24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59

Orthology mapping

Gene orthology mapping for phylogenetic analysis was carried out using all-against-all blastp and three different sequence clustering algorithms, bidirectional best hit analysis (BDBH) (Contreras-Moreira & Vinuesa 2013), clusters of orthologous groups (COG) (Contreras-Moreira & Vinuesa 2013; Kristensen et al. 2010), and Markov cluster algorithm (OMCL) (Li et al. 2003) imposing a minimum pairwise sequence alignment coverage of 80% (SAC) and 50% of sequence identity (SI). From each of those analyses we retained all clusters characterized by the absence of inparalogs (sequences with best hits in its own genome), and the ubiquitous presence in all strains (core genes). The intersection of the clusters obtained from the three clustering algorithms was our final list of orthologous genes to be used in the phylogenetic analysis.

Orthology mapping to characterize the pangenome was performed with OMCL and following Tettelin *et al* (Tettelin et al. 2005) we imposed 50% of SAC and 80% of SI, we did not exclude clusters with inparalogs, and included all clusters independent of size. As in (Lefébure & Stanhope 2007) we only considered genes as taxon specific (clusters of size 1) if sequences were larger than 50 amino acids and had no blast hit with any other protein (E-value 1.0E-05). This was done to minimize the inclusion of truncated genes found at the end of contigs of open genomes. Still, because it is possible that assembly errors would truncate genes at the end of contigs precluding their annotation, we also performed an explicit test for the existence of non-annotated genes at both ends of the contigs of open genomes. This test was done as follows: for each genome, we extracted all nucleotide sequences hanging from the beginning of the contig to the beginning of the first annotated gene as well as the nucleotide sequences hanging from the end of the last annotated gene to the end of the contig. By definition, truncated genes should sit on those nucleotide positions. We BLASTN all those sequences against a database built with one representative sequence from each OMCL cluster (10137 in total). We considered a significant hit any BLASTN hit larger than 150 nucleotides, with more than 80% sequence identity in BLAST query/subject pairs, and 50% coverage in BLASTN pairwise alignments (alignment length/subject sequence length). This analysis lead to the identification of 510 false negatives that were corrected the final matrix (Supplementary dataset 2). Orthology mapping was done with get_homologues package (Contreras-Moreira & Vinuesa 2013).

Phylogenetic analysis of BmrA protein

For this analysis we extracted all 103632_bmrA homologs mapped in the KO (KEGG Orthology) database (K18104) and added homologous sequences from *B. subtilis* genomes not available in the

60 KEGG database totaling 498 sequences. We aligned this dataset with Mafft (G-ins-I, and Blosum 62)
61 and run maximum likelihood and bootstrap analyses as before. From the results we identified the most
62 inclusive well-supported clade that included *B. subtilis* sequences and its closest relatives. To this
63 subset of sequences we added representatives of the remaining clades to build a simplified dataset that
64 was realigned and reanalyzed with maximum likelihood and bootstrap to produce the final result
65 presented here. Within *B. subtilis* we identified three homologous copies. One from the softcore
66 genome (*B. subtilis* str. 168; NCBI_ProteinID: NP_391362.1; Genome location: NC_000964:
67 3577745..3579514) that has been transmitted vertically within *B. subtilis* group. The other two copies
68 belong to the shell (*B. subtilis* str. XF-1; NCBI_ProteinID: AGE62337; Genome location:
69 NC_020832:531386..533119) and the cloud (*B. subtilis* str. BSn5; NCBI_ProteinID: ADV95129;
70 Genome location: NC_014976:2435115..2436872) genome and were recently acquired by *B. subtilis*
71 strains from unrelated organisms. Since initial acquisition these later two genes have been transferred
72 across strains sampled among Plant niches. These three gene copies are on average 58% divergent
73 (between-group mean distance measured as the proportion of amino acid differences) hence too
74 divergent to be detected as orthologous copies by the algorithm used in this study.

75
76
77

78 **Supplementary results**

79

80 **Shared genes and core genome extrapolation (fig. 2d)**

81 The average values of the shared genes were extrapolated by fitting the exponential decay function:

$$82 \quad F_c(n) = 2119 \exp\left[\frac{-n}{18.45}\right] + 1659$$

83

$$84 \quad [1]$$

85 where n is the number of distinct genomes, and the size of the core genome as $n \rightarrow \infty$ is 1659.

86 Estimations were performed with get_homologues package (Contreras-Moreira & Vinuesa 2013)
87 following Tettelin *et al* (Tettelin et al. 2005).

88

89 **Strain-specific genes and pangenome extrapolation (fig. 2e)**

90 The average number of specific genes per strain closely follows an exponential decay function with
91 function:

$$92 \quad P(n) = 4194 + 57.0(n - 1) + 219 \exp\left[\frac{-2}{7.80}\right] \cdot \left[\frac{1 - \exp\left(\frac{-(n-1)}{7.80}\right)}{1 - \exp\left(\frac{-1}{7.80}\right)} \right]$$

93

$$[2]$$

94 where n is the number of distinct genomes, and 57.0 is the number of new genes asymptotically
95 predicted for each new sequenced genome. For $n=42$ genomes the extrapolated pangenome size is
96 $P(n) \approx 7932.9$. Estimations were performed with get_homologues package (Contreras-Moreira &
97 Vinuesa 2013) following Tettelin *et al* (Tettelin et al. 2005).

98

99 **Shared and strain-specific genes for core and pangenome extrapolation using only *B. subtilis*** 100 **wild strains (32 genomes) (fig. S3)**

101

102 Fitted functions are,

$$103 \quad F_c(n) = 2097 \exp\left[\frac{-n}{12.77}\right] + 1803$$

104

$$105 \quad [3]$$

106 where n is the number of distinct genomes, and the size of the core genome as $n \rightarrow \infty$ is 1803(SE
107 225.20), and

$$108 \quad P(n) = 4169 + 73.0(n - 1) + 276 \exp\left[\frac{-2}{6.35}\right] \cdot \left[\frac{1 - \exp\left(\frac{-(n-1)}{6.35}\right)}{1 - \exp\left(\frac{-1}{6.35}\right)} \right]$$

109

$$[4]$$

110 where n is the number of distinct genomes, and 73.0 is the number of new genes asymptotically
111 predicted for each new sequenced genome. For $n=32$ genomes the extrapolated pangenome size is
112 $P(n) \approx 7806.5$ (SE 202.77). Estimations were performed with get_homologues package (Contreras-
113 Moreira & Vinuesa 2013) following Tettelin *et al* (Tettelin et al. 2005).

114

115

116

117 **Supplementary references**

118

119 Contreras-Moreira B, Vinuesa P. 2013. GET_HOMOLOGUES, a versatile software package for
120 scalable and robust microbial pangenome analysis. *Appl Environ Microbiol.* 79:7696–7701. doi:
121 10.1128/AEM.02411-13.

122 Kristensen DM et al. 2010. A low-polynomial algorithm for assembling clusters of orthologous groups
123 from intergenomic symmetric best matches. *Bioinformatics.* 26:1481–1487. doi:
124 10.1093/bioinformatics/btq229.

125 Krzywinski M et al. 2009. Circos: An information aesthetic for comparative genomics. *Genome Res.*
126 19:1639–1645. doi: 10.1101/gr.092759.109.

127 Lefébure T, Stanhope MJ. 2007. Evolution of the core and pan-genome of *Streptococcus*: positive
128 selection, recombination, and genome composition. *Genome Biol.* 8:R71. doi: 10.1186/gb-2007-8-5-

129 r71.

130 Li L, Stoeckert CJ, Roos DS. 2003. OrthoMCL: identification of ortholog groups for eukaryotic
131 genomes. *Genome Res.* 13:2178–2189. doi: 10.1101/gr.1224503.

132 Tettelin H et al. 2005. Genome analysis of multiple pathogenic isolates of *Streptococcus agalactiae*:
133 Implications for the microbial ‘pan-genome’. *P Natl Acad Sci Usa.* 102:13950–13955. doi:
134 10.2307/3376809?ref=search-gateway:46f22eed2f66acb6f02d383ffa46dc1e.

135

Supplementary figure legends

136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171

Fig. S1. Phylogenetic analysis of *B. subtilis* group with distant outgroups. Maximum-likelihood (LG+I+G) and bootstrap analyses were carried out on a concatenated dataset of 398 core proteins (104056 AA). This analysis was done to confirm the taxonomic identification of genomes available at the NCBI database and to confirm the outgroup selection. Five genomes (in red) were excluded from further analyses as they are more closely related to other species than to *B. subtilis*. Numbers on the branches are bootstrap support. Due to the small branch lengths at the tips of the tree we only show bootstrap support for the most basal branches. The scale bar indicates expected number of amino acid substitutions per site.

Fig. S2. Phylogenetic analysis of *B. subtilis* group and its closest outgroups. Maximum-likelihood (GTR+I+G+X) and bootstrap analyses were carried out on a concatenated dataset of 685 core genes (520227 nt). Numbers on the branches are bootstrap support. Due to the small branch lengths at the tips of the tree we only show bootstrap support for the most basal branches. The scale bar indicates expected number of nucleotide substitutions per site. Arrow indicates the branch that separates *B. subtilis* from other species.

Fig. S3. Core and pangenome in wild strains. Quantification of the core (a) and the pangenome (b) using only wild strains of *B. subtilis*. The number of shared (a) and novel (b) genes estimated after 100 random samples of the 32 genomes is plotted against the number of n strains sequentially added. Line in red is the fitted curve following exponential functions as in Tettelin (Tettelin et al. 2005).

Fig. S4. Distribution of intact prophages in *B. subtilis* shows the existence of three high frequency elements in the species pangenome. PBSX is ubiquitous among genomes in our dataset and SP β and *skin* elements are largely restricted to *B. s. subtilis*. Detection of intact prophages was performed using PHAST and identification resulted from BLAST searches of prophage sequences against NCBI nr/nt databases. Presence and absence (black/white squares) is plotted on the core genome tree of figure 1. Numbers on the terminal branches indicate total numbers of intact prophages detected in each genome. The tree is not drawn to scale. Colors in the taxa names reflect niche at the site of sampling following fig. 1.

Fig. S5. Variable pangenome is distributed throughout the chromosome. Genome atlas of *B. subtilis* strains with closed genomes. The circles from the center represent the location of the following elements: coding sequences (dark grey), tRNA (red) and prophages (intact: dark green;

172 questionable: green; incomplete: light green), core genome (black), softcore genome (yellow), shell
173 genome (orange), cloud genome (red). The outer circle represents relative frequency of each gene
174 cluster in the dataset. Images were created with Circos v. 0.67 (Krzywinski et al. 2009). The radial
175 histograms at the center of the genome atlas are rose diagrams showing the dispersion of cloud genes
176 around the genome after transforming each gene middle point position into radians. The genome was
177 divided into 100 bins and the radii of the sectors are equal to the square root of the relative frequencies
178 of the cloud genome making the area of each sector proportional to its frequency. The rose diagram
179 was scaled to the inner circle of the genome atlas.

180

181 **Fig. S6. Lateral gene acquisition of paralogue copies of a *B. subtilis* softcore gene from unrelated**
182 **organisms and posterior lateral transmission within the species.** Phylogenetic history of the *bmrA*
183 gene. This gene codes for a multidrug resistance ABC transporter ATP-binding protein, and its
184 evolutionary history is characterized by frequent lateral transference events between gram-positive
185 bacteria. In *B. subtilis* we identify three homologous copies, one from the softcore genome with
186 vertical transmission, and two xenologue copies from the shell and the cloud genome that are recent
187 LGT events from unrelated organisms. Since initial acquisition these later two genes have been
188 transferred across strains sampled among Plant niches. The phylogenetic tree was computed with
189 maximum-likelihood (LG+I+G), and numbers on the branches are bootstrap support. Analyses were
190 carried out in RaxML. The scale bar indicates expected number of nucleotide substitutions per site.

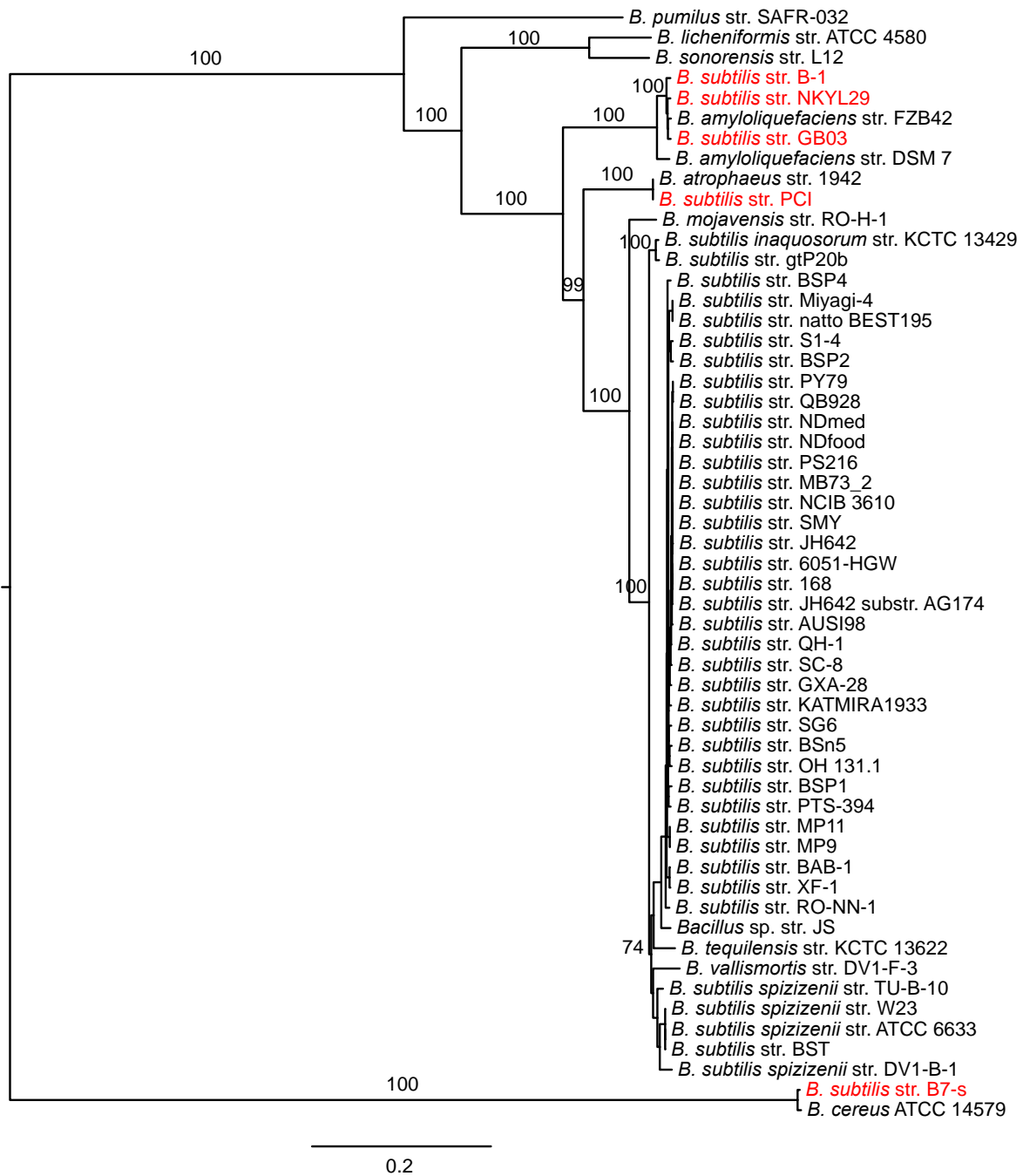
191

192

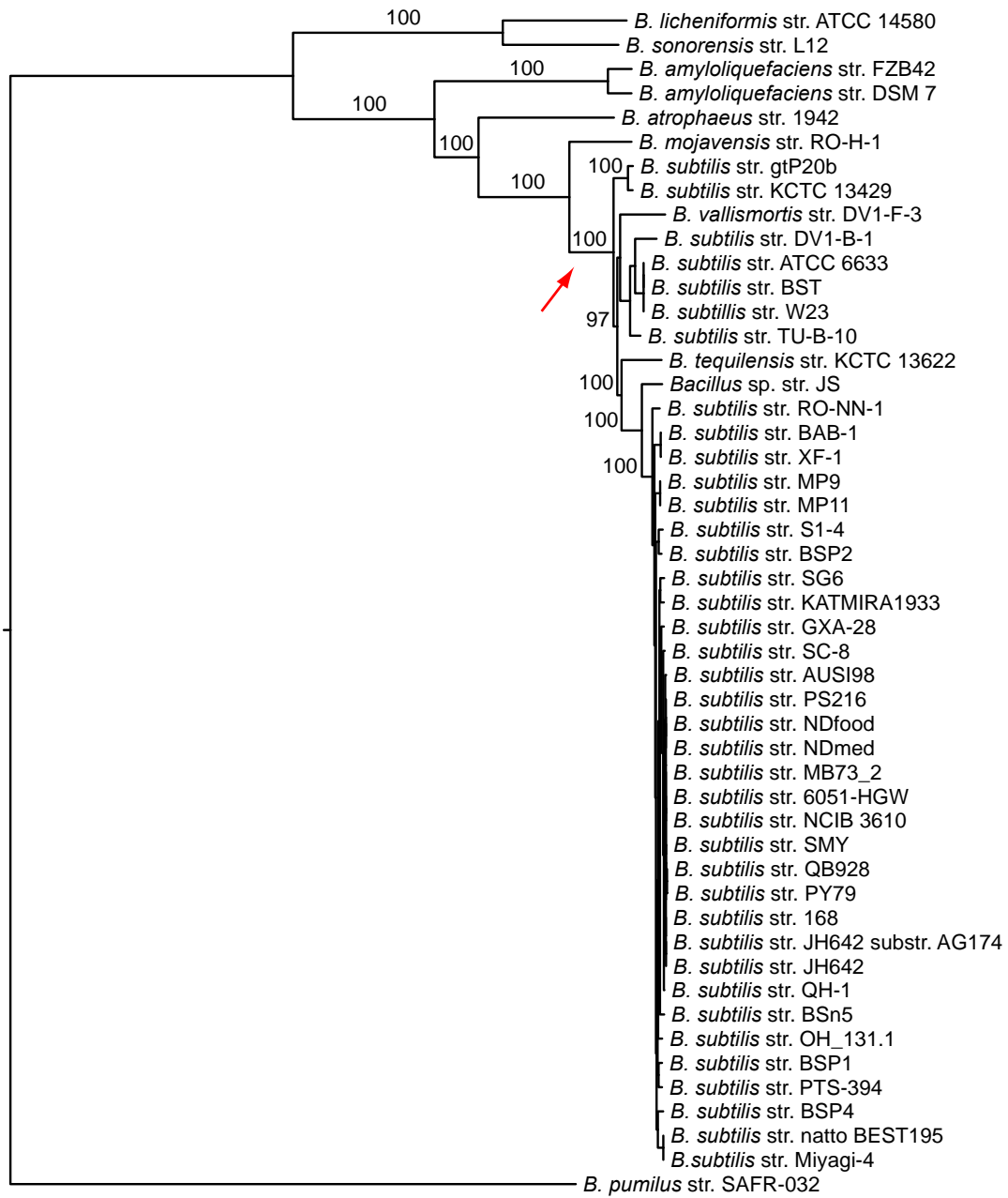
193
194
195

Supplementary figures

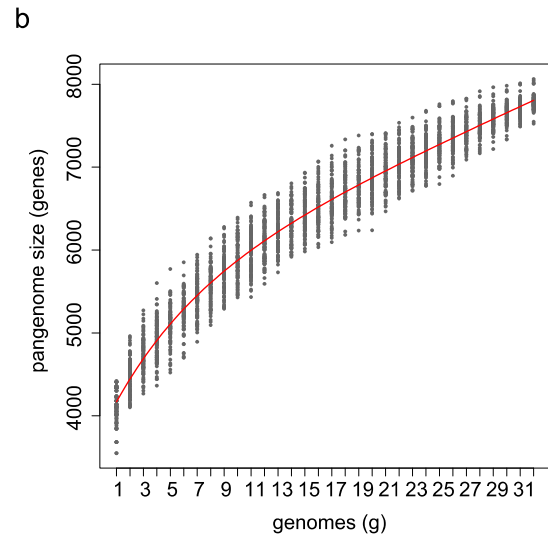
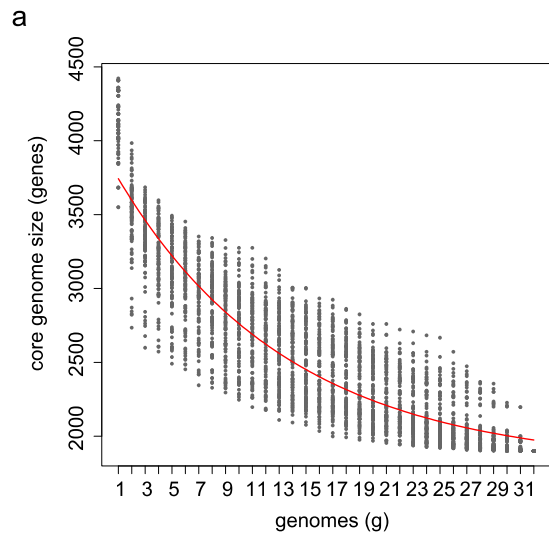
196 Fig. S1.



197
198 Fig. S2.

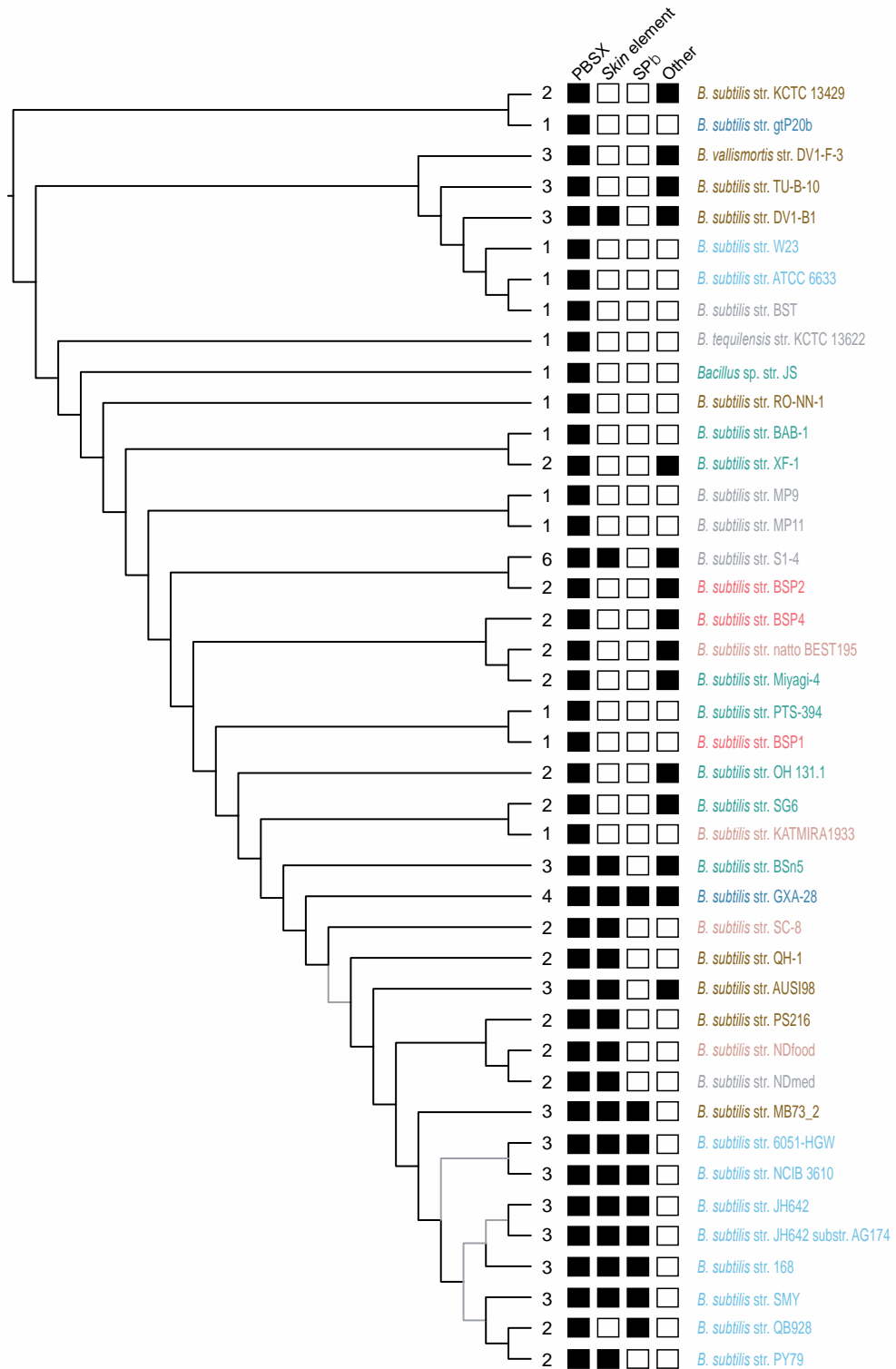


0.2



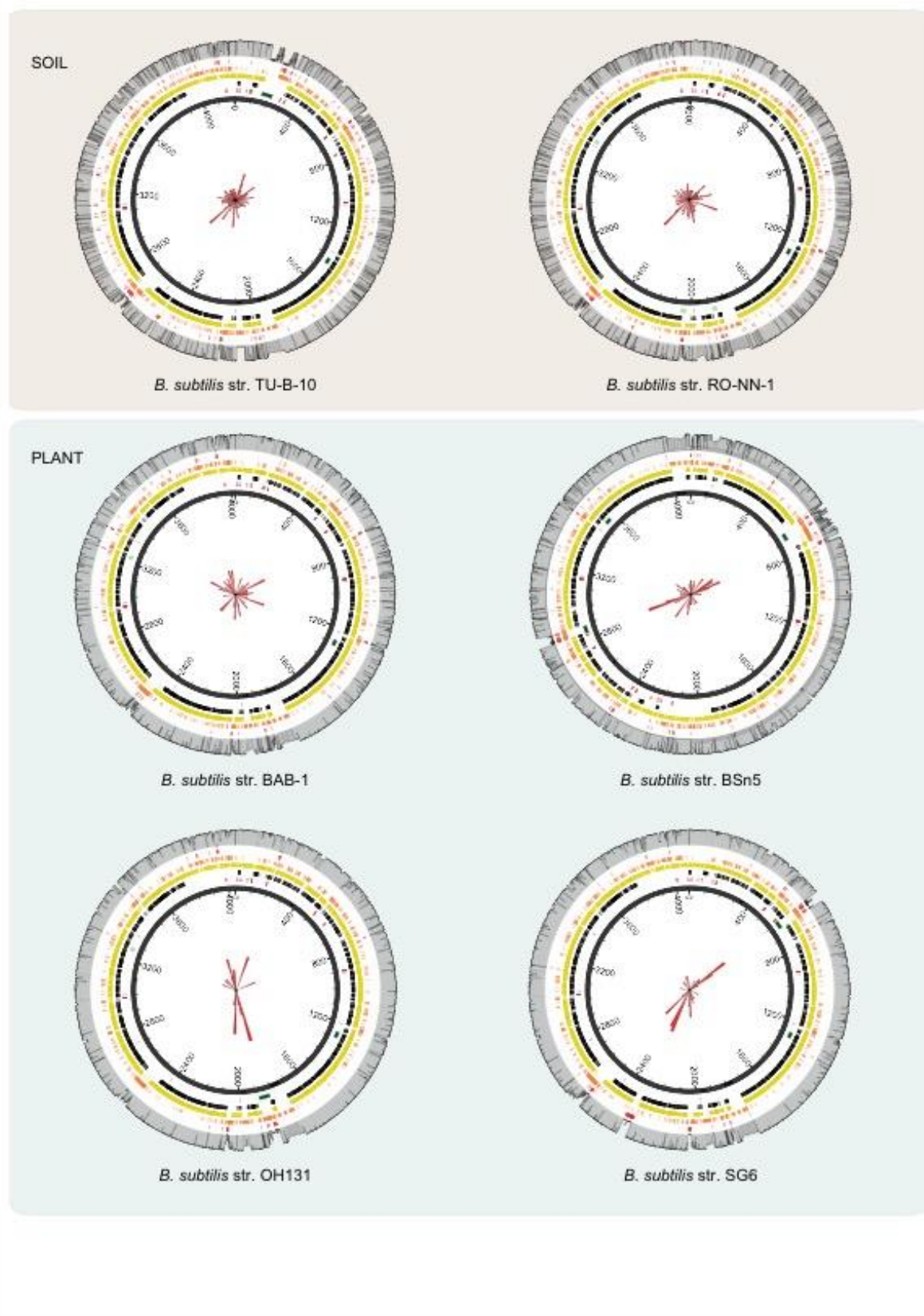
200
201
202

Fig. S3.



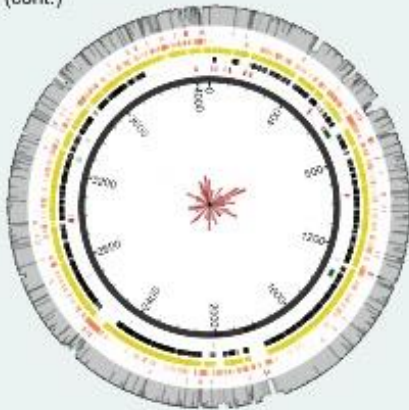
203
204
205
206

Fig. S4



207
208 **Fig. S5**
209

PLANT (cont.)

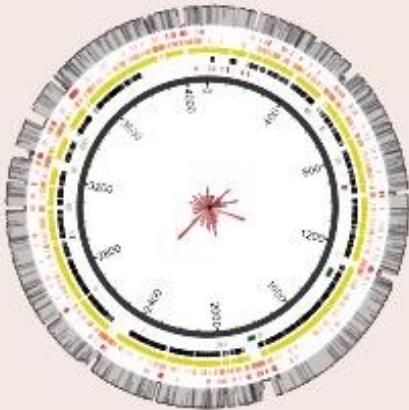


B. subtilis str. XF-1



Bacillus sp str. JS

FOOD



B. subtilis str. BEST195

LABORATORY



B. subtilis str. W23



B. subtilis str. 168



B. subtilis str. 6051-HGW

210
211
212

Fig. S5 (cont.)

LABORATORY (cont.)



B. subtilis str. JH642



B. subtilis str. JH642-AG174



B. subtilis str. NCIB3610



B. subtilis str. PY79



B. subtilis str. QB928



B. subtilis str. SMY

213
214
215
216

Fig. S5 (cont.)



217 Fig. S6
 218
 219