

1 **GET\_PHYLOMARKERS, a software package to select optimal orthologous**  
2 **clusters for phylogenomics and inferring pan-genome phylogenies, used for a**  
3 **critical geno-taxonomic revision of the genus *Stenotrophomonas***  
4

5  
6 Pablo Vinuesa<sup>1\*</sup>, Luz Edith Ochoa-Sánchez<sup>1</sup> and Bruno Contreras-Moreira<sup>2,3</sup>.  
7

8 Pablo Vinuesa's ORCID: <http://orcid.org/0000-0001-6119-2956>  
9 Bruno Contrera-Moreira's ORCID: <https://orcid.org/0000-0002-5462-907X>  
10

11 <sup>1</sup>Centro de Ciencias Genómicas, Universidad Nacional Autónoma de México, Cuernavaca, Morelos,  
12 Mexico.  
13

14 <sup>2</sup>Estación Experimental de Aula Dei – Consejo Superior de Investigaciones Científicas, Zaragoza,  
15 Spain.

16 <sup>3</sup>Fundación ARAID, Zaragoza, Spain.

17 \*Correspondence: Pablo Vinuesa, [vinuesa@ccg.unam.mx](mailto:vinuesa@ccg.unam.mx)  
18  
19

20 Running title: GET\_PHYLOMARKERS: open-source software for phylogenomics and genome  
21 taxonomy  
22  
23

24 Number of figures: 7

25 Number of tables: 3

26 Word count of main text: 10,232  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44

45  
46  
47

48 **Abstract. (333 words)**

49 The massive accumulation of genome-sequences in public databases promoted the proliferation of  
50 genome-level phylogenetic analyses in many areas of biological research. However, due to diverse  
51 evolutionary and genetic processes, many loci have undesirable properties for phylogenetic  
52 reconstruction. These, if undetected, can result in erroneous or biased estimates, particularly when  
53 estimating species trees from concatenated datasets. To deal with these problems, we developed  
54 GET\_PHYLOMARKERS, a pipeline designed to identify high-quality markers to estimate robust  
55 genome phylogenies from the orthologous clusters, or the pan-genome matrix (PGM), computed by  
56 GET\_HOMOLOGUES. In the first context, a set of sequential filters are applied to exclude  
57 recombinant alignments and those producing anomalous or poorly resolved trees. Multiple sequence  
58 alignments and maximum likelihood (ML) phylogenies are computed in parallel on multi-core  
59 computers. A ML species tree is estimated from the concatenated set of top-ranking alignments at the  
60 DNA or protein levels, using either FastTree or IQ-TREE (IQT). The latter is used by default due to its  
61 superior performance revealed in an extensive benchmark analysis. In addition, parsimony and ML  
62 phylogenies can be estimated from the PGM.

63 We demonstrate the practical utility of the software by analyzing 170 *Stenotrophomonas* genome  
64 sequences available in RefSeq and 10 new complete genomes of environmental *S. maltophilia* complex  
65 (Smc) isolates reported herein. A combination of core-genome and PGM analyses was used to revise  
66 the molecular systematics of the genus. An unsupervised learning approach that uses a goodness of  
67 clustering statistic identified 20 groups within the Smc at a core-genome average nucleotide identity of  
68 95.9% that are perfectly consistent with strongly supported clades on the core- and pan-genome trees.  
69 In addition, we identified 14 misclassified RefSeq genome sequences, 12 of them labeled as *S.*  
70 *maltophilia*, demonstrating the broad utility of the software for phylogenomics and geno-taxonomic  
71 studies. The code, a detailed manual and tutorials are freely available for Linux/UNIX servers under the  
72 GNU GPLv3 license at [https://github.com/vinuesa/get\\_phylomarkers](https://github.com/vinuesa/get_phylomarkers). A docker image bundling  
73 GET\_PHYLOMARKERS with GET\_HOMOLOGUES is available at  
74 [https://hub.docker.com/r/csicunam/get\\_homologues/](https://hub.docker.com/r/csicunam/get_homologues/), which can be easily run on any platform.

75  
76

77 **Keywords.** Phylogenetics, genome-phylogeny, maximum-likelihood, species-tree, species delimitation,  
78 *Stenotrophomonas maltophilia* complex, Mexico.

## 79 INTRODUCTION

80

81 Accurate phylogenies represent key models of descent in modern biological research. They are applied  
82 to the study of a broad spectrum of evolutionary topics, ranging from the analysis of populations up to  
83 the ecology of communities (Dornburg et al., 2017). The way microbiologists describe and delimit  
84 species is undergoing a major revision in the light of genomics (Rosselló-Móra and Amann, 2015;  
85 Vandamme and Peeters, 2014), as reflected in the emerging field of microbial genomic taxonomy  
86 (Konstantinidis and Tiedje, 2007; Thompson et al., 2009, 2013). Current geno-taxonomic practice is  
87 largely based on the estimation of (core-)genome phylogenies (Ciccarelli et al., 2006; Daubin et al.,  
88 2002; Lerat et al., 2003; Tettelin et al., 2005; Wu and Eisen, 2008) and the computation of diverse  
89 overall genome relatedness indices (OGRIs) (Chun and Rainey, 2014), such as the popular genomic  
90 average nucleotide identity (gANI) values (Goris et al., 2007; Konstantinidis and Tiedje, 2005; Richter  
91 and Rossello-Mora, 2009). These indices are rapidly and effectively replacing the traditional DNA-  
92 DNA hybridization values used for species delimitation in the pre-genomic era (Stackebrandt et al.,  
93 2002; Stackebrandt and Goebel, 1994; Vandamme et al., 1996).

94

95 The ever-increasing volume of genome sequences accumulating in public sequence repositories  
96 provides a huge volume of data for phylogenetic analysis. This significantly improves our capacity to  
97 understand the evolution of species and any associated traits (Dornburg et al., 2017). However, due to  
98 diverse evolutionary forces and processes, many loci in genomes have undesirable properties for  
99 phylogenetic reconstruction. If undetected, these can lead to erroneous or biased estimates (Parks et al.,  
100 2018; Shen et al., 2017), although, ironically, with strong branch support (Kumar et al., 2012). Their  
101 impact is particularly strong in concatenated datasets (Degnan and Rosenberg, 2009; Kubatko and  
102 Degnan, 2007), which are standard in microbial phylogenomics (Wu and Eisen, 2008). Hence, robust  
103 phylogenomic inference requires the selection of well-suited markers for the task (Vinuesa, 2010).

104

105 For this study we developed GET\_PHYLOMARKERS, an open-source and easy-to-use software  
106 package designed with the aim of inferring robust genome-level phylogenies and providing tools for  
107 microbial genome taxonomy. We describe the implementation details of the pipeline and how it  
108 integrates with GET\_HOMOLOGUES (Contreras-Moreira and Vinuesa, 2013; Vinuesa and Contreras-  
109 Moreira, 2015). The latter is a popular and versatile genome-analysis software package designed to  
110 identify robust clusters of homologous sequences. It has been widely used in microbial pan-genomics  
111 and comparative genomics (Lira et al., 2017; Nourdin-Galindo et al., 2017; Sandner-Miranda et al.,  
112 2018; Savory et al., 2017), including recent bacterial geno-taxonomic (Gauthier et al., 2017; Gomila et  
113 al., 2017), and plant pan-genomic studies (Contreras-Moreira et al., 2017; Gordon et al., 2017).  
114 Regularly updated auxiliary scripts bundled in the GET\_HOMOLOGUES package compute diverse  
115 OGRIs, at the protein, CDS and transcript levels, provide graphical and statistical tools for a range of  
116 pan-genome analyses, including inference of pan-genome phylogenies under the parsimony criterion.  
117 GET\_PHYLOMARKERS was designed to work both at the core-genome and pan-genome levels,  
118 using either the homologous gene clusters or the pan-genome matrix computed by  
119 GET\_HOMOLOGUES. In the first context, it identifies single-copy orthologous gene families with  
120 optimal attributes (listed further down) and concatenates them to estimate a genomic species tree. In the  
121 second scenario, it uses the pan-genome matrix (PGM) to estimate phylogenies under the maximum  
122 likelihood (ML) and parsimony optimality criteria. In addition, we implemented unsupervised learning  
123 methods that automatically identify species-like genome clusters based on the statistical analysis of the  
124 PGM and core-genome average nucleotide identity matrices (cgANIb).

125 To demonstrate these capabilities and benchmark performance, we applied the pipeline to critically  
126 evaluate the molecular systematics and taxonomy of the genus *Stenotrophomonas*. Species delimitation  
127 is problematic and far from resolved in this genus (Ochoa-Sánchez and Vinuesa, 2017), despite recent  
128 efforts using genomic approaches with a limited number of genome sequences (Lira et al., 2017; Patil  
129 et al., 2016; Yu et al., 2016).

130

131 The genus *Stenotrophomonas* (Gammaproteobacteria, *Xanthomonadales*, *Xanthomonadaceae*)  
132 (Palleroni, 2005; Palleroni and Bradbury, 1993) groups ubiquitous, aerobic, non-fermenting bacteria  
133 that thrive in diverse aquatic and edaphic habitats, including human-impacted ecosystems (Ryan et al.,  
134 2009). As of March 2018, 14 validly described species were listed in Jean Euzéby's list of prokaryotic  
135 names with standing in nomenclature (<http://www.bacterio.net/stenotrophomonas.html>). By far, its  
136 best-known species is *S. maltophilia*. It is considered a globally emerging, multidrug-resistant (MDR)  
137 and opportunistic pathogen (Brooke, 2012; Chang et al., 2015). *S. maltophilia*-like organisms display  
138 high genetic, ecological and phenotypic diversity (Valdezate et al., 2004; Vasileuskaya-Schulz et al.,  
139 2011), forming the so-called *S. maltophilia* complex (Smc) (Berg and Martinez, 2015; Svensson-  
140 Stadler et al., 2012). Heterogeneous resistance and virulence phenotypes have been reported for  
141 environmental isolates of diverse ecological origin classified as *S. maltophilia* (Adamek et al., 2011;  
142 Deredjian et al., 2016). We have recently shown that this phenotypic heterogeneity largely results from  
143 problems in species delimitations within the Smc (Ochoa-Sánchez and Vinuesa, 2017). We analyzed  
144 the genetic diversity in a collection of 108 *Stenotrophomonas* isolates recovered from several water  
145 bodies in Morelos, Central Mexico, based on sequence data generated for the 7 loci used in the  
146 Multilocus Sequence Typing (MLST) scheme available for *S. maltophilia* at <https://pubmlst.org>. We  
147 assembled a large set of reference sequences retrieved from the MLST database (Kaiser et al., 2009;  
148 Vasileuskaya-Schulz et al., 2011) and from selected genome sequences (Crossman et al., 2008;  
149 Davenport et al., 2014; Lira et al., 2012; Patil et al., 2016; Vinuesa and Ochoa-Sánchez, 2015),  
150 encompassing 11 out of the 12 validly described species at the time. State-of-the-art phylogenetic and  
151 population genetics methods, including the multispecies coalescent model coupled with Bayes factor  
152 analysis and Bayesian clustering of the multilocus genotypes consistently resolved five conservatively-  
153 defined genospecies within the Smc clade, which were named *S. maltophilia* and Smc1-Smc4. The  
154 approach also delimited Smc5 as a sister clade of *S. rhizophila*. Importantly, we showed that *i*) only  
155 members of the Smc clade that we designed as *S. maltophilia* were truly MDR and *ii*) that *S.*  
156 *maltophilia* was the only species that consistently expressed metallo-beta-lactamases (Ochoa-Sánchez  
157 and Vinuesa, 2017). Strains of the genospecies Smc1 and Smc2 were only recovered from the Mexican  
158 rivers and displayed significantly lower resistance levels than sympatric *S. maltophilia* isolates,  
159 revealing well-defined species-specific phenotypes.

160

161 Given this context, the present study was designed with two major goals. The first one was to develop  
162 GET\_PHYLOMARKERS, a pipeline for the automatic and robust estimation of genome phylogenies  
163 using state-of-the art methods. The emphasis of the pipeline is on selecting top-ranking markers for the  
164 task, based on the following quantitative/statistical criteria: *i*) they should not present signs of  
165 recombination, *ii*) the resulting gene trees should not be anomalous or deviating from the distribution of  
166 tree topologies and branch lengths expected under the multispecies coalescent model and *iii*) they  
167 should have a strong phylogenetic signal. The top-scoring markers are concatenated to estimate the  
168 species phylogeny under the maximum likelihood optimality criterion using either FastTree (Price et  
169 al., 2010) or IQ-TREE (Nguyen et al., 2015). The second aim was to apply GET\_PHYLOMARKERS  
170 to challenge and refine the species delimitations reported in our previous MLSA study (Ochoa-Sánchez  
171 and Vinuesa, 2017) using a genomic approach, focusing on resolving the geno-taxonomic structure of

172 the Smc and *S. maltophilia sensu lato* clades. For this purpose we sequenced five strains from the new  
173 genospecies Smc1 and Smc2 and analyzed them together with all reference genome sequences  
174 available for the genus *Stenotrophomonas* as of August 2017 using the methods implemented in  
175 GET\_PHYLOMARKERS. The results were used to critically revise the molecular systematics of the  
176 genus in light of genomics, identify misclassified genome sequences, suggest correct classifications for  
177 them and discover multiple novel genospecies within *S. maltophilia*.

178

179

## 180 MATERIALS AND METHODS

181

### 182 Genome sequencing, assembly and annotation

183 Ten *Stenotrophomonas* strains from our collection were selected (Table 1) for genome sequencing  
184 using a MiSeq instrument (2x300 bp) at the Genomics Core Sequencing Service provided by Arizona  
185 State University (DNASU). They were all isolated from rivers in the state of Morelos, Central Mexico,  
186 and classified as genospecies 1 (Smc1) or 2 (Smc2), as detailed in a previous publication (Ochoa-  
187 Sánchez and Vinuesa, 2017). Adaptors at the 5'-ends and low quality residues at the 3' ends of reads  
188 were trimmed-off using ngsShoRT v2.1 (Chen et al., 2014) and passed to Spades v3.10.1 (Bankevich et  
189 al., 2012) for assembly (with options --careful -k 33,55,77,99,127,151). The resulting assembly  
190 scaffolds were filtered to remove those with low coverage (< 7X) and short length (< 500 nt). All  
191 complete genome sequences available in RefSeq for *Stenotrophomonas* spp. were used as references  
192 for automated ordering of assembly scaffolds using MeDuSa v1.6 (Bosi et al., 2015). A final assembly  
193 polishing step was performed by remapping the quality-filtered sequence reads on the ordered scaffolds  
194 using BWA (Li and Durbin, 2009) and passing the resulting sorted binary alignments to SAMtools (Li  
195 et al., 2009) for indexing. The indexed alignments were used by Pilon 1.21 (Walker et al., 2014) for gap  
196 closure and filling, correction of indels and single nucleotide polymorphisms (SNPs), as previously  
197 described (Vinuesa and Ochoa-Sánchez, 2015). The polished assemblies were annotated with NCBI's  
198 Prokaryotic Genome Annotation Pipeline (PGAP v4.2) (Angiuoli et al., 2008). BioProject and  
199 BioSample accession numbers are provided in Table S1.

200

### 201 Reference genomes

202 On August 1<sup>st</sup>, 2017, a total of 169 annotated *Stenotrophomonas* genome sequences were available in  
203 RefSeq, 134 of which were labeled as *S. maltophilia*. The corresponding GenBank files were retrieved,  
204 as well as the corresponding table with assembly metadata. Seven complete *Xanthomonas* spp.  
205 genomes were also downloaded to use them as outgroup sequences. In January 2018, the genome  
206 sequence of *S. bentonitica* strain VV6 was added to RefSeq and included in the revised version of this  
207 work to increase the taxon sampling.

208

### 209 Computing consensus core- and pan-genomes with GET\_HOMOLOGUES

210 We used GET\_HOMOLOGUES (v05022018) (Contreras-Moreira and Vinuesa, 2013) to compute  
211 clusters of homologous gene families from the input genome sequences, as previously detailed  
212 (Vinuesa and Contreras-Moreira, 2015). Briefly, the source GenBank-formatted files were passed to  
213 get\_homologues.pl and instructed to compute homologous gene clusters by running either our heuristic  
214 (fast) implementation of the bidirectional best-hit (BDBH) algorithm ('-b') to explore the complete  
215 dataset, or the full BDBH, Clusters of Orthologous Groups - triangles (COGtriangles), and OrthoMCL  
216 (Markov Clustering of orthologues, OMCL) algorithms for the different sets of selected genomes, as  
217 detailed in the relevant sections and explained in the GET\_HOMOLOGUES's online manual (ead-  
218 csic-compbio.github.io/get\_homologues/manual/manual.html). PFAM-domain scanning was enabled



219 for the latter runs (-D flag). BLASTP hits were filtered by imposing a minimum of 90% alignment  
220 coverage (-C 90). The directories holding the results from the different runs were then passed to the  
221 auxiliary script `compare_clusters.pl` to compute either the consensus core genome (-t  
222 `number_of_genomes`) or pan-genome clusters (-t 0). The commands to achieve this can be found in the  
223 online tutorial [https://vinuesa.github.io/get\\_phylomarkers/#get\\_homologues-get\\_phylomarkers-](https://vinuesa.github.io/get_phylomarkers/#get_homologues-get_phylomarkers-tutorials)  
224 [tutorials](https://vinuesa.github.io/get_phylomarkers/#get_homologues-get_phylomarkers-tutorials) provided with the distribution.

## 226 **Overview of the computational steps performed by the GET\_PHYLOMARKERS pipeline**

227 Figure 1 presents a flow-chart that summarizes the computational steps performed by the pipeline,  
228 which are briefly described below. For an in-depth description of each step and associated parameters,  
229 as well as for a full version of the pipeline's flow-chart, the reader is referred to the online manual  
230 ([https://vinuesa.github.io/get\\_phylomarkers/](https://vinuesa.github.io/get_phylomarkers/)). The pipeline is primarily intended to run DNA-based  
231 phylogenies ('-R 1 -t DNA') on a collection of genomes from different species of the same genus or  
232 family. However, it can also select optimal markers for population genetics ('-R 2 -t DNA'), when the  
233 source genomes belong to the same species (not shown here). For more divergent genome sequences  
234 the pipeline should be run using protein sequences ('-R 1 -t PROT'). The analyses are started from the  
235 directory holding single-copy core-genome clusters generated either by '`get_homologues.pl -e -t`  
236 `number_of_genomes`' or by '`compare_clusters.pl -t number_of_genomes`'. Note that both the protein  
237 (`faa`) and nucleotide (`fna`) FASTA files for the clusters are required, as detailed in the online tutorial  
238 ([https://vinuesa.github.io/get\\_phylomarkers/#get\\_homologues-get\\_phylomarkers-tutorials](https://vinuesa.github.io/get_phylomarkers/#get_homologues-get_phylomarkers-tutorials)). The former  
239 are first aligned with `clustal-omega` (Sievers et al., 2012) and then used by `pal2nal` (Suyama et al.,  
240 2006) to generate codon alignments. These are subsequently scanned with the `Phi-test` (Bruen et al.,  
241 2005) to identify and discard those with significant evidence for recombinant sequences. Maximum-  
242 likelihood phylogenies are inferred for each of the non-recombinant alignments using by default `IQ-`  
243 `TREE v.1.6.2` (Nguyen et al., 2015), which will perform model selection with `ModelFinder`  
244 (Kalyaanamoorthy et al., 2017) using a subset of models and the '-fast' flag enabled for rapid  
245 computation, as detailed in the online manual. Alternatively, `FastTree v2.1.10` (Price et al., 2010) can  
246 be executed using the '-A F' option, which will estimate phylogenies under the GTR+Gamma model.  
247 `FastTree` was compiled with double-precision enabled for maximum accuracy (see the manual for  
248 details). The resulting gene trees are screened to detect 'outliers' with help of the R package `kdetrees`  
249 (v.0.1.5) (Weyenberg et al., 2014, 2017). It implements a non-parametric test based on the distribution  
250 of tree topologies and branch lengths expected under the multispecies coalescent, identifying those  
251 phylogenies with unusual topologies or branch lengths. The stringency of the test can be controlled  
252 with the `-k` parameter (inter-quartile range multiplier for outlier detection, by default set to the standard  
253 1.5). In a third step, the phylogenetic signal of each gene-tree is computed based on mean branch  
254 support values (Vinuesa et al., 2008), keeping only those above a user-defined mean Shimodaira-  
255 Hasegawa-like (SH-`alrt`) bipartition support (Anisimova and Gascuel, 2006) threshold ('-m 0.75' by  
256 default). To make all the previous steps as fast as possible, they are run in parallel on multi-core  
257 machines using `GNU parallel` (Tange, 2011). The set of alignments passing all filters are concatenated  
258 and subjected to maximum-likelihood (ML) tree searching using by default `IQ-TREE` with model  
259 fitting to estimate the genomic species-tree.

260 The complete GET\_PHYLOMARKERS pipeline is launched with the master script  
261 `run_get_phylomarkers_pipeline.sh`, which calls a subset of auxiliary Bash, Perl and R programs to  
262 perform specific tasks. This architecture allows the user to run the individual steps separately, which  
263 adds convenient flexibility for advanced users (examples provided in the Supplementary Materials).  
264 The pipeline is highly customizable, and the reader is referred to the latest version of the online manual  
265 for the details of each option. However, the default values should produce satisfactory results for most

266 purposes, as these were carefully selected based on the benchmark analysis presented in this work. All  
267 the source code is freely available under the GNU GENERAL PUBLIC LICENSE V3 from  
268 [https://github.com/vinuesa/get\\_phylomarkers](https://github.com/vinuesa/get_phylomarkers). Detailed installation instructions are provided  
269 ([https://github.com/vinuesa/get\\_phylomarkers/blob/master/INSTALL.md](https://github.com/vinuesa/get_phylomarkers/blob/master/INSTALL.md)), along with a hands-on  
270 tutorial ([https://vinuesa.github.io/get\\_phylomarkers/](https://vinuesa.github.io/get_phylomarkers/)). The software has been extensively tested on  
271 diverse Linux distributions (CentOS, Ubuntu and RedHat). In addition, a docker image bundling  
272 GET\_HOMOLOGUES and GET\_PHYLOMARKERS is available at  
273 [https://hub.docker.com/r/csicunam/get\\_homologues/](https://hub.docker.com/r/csicunam/get_homologues/). We recommend running the docker image to  
274 avoid potential trouble with the installation and configuration of diverse dependencies (second party  
275 binaries, as well as Perl and R packages), making it easy to install on any architecture, including  
276 Windows, and to reproduce analyses with exactly the same software.  
277

### 278 **Estimating maximum likelihood and parsimony pan-genome trees from the pan-genome matrix** 279 **(PGM).**

280 The GET\_PHYLOMARKERS package contains auxiliary scripts to perform diverse clustering and  
281 phylogenetic analyses based on the pangenome\_matrix\_t0.\* files returned by the compare\_clusters.pl  
282 script (options '-t 0 -m') from the GET\_HOMOLOGUES suite. In this work, consensus PGMs (Vinuesa  
283 and Contreras-Moreira, 2015) were computed as explained in the online tutorial  
284 ([https://vinuesa.github.io/get\\_phylomarkers/#get\\_homologues-get\\_phylomarkers-tutorials](https://vinuesa.github.io/get_phylomarkers/#get_homologues-get_phylomarkers-tutorials)). These  
285 represent the intersection of the clusters generated by the COGtriangles and OMCL algorithms. Adding  
286 the -T flag to the previous command instructs compare\_clusters.pl to compute a Wagner (multistate)  
287 parsimony tree from the pan-genome matrix, launching a tree search with 50 taxon jumbles with pars  
288 from the PHYLIP (Felsenstein, 2004b) package (v.3.69). A more thorough and customized ML or  
289 parsimony analysis of the PGM can be performed with the aid of the auxiliary script  
290 estimate\_pangenome\_phylogenies.sh, bundled with GET\_PHYLOMARKERS. By default this script  
291 performs a ML tree-search using IQ-TREE v1.6.2 (Nguyen et al., 2015). It will first call ModelFinder  
292 (Kalyaanamoorthy et al., 2017) using the JC2 and GTR2 base models for binary data, the latter  
293 accounting for unequal state frequencies. The best fitting base model + ascertain bias correction +  
294 among-site rate variation parameters are selected using the Akaike Information Criterion (AIC). IQ-  
295 TREE (Nguyen et al., 2015) is then called to perform a ML tree search under the selected model with  
296 branch support estimation. These are estimated using approximate Bayesian posterior probabilities  
297 (aBypp), a popular single branch test (Guindon et al., 2010), as well as the recently developed ultrafast-  
298 bootstrap2 (UFBoot2) test (Hoang et al., 2017). In addition, the user may choose to run a parsimony  
299 analysis with bootstrapping on the PGM, as detailed in the online manual and illustrated in the tutorial.  
300 Note however, that the parsimony search with bootstrapping is much slower than the default ML  
301 search.  
302

### 303 **Unsupervised learning methods for the analysis of pairwise average nucleotide (ANI) and** 304 **aminoacid (AAI) identity matrices**

305 The GET\_HOMOLOGUES distribution contains the plot\_matrix\_heatmap.sh script which generates  
306 ordered heatmaps with attached row and column dendrograms from squared tab-separated numeric  
307 matrices. These can be presence/absence PGM matrices or similarity / identity matrices, as those  
308 produced with the get\_homologues -A option. Optionally, the input cgANIb matrix can be converted to  
309 a distance matrix to compute a neighbor joining tree, which makes the visualization of relationships in  
310 large ANI matrices easier. Recently added functionality includes reducing excessive redundancy in the  
311 tab-delimited ANI matrix file (-c max\_identity\_cut-off\_value) and sub-setting the matrix with regular  
312 expressions, to focus the analysis on particular genomes extracted from the full cgANIb matrix. From

313 version 1.0 onwards, the mean silhouette-width (Rousseeuw, 1987) goodness of clustering statistics to  
314 determine the optimal number of clusters automatically. The script currently depends on the R packages  
315 ape (Popescu et al., 2012), dendextend (<https://cran.r-project.org/package=dendextend>), factoextra  
316 (<https://cran.r-project.org/package=factoextra>) and gplots ([https://CRAN.R-](https://CRAN.R-project.org/package=gplots)  
317 [project.org/package=gplots](https://CRAN.R-project.org/package=gplots)).  
318  
319

## 320 RESULTS

321

### 322 **Ten new complete genome assemblies for the Mexican environmental *Stenotrophomonas*** 323 ***maltophilia* complex isolates previously classified as genospecies 1 (Smc1) and 2 (Smc2).**

324 In this study we report the sequencing and assembly of five isolates each from the genospecies 1  
325 (Smc1) and 2 (Smc2) recovered from rivers in Central Mexico, previously reported in our extensive  
326 MLSA study of the genus *Stenotrophomonas* (Ochoa-Sánchez and Vinuesa, 2017). All assemblies  
327 resulted in a single chromosome with gaps. No plasmids were detected. A summary of the annotated  
328 features for each genome are presented in Table 1. Assembly details for each genome are provided in  
329 supplementary Table S1.  
330

### 331 **Rapid phylogenetic exploration of *Stenotrophomonas* genome sequences available at NCBI's** 332 **RefSeq repository running GET\_PHYLOMARKERS in fast runmode**

333 A total of 170 *Stenotrophomonas* and 7 *Xanthomonas* reference genomes were retrieved from RefSeq  
334 (see methods). Figure 2A depicts parallel density plots showing the distribution of the number of  
335 fragments for the *Stenotrophomonas* assemblies at the Complete ( $n = 16$ ), Chromosome ( $n = 3$ ),  
336 Scaffold ( $n = 63$ ) and Contig ( $n = 88$ ) finishing levels. The distributions have conspicuous long tails,  
337 with an overall mean and median number of fragments of  $\sim 238$  and  $\sim 163$ , respectively. The table insets  
338 in Fig. 2A provide additional descriptive statistics of the distributions. A first GET\_HOMOLOGUES  
339 run was launched using this dataset ( $n = 177$ ) with two objectives: *i*) to test its performance with a  
340 relatively large set of genomes and *ii*) to get an overview of their evolutionary relationships to select a  
341 non-redundant set of those with the best assemblies. For this analysis, GET\_HOMOLOGUES was run  
342 in its “fast-BDBH” mode (-b), on 60 cores (-n 60; AMD Opteron™ Processor 6380, 2500.155 MHz),  
343 and imposing a stringent 90% coverage cut-off for BLASTP alignments (-C 90), excluding  
344 inparalogues (-e). This analysis took 1h:32m:13s to complete and identified 132 core genes. These  
345 were fed into the GET\_PHYLOMARKERS pipeline, which was executed using a default FastTree  
346 search with the following command line: run\_get\_phylomarkers\_pipeline.sh -R 1 -t DNA -A F, which  
347 took 8m:1s to complete on the same number of cores. Only 79 alignments passed the Phi  
348 recombination test. Thirteen of them failed to pass the downstream kdetree test. The phylogenetic  
349 signal test excluded nine additional loci with average SH-*alrt* values  $< 0.70$ . Only 57 alignments  
350 passed all filters and were concatenated into a supermatrix of 38,415 aligned residues, which were  
351 collapsed to 19,129 non-gapped and variable sites. A standard FastTree maximum-likelihood tree-  
352 search was launched with the command: run\_get\_phylomarkers\_pipeline.sh -R 1 -t DNA -A F'. The  
353 resulting phylogeny ( $\ln L = -475237.540$ ) is shown in supplementary Figure S1. Based on this tree and  
354 the level of assembly completeness for each genome (Fig. 2A), we decided to discard those with  $> 300$   
355 contigs (Fig. 2B). This resulted in the loss of 19 genomes labeled as *S. maltophilia*. However, we  
356 retained *S. pictorum* JCM 9942, a highly fragmented genome with 829 contigs (Patil et al., 2016) to  
357 maximize taxon sampling. Several *S. maltophilia* subclades contained identical sequences (Fig. S1) and  
358 were trimmed, retaining only the assembly with the lowest numbers of scaffolds or contigs.  
359



## 360 **Selection of a stringently defined set of orthologous genes using GET\_HOMOLOGUES**

361 After the quality and redundancy filtering described in the previous section, 109 reference genomes  
362 (102 *Stenotrophomonas* + 7 *Xanthomonas*) were retained for more detailed investigation. Table S2  
363 provides an overview of them. To this set we added the 10 new genomes reported in this study (Table  
364 1). Figure 2B depicts parallel density plots summarizing the distribution of number of contigs/scaffolds  
365 in the selected reference genomes and the new genomes for the Mexican environmental Smc isolates  
366 previously classified as genospecies 1 (Smc1) and 2 (Smc2) (Ochoa-Sánchez and Vinuesa, 2017). A  
367 high stringency consensus core-genome containing 239 gene families was computed as the intersection  
368 of the clusters generated by the BDBH, COG-triangles and OMCL algorithms (Fig 3A).  
369

## 370 **GET\_PHYLOMARKERS in action: benchmarking the performance of FastTree and IQ-TREE** 371 **to select top-scoring markers for phylogenomics**

372 The set of 239 consensus core-genome clusters (Fig. 3A) was used to launch multiple instances of the  
373 GET\_PHYLOMARKERS pipeline to evaluate the phylogenetic performance of FastTree (FT; v2.1.10)  
374 and IQ-TREE (IQT; v1.6.2), two popular fast maximum-likelihood (ML) tree searching algorithms.  
375 Our benchmark was designed to compare: *i*) the execution times of the FT vs. IQT runs under default  
376 (FTdef, IQTdef) and thorough (FThigh, IQThigh) search modes (see methods and online manual for  
377 their parameterization details); *ii*) the phylogenetic resolution (average support values) of gene trees  
378 estimated by FT and IQT under both search modes; *iii*) the rank of lnL scores of the gene trees found in  
379 those searches for each locus; *iv*) the distribution of consensus values of each node in majority rule  
380 consensus trees computed from the gene trees found by each search type; *v*) the distribution of edge-  
381 lengths in the species-trees computed by each search type. The results of these analyses are  
382 summarized in Table 2 and in Figure 3. The first steps of the pipeline (Fig. 1) comprise the generation  
383 of codon alignments and their analysis to identify potential recombination events. Only 127 alignments  
384 (53.14 %) passed the Phi-test (Table 2). Phylogenetic analyses start downstream of the recombination  
385 test (Fig. 1). The computation times required by the two algorithms and search intensity levels were  
386 significantly different (Kruskal-Wallis,  $p < 2.2e-16$ ), FastTree being always the fastest, and displaying  
387 the lowest dispersion of compute times across trees (Fig. 3B). This is not surprising, as IQT searches  
388 involved selecting the best substitution model among a range of base models (see methods and online  
389 manual) and fitting additional parameters (+G+ASC+I+F+R) to account for heterogeneous base  
390 frequencies and rate-variation across sites. In contrast, FT searches just estimated the parameter values  
391 for the general time-reversible (GTR) model, and among-site rate variation was modeled fitting a  
392 gamma distribution with 20 rate categories (+G), as summarized in Table 2. Similar numbers of  
393 “outlier” trees (range 18:22) were detected by the kdetrees-test in the four search types (Table 2).  
394 However, the distributions of SH-*alrt* support values are strikingly different for both search algorithms  
395 (Wilcoxon,  $p < 2.2e-16$ ), revealing that gene-trees found by IQT have a much lower average support  
396 than those found by FT (Fig 3C). Consequently, the former searches were significantly more efficient  
397 to identify gene trees with low average branch support values (Table 2 and Fig. 3C). This result is in  
398 line with the well-established fact that poorly fitting and under-parameterized models produce less  
399 reliable tree branch lengths and overestimate branch support (Posada and Buckley, 2004), implying that  
400 the FT phylogenies may suffer from clade over-credibility. These results demonstrate that: *i*) FT-based  
401 searches are significantly faster than those performed with IQT, and *ii*) that IQT has a significantly  
402 higher discrimination power for phylogenetic signal than FT. Due to the fact that the number of top-  
403 scoring alignments selected by the two algorithms for concatenation is notably different (Table 2), the  
404 lnL scores of the resulting species-trees are not comparable (Table 2). Therefore, in order to further  
405 evaluate the quality of the gene-trees found by the four search strategies, we performed an additional  
406 benchmark under highly standardized conditions, based on the 105 optimal alignments that passed the

407 kdetrees-test in the IQThigh search (Table 2). Gene trees were estimated for each of these alignments  
408 using the four search strategies (FTdef, IQTdef, FThigh and IQThigh) and their lnL scores ranked for  
409 each gene tree. An association analysis (deviation from independence in a multi-way chi-squared test)  
410 was performed on the lnL ranks (1 to 4, coding for highest to lowest lnL scores, respectively) attained  
411 by each search type for each gene tree. As shown in Fig. 3D, the IQThigh search was the winner,  
412 attaining the first rank (highest lnL score) in 76/105 of the searches (72.38%), way ahead of the number  
413 of FThigh (26%), and IQTdef (0.009%) searches that ranked in the first position (highest lnL score for  
414 a particular alignment). A similar analysis performed on the full set of input alignments ( $n = 239$ )  
415 indicated that when operating on an unfiltered set, the difference in performance was even more  
416 striking, with IQT-based searches occupying  $> 97\%$  of the first rank positions (data not shown). These  
417 results highlight two points: *i*) the importance of proper model selection and thorough tree searching in  
418 phylogenetic inference and *ii*) that IQT generally finds better trees than FT. Finally, we evaluated  
419 additional phylogenetic attributes of the species-trees computed by each search type, either as the  
420 majority rule consensus (mjrc) tree of top-scoring gene-trees, or as the tree estimated from the  
421 supermatrices of concatenated alignments. Figure 3E shows the distribution of mjrc values of the mjrc  
422 trees computed by each search type, which can be interpreted as a proxy for the level phylogenetic  
423 congruence among the source trees. These values were significantly higher for the IQT than in the FT  
424 searches (Kruskal-Wallis,  $p = 0.027$ ), with a higher number of 100% mjrc clusters found in the former  
425 than in the latter type of trees (Fig. 3E). An analysis of the distribution of edge-lengths of the species-  
426 trees inferred from the concatenated alignments revealed that those found in IQT searches had  
427 significantly (Kruskal-Wallis,  $p = 1e-07$ ) shorter edges (branches) than those estimated by FT (Fig. 3F).  
428 This highlights again the importance of adequate substitution models for proper edge-length estimation.  
429 Tree-lengths (sum of edge lengths) of the species-trees found in IQT-based searches are about 0.63  
430 times shorter than those found by FT (Fig. S2). As a final exercise, we computed the Robinson-Foulds  
431 (RF) distances of each gene tree found in a given search type to the species tree inferred from the  
432 corresponding supermatrix. The most striking result of this analysis was that no single gene-tree had  
433 the same topology as the species tree inferred from the concatenated top-scoring alignments (Fig. S3).  
434

### 435 **Effect of tree-search intensity on the quality of the species trees found by IQT-REE and FastTree**

436 Given the astronomical number of different topologies that exist for 119 terminals, we decided to  
437 evaluate the effect of tree-search thoroughness on the quality of the trees found by FT and IQT,  
438 measured as their log-likelihood (lnL) score. To make the results comparable across search algorithms,  
439 we used the supermatrix of 55 top-scoring markers (25,896 variable, non-gapped sites) selected by the  
440 IQThigh run (Table 2). One thousand FT searches were launched from the same number of random  
441 topologies computed with the aid of a custom Perl script. In addition, a standard FT search was started  
442 from the default BioNJ tree. All these searches were run in “thorough” mode (-quiet -nt -gtr -bionj  
443 -slow -slownni -gamma -mlacc 3 -spr 16 -sprlength 10) on 50 cores. The resulting lnL profile for this  
444 search is presented in Figure 4A, which reached a maximal score of -717195.373. This is 121.281 lnL  
445 units better than the score of the best tree found in the search started from the BioNJ seed tree (lnL  
446 -717316.654, lower discontinuous blue line). In addition, 50 independent tree searches were run with  
447 IQ-TREE under the best fitting model previously found (Table 2), using the shell loop command (# 5)  
448 provided in the Supplementary Material. The corresponding lnL profile of this search is shown in Fig.  
449 4B, which found a maximum-scoring tree with a score of -707932.468. This is only 8.105 lnL units  
450 better than the worst tree found in that same search (Fig. 4B). Importantly, the best tree found in the  
451 IQT-search is 9262.905 lnL units better than that of the best tree found in the FT search, despite the much  
452 higher number of seed trees used for the latter. This result clearly demonstrates the superiority of the  
453 IQ-TREE algorithm for ML tree searching. Based on this evidence, and that presented in the previous

454 section (Table 2; Fig. 3), IQ-TREE was chosen as the default tree-search algorithm used by  
455 GET\_PHYLOMARKERS. The Robinson-Foulds distance between both trees was 46.  
456

### 457 **A robust genomic species phylogeny for the genus *Stenotrophomonas*: taxonomic implications and** 458 **identification of multiple misclassified genomes**

459 Figure 5 displays the best ML phylogeny found in the IQ-TREE search (Fig. 4B) described in the  
460 previous section. This is a highly resolved phylogeny. All bipartitions have an approximate Bayesian  
461 posterior probability (aBypp)  $p \geq 0.95$ . It was rooted at the branch subtending the *Xanthomonas* spp.  
462 clade, used as an outgroup. A first taxonomic inconsistency revealed by this phylogeny is the placement  
463 of *S. panacihumi* within the latter clade, making the genus *Stenotrophomonas* paraphyletic. It is worth  
464 noting that *S. panacihumi* is a non-validly described, and poorly characterized species (Yi et al., 2010).  
465 The genus *Stenotrophomonas*, as currently defined, and excluding *S. panacihumi*, consists of two major  
466 clades, labeled as I and II in Fig. 5, as previously defined (Ochoa-Sánchez and Vinuesa, 2017).  
467 Clade I groups environmental isolates, recovered from different ecosystems, mostly soils and plant  
468 surfaces, classified as *S. ginsengisoli* (Kim et al., 2010), *S. koreensis* (Yang et al., 2006), *S.*  
469 *daejeonensis* (Lee et al., 2011), *S. nitritireducens* (Finkmann et al., 2000), *S. acidaminiphila* (Labat et  
470 al., 2002), *S. humi* and *S. terrae* (Heylen et al., 2007). The recently described *S. pictorum* (Ouattara et  
471 al., 2017) is also included in clade I. These are all rather poorly studied species, for which only one or a  
472 few strains have been considered in the corresponding species description or to study particular aspects  
473 of their biology. None of these species have been reported as opportunistic pathogens, but some contain  
474 promising strains for plant growth-promotion and bio-remediation. Particularly notorious are the  
475 disproportionally long terminal branches (heterotachy) of *S. ginsengisoli* and *S. koreensis* (Fig. 5). The  
476 potential distortion of these long branches on the estimated phylogeny needs to be evaluated in future  
477 work.

478 Clade II contains the species *S. rhizophila* (Berg et al., 2002), *S. chelatiphaga* (Kaparullina et al.,  
479 2009), the recently described *S. bentonitica* (Sánchez-Castro et al., 2017), along with multiple species  
480 and genospecies lumped in the *S. maltophilia* complex (Smc; shaded area in Fig. 5) (Berg and  
481 Martinez, 2015; Svensson-Stadler et al., 2012). The Smc includes the validly described *S. maltophilia*  
482 (Palleroni and Bradbury, 1993) and *S. pavanii* (Ramos et al., 2011) (collapsed subclades Sm6 and Sm2,  
483 respectively, located within the clade labeled as *S. maltophilia sensu lato* in Figure 5), along with at  
484 least four undescribed genospecies (Sgn1-Sgn4) recently identified in our MLSA study of the genus  
485 (Ochoa-Sánchez and Vinuesa, 2017). In light of this phylogeny, we discovered 14 misclassified RefSeq  
486 genome sequences (out of 119; ~11.76 %), 12 of them labeled as *S. maltophilia*. These genomes are  
487 highlighted with black arrows in Figure 5. The phylogeny also supports the classification, either as a  
488 validly published species, or as new genospecies, of 8 (~ 6.72 %) additional RefSeq genomes (gray  
489 arrows) lacking a species assignment in the RefSeq record, as summarized in Table 3. In addition, the  
490 phylogeny resolved 13 highly supported lineages (aBypp > 0.95) within the *S. maltophilia sensu lato*  
491 (Smsl) cluster, shown as collapsed clades. They have a core-genome average nucleotide identity > 96  
492 % (Fig. 5). These lineages may represent 11 additional species in the Smsl clade, as detailed in  
493 following sections. Supplementary Figure S4 shows the non-collapsed version of the species-tree  
494 displayed in Figure 5.

495 No genome sequences, nor MLSA data are available for the recently described *S. tumulicola* (Handa et  
496 al., 2016).

497

### 498 **Pan-genome phylogenies for the genus *Stenotrophomonas* recover the same species clades as the** 499 **core-genome phylogeny**

500 A limitation of core-genome phylogenies is that they are estimated from the small fraction of single-  
501 copy genes shared by all organisms under study. Genes encoding adaptive traits relevant for niche-  
502 differentiation and subsequent speciation events typically display a lineage-specific distribution. Hence,  
503 phylogenetic analysis of pan-genomes, based on their differential gene-composition profiles, provide a  
504 complementary, more resolved and often illuminating perspective on the evolutionary relationships  
505 between species.

506 A consensus pan-genome matrix (PGM) containing 29,623 clusters was computed from the intersection  
507 of the clusters generated by the COG-triangles and OMCL algorithms (Figure 6). This PGM was  
508 subjected to ML tree searching using the binary and morphological models implemented in IQ-TREE  
509 for phylogenetic analysis of discrete characters with the aid of the `estimate_pangenome_phylogenies.sh`  
510 script bundled with GET\_PHYLOMARKERS (Fig. 1). As shown in the tabular inset of Figure 6, the  
511 binary GTR2+FO+R4 model was by large the best-fitting one (with the smallest AIC and BIC values).  
512 Twenty five independent IQ-TREE searches were performed on the consensus PGM with the best-  
513 fitting model. The best tree found is presented in Figure 6, rooted with the *Xanthomonas* spp. outgroup  
514 sequences. It depicts the evolutionary relationships of the 119 genomes based on their gene content  
515 (presence-absence) profiles. The numbers on the nodes indicate the approximate Bayesian posterior  
516 probabilities (aBypp) / UFBoot2 support values (see methods). The same tree, but without collapsing  
517 clades, is presented in the supplementary figure S5. This phylogeny resolves exactly the same species-  
518 like clades highlighted on the core-genome phylogeny presented in Figure 5, which are also grouped in  
519 the two major clades I and II. These are labeled with the same names and color-codes, for easy cross-  
520 comparison. However, there are some notorious differences in the phylogenetic relationships between  
521 species on both trees, like the placement of *S. panacihumi* outside of the *Xanthomonas* clade, and the  
522 sister relation of genospecies 3 (Sgnp3) to the *S. maltophilia sensu lato* clade. These same relationships  
523 were found in a multi-state (Wagner) parsimony phylogeny of the PGM shown in Supplementary  
524 Figure S6. In summary, all core-genome and pan-genome analyses presented consistently support our  
525 previous claim that the five genospecies defined in our MLSA study represent distinct species and  
526 support the existence of multiple cryptic species within the Smsl clade, as defined in Figure 5.

### 528 **Application of non-supervised learning approaches to BLAST-based core-genome average** 529 **nucleotide distance (cgANDb) and Gower pan-genome distances (pgGdist) provide statistically-** 530 **consistent results for prokaryotic species delimitation**

531 The final goal of any geno-taxonomic study is to identify species-like clusters. These should consist of  
532 monophyletic groups identified on genome trees that display average genome identity (gANI) values >  
533 94 %, based on a widely accepted cutoff-value (Rosselló-Mora and Amann, 2015). In this section we  
534 searched for such species-clusters within the taxonomically problematic *Stenotrophomonas maltophilia*  
535 complex (Smc). Our core- and pan-genome phylogenies consistently identified potential species-clades  
536 within the Smc that grouped exactly the same strains (compare Figs. 5 and 6). We additionally  
537 performed a cluster analysis of core-genome ANI values computed from the pairwise BLASTN  
538 alignments (cgANIb) used to define OMCL core-genome clusters for the 86 Smc genomes analyzed in  
539 this study. The resulting cgANIb matrix was then converted to a distance matrix (cgANDb = 100 % -  
540 cgANIb) and clustered with the aid of the `plot_matrix_heatmap.sh` script from the  
541 GET\_HOMOLOGUES suite. Figure 7 shows the resulting tree, which resolves 16 clusters within the  
542 Smc at a conservative cgANDb cutoff value of 5% (cgANIb = 95%). At this distance level, the four  
543 genospecies labeled as Sgn1-Sgn4 on Figure 5 are resolved as five clusters because the most divergent  
544 Sgn1 genome (ESTM1D\_MKCIP4\_1) is split as a separate lineage. This is the case also at cgANDb =  
545 6 (Fig. 7), reason why this strain most likely represents a sixth genospecies. All these genospecies are  
546 very distantly related to the large *S. maltophilia s. lato* cluster, which gets split into 11 sub-clusters at



547 the conservative cgANDB = 5 % cutoff. Thirteen clusters are resolved at the 4 % threshold, and a  
548 minimum of seven at the 6 % level (cgANIb = 94%), as shown by the dashed lines (Fig. 7). These  
549 results strongly suggest that the *S. maltophilia sensu lato* clade (Fig. 5) actually comprises multiple  
550 species. The challenging question is how many? In an attempt to find a statistically-sound answer, we  
551 applied an unsupervised learning approach based on the evaluation of different goodness of clustering  
552 statistics to determine the optimal number of clusters ( $k$ ) for the cgANDB matrix. The gap-statistic and  
553 a parametric, model-based cluster analysis yielded  $k$  values  $\geq 35$  (data not shown). These values seem  
554 too high for this dataset, as they correspond to a gANI value  $> 98\%$ . However, the more conservative  
555 average silhouette width (ASW) method (Kaufman and Rousseeuw, 1990) identified an optimal  $k = 19$   
556 (inset in Fig 7) for the complete set of Smc genomes. This number of species-like clusters is much  
557 more reasonable for this data set, as it translates to a range of cgANDB between 4.5 and 4.7 (cgANIb  
558 range: 95.5% - 95.3%). Close inspection of the ASW profile reveals that the first peak is found at  $k =$   
559 13, which has an almost identical ASW as that of the maximal value and maps to a cgANDg = 5.7  
560 (cgANIb of 94.3%). In summary, the range of reasonable numbers of clusters proposed by the ASW  
561 statistic ( $k = 13$  to  $k = 19$ ) corresponds to cgANDB values in the range of 5.7% - 4.5% (cgANIb range:  
562 94.3% - 95.5%), which fits well with the new gold-standard for species delimitation (gANI  $> 94\%$ ),  
563 established in influential works (Konstantinidis and Tiedje, 2005; Richter and Rossello-Mora, 2009).  
564 We noted however, that at a cgANDB = 4.1% (cgANIb = 95.9 %) the strain composition of the clusters  
565 was 100% concordant with the monophyletic subclades shown in the core-genome (Fig. 5) and pan-  
566 genome (Fig. 6) phylogenies. Importantly, at this cutoff, the length of the branches subtending each  
567 cluster is maximal, both on the core-genome phylogeny (Fig. 5) and on the cgANDB cladogram (Fig.  
568 7). Based on the combined and congruent evidence provided by these complementary approaches, we  
569 can safely conclude that: *i*) the Smc genomes analyzed herein may actually comprise up to 19 or 20  
570 different species-like lineages, and *ii*) that only the strains grouped in the cluster labeled as Sm6 in  
571 Figs. 5, 6 and 7 should be called *S. maltophilia*. The latter is the most densely sampled species-like  
572 cluster ( $n = 19$ ) and includes ATCC 13637<sup>T</sup>, the type strain of the species.  
573

#### 574 **On the ecology and other biological attributes of the species-like clusters in the** 575 ***Stenotrophomonas maltophilia* complex**

576 In this final section we present a brief summary of the ecological attributes reported for selected  
577 members of the species-like clusters resolved within the Smc (Figs 5 and 7). The four unnamed  
578 genospecies (Sgn1-Sgn4) group mainly environmental isolates. This is consistent with our previous  
579 evolutionary and ecological analyses of a comprehensive multilocus dataset of the genus (Ochoa-  
580 Sánchez and Vinuesa, 2017). In that study only Mexican environmental isolates were found to be  
581 members of the newly discovered genospecies Sgn1 and Sgn2 (named as Smc1 and Smc2,  
582 respectively). In this work we discovered that the recently sequenced maize root isolate AA1 (Niu et  
583 al., 2017), misclassified as *S. maltophilia*, clusters tightly with the Sgn1 strains (Fig. 5). The *S.*  
584 *maltophilia sensu lato* clade is split into 13 or 14 groups based on cgANDB (Fig 7). Sm6 forms the  
585 largest cluster, grouping mostly clinical isolates related to the type strain *S. maltophilia* ATCC 13637<sup>T</sup>,  
586 like the model strain K279a (Crossman et al., 2008), ISMMS4 (Pak et al., 2015), 862\_SMAL,  
587 1149\_SMAL and 1253\_SMAL (Roach et al., 2015), as well as EPM1 (Sassera et al., 2013), recovered  
588 from the human parasite *Giardia duodenalis*. However, this group also comprises some environmental  
589 isolates like BurE1, recovered from a bulk soil sample (Youenou et al., 2015). In summary, cluster Sm6  
590 holds the *bona fide S. maltophilia* strains (*sensu stricto*), which may be well-adapted to associate with  
591 different eukaryotic hosts and cause opportunistic infections in humans. Cluster Sm4a contains the  
592 model strain D574 (Lira et al., 2012) along with four other clinical isolates (Conchillo-Solé et al., 2015)  
593 and therefore may represent a second clade enriched in strains with high potential to cause



594 opportunistic pathogenic infections in humans. Noteworthy, this group is distantly related to Sm6 (Figs.  
595 5 and 7). Cluster Sm4b is closely related to Sm4a based on the pan-genome phylogeny and the  
596 cgANDd cladogram (Figs. 6 and 7). It groups the Brazilian rhizosphere-colonizing isolate JV3, the  
597 Chinese highly metal tolerant strain TD3 (Ge and Ge, 2016) and strain As1, isolated from the Asian  
598 malaria vector *Anopheles stephensi* (Hughes et al., 2016). The lineage Sm3 holds eight isolates of  
599 contrasting origin, including the Chinese soil isolate DDT-1, capable of using DDT as the sole source  
600 of carbon and energy (Pan et al., 2016), as well as clinical isolates like 1162\_SMAL (Roach et al.,  
601 2015) and AU12-09, isolated from a vascular catheter (Zhang et al., 2013), and environmental isolates  
602 like SmF22, Sm32COP and SmSOFb1, isolated from different manures in France (Bodilis et al., 2016).  
603 Cluster Sm2 groups the *S. pavanii* strains, including the type strain DSM\_25135<sup>T</sup>, isolated from the  
604 stems of sugar cane in Brazil (Ramos et al., 2011), together with the clinical isolates ISMMS6 and  
605 ISMMS7, that carry mutations conferring quinolone resistance and causing bacteremia (Pak et al.,  
606 2015), and strain C11, recovered from pediatric cystic fibrosis patients (Ormerod et al., 2015). Cluster  
607 Sm5 includes two strains recovered from soils, ATCC 19867 which was first classified as  
608 *Pseudomonas hibiscicola*, and later reclassified as *S. maltophilia* based on MLSA studies  
609 (Vasileuskaya-Schulz et al., 2011), and the African strain BurA1, isolated from bulk soil samples  
610 collected in sorghum fields in Burkina Faso (Youenou et al., 2015). Cluster Sm9 holds clinical isolates,  
611 like 131\_SMAL, 424\_SMAL and 951\_SMAL (Roach et al., 2015). Its sister group is Sm10. It holds 9  
612 strains of contrasting geographic and ecological provenances, ranging from Chinese soil and plant-  
613 associated bacteria like the rice-root endophyte RR10 (Zhu et al., 2012), the grassland-soil tetracycline  
614 degrading isolate DT1 (Naas et al., 2008), and strain B418, isolated from a barley rhizosphere and  
615 displaying plant-growth promotion properties (Wu et al., 2015), to clinical isolates (22\_SMAL,  
616 179\_SMAL, 453\_SMAL, 517\_SMAL) collected and studied in the context of a large genome  
617 sequencing project carried out at the University of Washington Medical Center (Roach et al., 2015).  
618 Cluster Sm11 tightly groups the well-characterized poplar endophyte R551-3, which is a model plant-  
619 growth-promoting bacterium (Alavi et al., 2014; Ryan et al., 2009; Taghavi et al., 2009) and SB01,  
620 cultured from the gut of the olive fruit fly *Bactrocera oleae* (Blow et al., 2016). Cluster Sm 12 contains  
621 the environmental strain SKA14 (Adamek et al., 2014), along with the clinical isolates ISMMS3 (Pak  
622 et al., 2015) and 860\_SMAL (Roach et al., 2015). Sm1, Sm7 and Sm8 each hold a single strain.  
623 The following conclusions can be drawn from this analysis: *i*) the species-like clusters within the *S.*  
624 *maltophilia sensu lato* (Smsl) clade (Fig. 5) are enriched in opportunistic human pathogens, when  
625 compared with the Smc clusters Sgn1-Sgn4; *ii*) most Smsl clusters also contain diverse non-clinical  
626 isolates isolated from a wide range of habitats, demonstrating the great ecological versatility found  
627 even within specific Smsl clusters like Sm3 or Sm10; *iii*) taken together, these observations strongly  
628 suggest that the Smc species-like clusters are all of environmental origin, with the potential for the  
629 opportunistic colonization of diverse human organs. This potential may be particularly high in certain  
630 lineages, like in *S. maltophilia sensu stricto* (Sm6) or Sm4a, both enriched in clinical isolates.  
631 However, a much denser sampling of genomes and associated phenotypes is required for all clusters to  
632 be able to identify statistically sound associations between them.

633  
634

## 635 DISCUSSION

636

637 In this study we developed and benchmarked GET\_PHYLOMARKERS, an open-source,  
638 comprehensive, and easy-to-use software package for phylogenomics and microbial genome taxonomy.  
639 Programs like amphora (Wu and Eisen, 2008) or phylosift (Darling et al., 2014) allow users to infer  
640 genome-phylogenies from huge genomic and metagenomic datasets by scanning new sequences against

641 a reference database of conserved protein sequences to establish the phylogenetic relationships between  
642 the query sequences and database hits. The first program searches the input data for homologues to a  
643 set of 31 highly conserved proteins used as phylogenetic markers. Phylosift is more oriented towards  
644 the phylogenetic analysis of metagenome community composition and structure. Other approaches  
645 have been developed to study large populations of a single species. These are based on the  
646 identification of single nucleotide polymorphisms in sequence reads produced by high-throughput  
647 sequencers, using either reference-based or reference-free approaches, and subjecting them to  
648 phylogenetic analysis (Timme et al., 2013). The GET\_PHYLOMARKERS software suite was designed  
649 with the aim of identifying orthologous clusters with optimal attributes for phylogenomic analysis and  
650 accurate species-tree inference. It also provides tools to infer phylogenies from pan-genomes, as well as  
651 non-supervised learning approaches for the analysis of overall genome relatedness indices (OGRIs) for  
652 geno-taxonomic studies of multiple genomes. These attributes make GET\_PHYLOMARKERS unique  
653 in the field.

654  
655 It is well-established that the following factors strongly affect the accuracy of genomic phylogenies: *i*)  
656 correct orthology inference; *ii*) multiple sequence alignment quality; *iii*) presence of recombinant  
657 sequences; *iv*) loci producing anomalous phylogenies, which may result for example from horizontal  
658 gene transfer, differential loss of paralogues between lineages and *v*) amount of the phylogenetic signal.  
659 GET\_PHYLOMARKERS aims to minimize the negative impact of potentially problematic or poorly  
660 performing orthologous clusters by explicitly considering and evaluating these factors. Orthologous  
661 clusters were identified with GET\_HOMOLOGUES (Contreras-Moreira and Vinuesa, 2013) because of  
662 its distinctive capacity to compute high stringency clusters of single-copy orthologs. In this study we  
663 used a combination of BLAST alignment filtering imposing a high (90%) query coverage threshold,  
664 PFAM-domain composition scanning and calculation of a consensus core-genome from the  
665 orthologous gene families produced by three clustering algorithms (BDBH, COGtriangles and OMCL)  
666 to minimize errors in orthology inference. Multiple sequence alignments were generated with  
667 CLUSTAL-OMEGA (Sievers et al., 2012), a state-of-the-art software under constant development,  
668 capable of rapidly aligning hundreds of protein sequences with high accuracy, as reported in recent  
669 benchmark studies (Le et al., 2017; Sievers and Higgins, 2018). GET\_PHYLOMARKERS generates  
670 protein alignments and uses them to compute the corresponding DNA-alignments, ensuring that the  
671 codon structure is always properly maintained. Recombinant sequences have been known for a long  
672 time to strongly distort phylogenies because they merge independent evolutionary histories into a  
673 single lineage. Recombination erodes the phylogenetic signal and misleads classic treeing algorithms,  
674 which assume a single underlying history (Didelot and Maiden, 2010; Martin, 2009; Pease and Hahn,  
675 2013; Posada and Crandall, 2002; Schierup and Hein, 2000; Turrientes et al., 2014). Hence, the first  
676 filtering step in the pipeline is the detection of putative recombinant sequences using the very fast,  
677 sensitive and robust  $\phi(w)$  statistic (Bruen et al., 2005). The genus *Stenotrophomonas* has been  
678 previously reported to have high recombination rates (Ochoa-Sánchez and Vinuesa, 2017; Yu et al.,  
679 2016). It is therefore not surprising that the  $\phi(w)$  statistic detected significant evidence for  
680 recombination in up to 47% of the orthologous clusters. The non-recombinant sequences are  
681 subsequently subjected to maximum-likelihood phylogenetic inference to identify anomalous trees  
682 using the non-parametric *kdetrees* statistic (Weyenberg et al., 2014, 2017). The method estimates  
683 distributions of phylogenetic trees over the "tree space" expected under the multispecies-coalescent,  
684 identifying outlier trees based on their topologies and branch lengths in the context of this distribution.  
685 Since this test is applied downstream of the recombination analysis, only a modest, although still  
686 significant proportion (14%-17%) of outlier trees were detected (Table 2). The next step determines the  
687 phylogenetic signal content of each gene tree (Vinuesa et al., 2008). It has been previously established

688 that highly informative trees are less prone to get stuck in local optima (Money and Whelan, 2012).  
689 They are also required to properly infer divergence at the deeper nodes of a phylogeny (Salichos and  
690 Rokas, 2013), and to get reliable estimates of tree congruence and branch support in large concatenated  
691 datasets typically used in phylogenomics (Shen et al., 2017). We found that IQ-TREE-based searches  
692 allowed a significantly more efficient filtering of poorly resolved trees than FastTree. This is likely due  
693 to the fact that the former fits more sophisticated models (with more parameters) to better account for  
694 among-site rate variation. Under-parameterized and poorly fitting substitution models partly explain the  
695 apparent overestimation of bipartition support values done by FastTree. This is also the cause of the  
696 poorer performance of FastTree, which finds gene trees that generally have lower  $\ln L$  scores than those  
697 found by IQ-TREE. A recent comparison of the performance of four fast ML phylogenetic programs  
698 using large phylogenomic data sets identified IQ-TREE (Nguyen et al., 2015) as the most accurate  
699 algorithm. It consistently found the highest-scoring trees. FastTree (Price et al., 2010) was, by large, the  
700 fastest program evaluated, although at the price of being the less accurate one (Zhou et al., 2017). This  
701 is in line with our findings. We could show that the higher accuracy of IQ-TREE is particularly striking  
702 when using large concatenated datasets. As stated above, this is largely attributable to the much richer  
703 choice of models implemented in the former. ModelFinder (Kalyaanamoorthy et al., 2017) selected  
704 GTR+ASC+F+R6 model for the concatenated supermatrix, which is much richer in parameters than the  
705 GTR+CAT+Gamma20 model fitted by FastTree. The +ASC is an ascertainment bias correction  
706 parameter, which should be applied to alignments without constant sites (Lewis, 2001), such as the  
707 supermatrices generated by GET\_PHYLOMARKERS (see methods). The FreeRate model (+R)  
708 generalizes the +G model (fitting a discrete Gamma distribution to model among-site rate variation) by  
709 relaxing the assumption of Gamma-distributed rates (Yang, 1995). The FreeRate model typically fits  
710 data better than the +G model and is recommended for the analysis of large data sets (Soubrier et al.,  
711 2012).

712 The impact of substitution models in phylogenetics has been extensively studied (Posada and Buckley,  
713 2004). However, the better models implemented in IQ-TREE are not the only reason for its superior  
714 performance. A key aspect strongly impacting the quality of phylogenomic inference with large  
715 datasets is tree-searching. This has been largely neglected in most molecular systematic and  
716 phylogenetic studies of prokaryotes (Ochoa-Sánchez and Vinuesa, 2017; Vinuesa, 2010; Vinuesa et al.,  
717 2008). Due to the factorial increase of the number of distinct bifurcating topologies possible with every  
718 new sequence added to an alignment (Felsenstein, 2004a), searching the tree-space for large datasets is  
719 an NP-hard (non-deterministic polynomial-time) problem that necessarily requires heuristic algorithms.  
720 This implies that once an optimum is found, there is no way of telling whether it is the global one. The  
721 strategy to gain quantitative evidence about the quality of a certain tree is to compare its score in the  
722 context of other trees found in searches initiated from a pool of different seed trees. Due to the high  
723 dimensionality of the likelihood space, and the strict “hill-climbing” nature of ML tree search  
724 algorithms (Felsenstein, 2004a), they generally get stuck in local optima (Money and Whelan, 2012).  
725 The scores of the best trees found in each search can then be compared in the form of an “ $\ln L$  score  
726 profile”, as performed in our study. Available software implementations for fast ML tree searching use  
727 different branch-swapping strategies to try to escape from early encountered “local optima”. IQ-TREE  
728 implements a more efficient tree-searching strategy than FastTree, based on a combination of hill-  
729 climbing and stochastic nearest-neighbor interchange (NNI) operations, always keeping a pool of seed  
730 trees, which help to escape local optima (Nguyen et al., 2015). This was evident when the  $\ln L$  score  
731 profiles of both programs were compared. IQ-TREE found a much better scoring species tree despite  
732 the much higher number of independent searches performed with FastTree (50 vs. 1001) using its most  
733 intensive branch-swapping regime. An important finding of our study is the demonstration that the  $\ln L$

734 search profile of IQ-TREE is much shallower than that of FastTree. This suggests that the former finds  
735 trees much closer to the potential optimum than the latter. It has been shown that the highest-scoring  
736 (best) trees tend to have shorter branches, and overall tree-length, than those stuck in worse local  
737 optima (Money and Whelan, 2012). In agreement with this report, the best species-tree found by IQ-  
738 TREE has a notoriously shorter total length and significantly shorter edges than those of the best  
739 species-tree found by FastTree.

740  
741 Our extensive benchmark analysis conclusively demonstrated the superior performance of IQ-TREE.  
742 Based on this evidence, it was chosen as the default search algorithm for GET\_PHYLOMARKERS.  
743 However, it should be noted that topological differences between the best trees found by both programs  
744 were minor, not affecting the composition of the major clades in the corresponding species trees. It is  
745 therefore safe to conclude that the reclassification of *Stenotrophomonas* genome sequences proposed in  
746 Table 3 is robust. They are consistently supported by the species-trees estimated with both programs.  
747 This result underlines the utility of GET\_PHYLOMARKERS to identify misclassified genomes in  
748 public sequence repositories, a problem found in many genera (Gomila et al., 2017; Sangal et al.,  
749 2016). GET\_PHYLOMARKERS is unique in its ability to combine core-genome phylogenomics with  
750 ML and parsimony phylogeny estimation from the pan-genome matrix. In line with other recent studies  
751 (Caputo et al., 2015; Tu and Lin, 2016), we demonstrate that pan-genome analyses are valuable in the  
752 context of microbial molecular systematics and taxonomy. All genomes found to be misclassified based  
753 on the phylogenomic analysis of core-genomes were corroborated by the ML and parsimony analyses  
754 of the PGM. Furthermore, the combined evidence gained from these independent approaches  
755 consistently revealed that the Smc contains up to 20 monophyletic and strongly supported species-like  
756 clusters. These are defined at the cgANIb 95.9% threshold, and include the previously identified  
757 genospecies Smc1-Smc4 (Ochoa-Sánchez and Vinuesa, 2017), and up to 13 genospecies within the *S.*  
758 *maltophilia sensu lato* clade. This threshold fits well with the currently favored gANI > 94% cutoff for  
759 species delimitation (Konstantinidis and Tiedje, 2005; Richter and Rossello-Mora, 2009). The  
760 consistency among all the different approaches strongly supports the proposed delimitations. We used  
761 an unsupervised learning procedure to determine the optimal number of clusters ( $k$ ) in the cgANDb  
762 matrix computed from the 86 Smc genomes analyzed. The average silhouette width goodness of  
763 clustering statistic proposed an optimal  $k = 19$ , which corresponds to a gANI = 95.5%. At this cutoff,  
764 13 (instead of 14) species-like clusters are delimited within the *S. maltophilia sensu lato* clade. This  
765 unsupervised learning method therefore seems promising to define the optimal number of clusters in  
766 ANI-like matrices using a statistical procedure. However, it should be critically and extensively  
767 evaluated in other geno-taxonomic studies to better understand its properties and possible limitations,  
768 before being broadly used.

769  
770 Current models of microbial speciation predict that bacterial species-like lineages should be identifiable  
771 by significantly reduced gene flow between them, even when recombination levels are high within  
772 species (Cadillo-Quiroz et al., 2012; Shapiro et al., 2012). Such lineages should also display  
773 differentiated ecological niches and phenotypes (Koeppel et al., 2008; Shapiro and Polz, 2015). In our  
774 previous comprehensive multilocus sequence analysis of species borders in the genus  
775 *Stenotrophomonas* (Ochoa-Sánchez and Vinuesa, 2017) we could show that those models fitted our  
776 data well. We found highly significant genetic differentiation and marginal gene-flow across strains  
777 from sympatric Smc1 and Smc2 lineages, as well as highly significant differences in the resistance  
778 profiles of *S. maltophilia sensu lato* isolates versus Smc1 and Smc2 isolates. We could also show that  
779 all three lineages have different habitat preferences (Ochoa-Sánchez and Vinuesa, 2017). The genomic  
780 analyses presented in this study for five Smc1 and Smc2 strains, respectively, fully support their



781 separate species status from a geno-taxonomic perspective. Given the recognized importance of gene  
782 gain and loss processes in bacterial speciation and ecological specialization (Caputo et al., 2015;  
783 Jeukens et al., 2017; Richards et al., 2014; Shapiro and Polz, 2015), as reported also in plants (Gordon  
784 et al., 2017), we think that the evidence gained from pan-genome phylogenies is particularly  
785 informative for microbial geno-taxonomic investigations. We believe they should be used to validate  
786 the groupings obtained by the classical gANI cutoff-based species delimitation procedure (Goris et al.,  
787 2007; Konstantinidis and Tiedje, 2005; Richter and Rossello-Mora, 2009) that dominates current geno-  
788 taxonomic research. It is well documented that pan-genome-based groupings tend to better reflect  
789 ecologically relevant phenotypic differences between groups (Caputo et al., 2015; Jeukens et al., 2017;  
790 Lukjancenko et al., 2010). We recommend that future geno-taxonomic studies search for a consensus of  
791 the complementary views of genomic diversity provided by OGRIs, core- and pan-genome  
792 phylogenies, as performed herein. GET\_PHYLOMARKERS is a useful and versatile tool for this task.  
793

794 In summary, in this study we developed a comprehensive and powerful suite of open-source  
795 computational tools for state-of-the art phylogenomic and pan-genomic analyses. Their application to  
796 critically analyze the geno-taxonomic status of the genus *Stenotrophomonas* provided compelling  
797 evidence that the taxonomically ill-defined *S. maltophilia* complex holds many cryptic species.  
798 However, we refrain at this point from making formal taxonomic proposals for them because we have  
799 not yet performed the above-mentioned population genetic analyses to demonstrate the genetic  
800 cohesiveness of the individual species and their differentiation from closely related ones. This will be  
801 the topic of a follow-up work in preparation. We think that comparative genomic analyses designed to  
802 identify lineage-specific genetic differences that may underlie niche-differentiation of species are the  
803 most powerful and objective criteria to delimit species in any taxonomic group (Ochoa-Sánchez and  
804 Vinuesa, 2017; Vinuesa et al., 2005).  
805  
806

## 807 **AUTHOR CONTRIBUTIONS**

808

809 PV designed the project, wrote the bulk of the code, assembled the genomes, performed the analyses  
810 and wrote the paper. LEOS isolated the strains sequenced in this study and performed all wet-lab  
811 experiments. BCM was involved in the original design of the project, contributed code, and set up the  
812 docker image. All authors read and approved the final version of the manuscript.  
813

## 814 **FUNDING**

815

816 We gratefully acknowledge the funding provided by DGAPA/PAPIIT-UNAM (grants IN201806-2,  
817 IN211814 and IN206318) and CONACyT-México (grants P1-60071, 179133 and FC-2105-2-879) to  
818 PV, as well as the Fundación ARAID, Consejo Superior de Investigaciones Científicas (grant  
819 200720I038 and Spanish MINECO (AGL2013-48756-R) to BCM.  
820

## 821 **ACKNOWLEDGEMENTS**

822

823 We thank Javier Rivera for excellent technical support with wet-lab experiments and José Alfredo  
824 Hernández and Víctor del Moral for support with server administration. Jason Steel from the DNASU  
825 Sequencing Core at The Biodesign Institute, Arizona State University, is acknowledged for generating  
826 the genome sequences of our samples. Dr. Claudia Silva is thanked for her critical reading of the



827 manuscript. We are thankful to GitHub (<https://github.com/>), docker (<https://hub.docker.com/>) and the  
828 open-source community at large, for providing great resources for software development.  
829

830

## 831 REFERENCES

832

833

834 Adamek, M., Linke, B., and Schwartz, T. (2014). Virulence genes in clinical and environmental *Stenotrophomonas maltophilia* isolates: A genome sequencing  
835 and gene expression approach. *Microb. Pathog.* 67–68, 20–30. doi:10.1016/j.micpath.2014.02.001.

836 Adamek, M., Overhage, J., Bathe, S., Winter, J., Fischer, R., and Schwartz, T. (2011). Genotyping of Environmental and Clinical *Stenotrophomonas*  
837 *maltophilia* Isolates and their Pathogenic Potential. *PLoS One* 6, 1–11. doi:10.1371/journal.pone.0027615.

838 Alavi, P., Starcher, M. R., Thallinger, G. G., Zachow, C., Müller, H., and Berg, G. (2014). *Stenotrophomonas* comparative genomics reveals genes and  
839 functions that differentiate beneficial and pathogenic bacteria. *BMC Genomics* 15, 482. doi:10.1186/1471-2164-15-482.

840 Angiuoli, S. V., Gussman, A., Klimke, W., Cochrane, G., Field, D., Garrity, G., et al. (2008). Toward an online repository of Standard Operating Procedures  
841 (SOPs) for (meta)genomic annotation. *OMICS* 12, 137–41. doi:10.1089/omi.2008.0017.

842 Anisimova, M., and Gascuel, O. (2006). Approximate likelihood-ratio test for branches: A fast, accurate, and powerful alternative. *Syst. Biol.* 55, 352–539.

843 Bankevich, A., Nurk, S., Antipov, D., Gurevich, A. A., Dvorkin, M., Kulikov, A. S., et al. (2012). SPAdes: a new genome assembly algorithm and its  
844 applications to single-cell sequencing. *J. Comput. Biol.* 19, 455–477. doi:10.1089/cmb.2012.0021.

845 Berg, G., Fritze, A., Hagemann, M., and Wolf, A. (2002). *Stenotrophomonas rhizophila* sp. nov., a novel plant-associated bacterium with antifungal  
846 properties. *Int. J. Syst. Evol. Microbiol.* 52, 1937–1944. doi:10.1099/00207713-52-6-1937.

847 Berg, G., and Martinez, J. L. (2015). Friends or foes: Can we make a distinction between beneficial and harmful strains of the *Stenotrophomonas*  
848 *maltophilia* complex? *Front. Microbiol.* 6. doi:10.3389/fmicb.2015.00241.

849 Blow, F., Vontas, J., and Darby, A. C. (2016). Draft Genome Sequence of *Stenotrophomonas maltophilia* SBo1 Isolated from *Bactrocera oleae*. *Genome*  
850 *Announc.* 4, e00905-16. doi:10.1128/genomeA.00905-16.

851 Bodilis, J., Youenou, B., Briolay, J., Brothier, E., Favre-Bonté, S., and Nazaret, S. (2016). Draft Genome Sequences of *Stenotrophomonas maltophilia*  
852 Strains Sm32COP, Sm41DVV, Sm46PAILV, SmF3, SmF22, SmSOFb1, and SmCVFa1, Isolated from Different Manures in France. *Genome*  
853 *Announc.* 4, e00841-16. doi:10.1128/genomeA.00841-16.

854 Bosi, E., Donati, B., Galardini, M., Brunetti, S., Sagot, M.-F., Lió, P., et al. (2015). MeDuSa: a multi-draft based scaffold. *Bioinformatics* 31, 2443–2451.  
855 doi:10.1093/bioinformatics/btv171.

856 Brooke, J. S. (2012). *Stenotrophomonas maltophilia*: an emerging global opportunistic pathogen. *Clin. Microbiol. Rev.* 25, 2–41. doi:25/1/2  
857 [pii]10.1128/CMR.00019-11.

858 Bruen, T. C., Philippe, H., and Bryant, D. (2005). A Simple and Robust Statistical Test for Detecting the Presence of Recombination. *Genetics* 172, 2665–  
859 2681. doi:10.1534/genetics.105.048975.

860 Cadillo-Quiroz, H., Didelot, X., Held, N. L., Herrera, A., Darling, A., Reno, M. L., et al. (2012). Patterns of gene flow define species of thermophilic  
861 Archaea. *PLoS Biol.* 10, e1001265. doi:10.1371/journal.pbio.1001265.

862 Caputo, A., Merhej, V., Georgiades, K., Fournier, P.-E., Croce, O., Robert, C., et al. (2015). Pan-genomic analysis to redefine species and subspecies based  
863 on quantum discontinuous variation: the *Klebsiella* paradigm. *Biol. Direct* 10, 55. doi:10.1186/s13062-015-0085-2.

864 Chang, Y.-T., Lin, C.-Y., Chen, Y.-H., and Hsueh, P.-R. (2015). Update on infections caused by *Stenotrophomonas maltophilia* with particular attention to  
865 resistance mechanisms and therapeutic options. *Front. Microbiol.* 6, 893. doi:10.3389/fmicb.2015.00893.

866 Chen, C., Khaleel, S. S., Huang, H., and Wu, C. H. (2014). Software for pre-processing Illumina next-generation sequencing short read sequences. *Source*  
867 *Code Biol. Med.* 9, 8. doi:10.1186/1751-0473-9-8.

868 Chun, J., and Rainey, F. A. (2014). Integrating genomics into the taxonomy and systematics of the Bacteria and Archaea. *Int. J. Syst. Evol. Microbiol.* 64,  
869 316–324. doi:10.1099/ijs.0.054171-0.

870 Ciccarelli, F. D., Doerks, T., von Mering, C., Creevey, C. J., Snel, B., and Bork, P. (2006). Toward automatic reconstruction of a highly resolved tree of life.  
871 *Science (80-. )* 311, 1283–1287.

- 872 Conchillo-Solé, O., Yero, D., Coves, X., Huedo, P., Martínez-Servat, S., Daura, X., et al. (2015). Draft Genome Sequence of *Stenotrophomonas maltophilia*  
873 Strain UV74 Reveals Extensive Variability within Its Genomic Group. *Genome Announc.* 3, e00611-15. doi:10.1128/genomeA.00611-15.
- 874 Contreras-Moreira, B., Cantalapiedra, C. P. C. P. C. P., García-Pereira, M. J. M. J. M. J., Gordon, S. P. S. P., Vogel, J. P. J. P., Igartua, E., et al. (2017).  
875 Analysis of plant pan-genomes and transcriptomes with GET\_HOMOLOGUES-EST, a clustering solution for sequences of the same species.  
876 *Front. plant Sci.* 8:184, 184. doi:10.3389/fpls.2017.00184.
- 877 Contreras-Moreira, B., and Vinuesa, P. (2013). GET\_HOMOLOGUES, a versatile software package for scalable and robust microbial pangenome analysis.  
878 *Appl. Environ. Microbiol.* 79, 7696–7701. doi:AEM.02411-13 [pii]10.1128/AEM.02411-13.
- 879 Crossman, L. C., Gould, V. C., Dow, J. M., Vernikos, G. S., Okazaki, A., Sebahia, M., et al. (2008). The complete genome, comparative and functional  
880 analysis of *Stenotrophomonas maltophilia* reveals an organism heavily shielded by drug resistance determinants. *Genome Biol.* 9, R74. doi:gb-  
881 2008-9-4-r74 [pii]10.1186/gb-2008-9-4-r74.
- 882 Daubin, V., Gouy, M., and Perriere, G. (2002). A phylogenomic approach to bacterial phylogeny: evidence of a core of genes sharing a common history.  
883 *Genome. Res.* 12, 1080–1090.
- 884 Darling, A. E., Jospin, G., Lowe, E., Matsen, F. A., Bik, H. M., and Eisen, J. A. (2014). PhyloSift: phylogenetic analysis of genomes and metagenomes.  
885 *PeerJ* 2, e243. doi:10.7717/peerj.243.
- 886 Davenport, K. W., Daligault, H. E., Minogue, T. D., Broomall, S. M., Bruce, D. C., Chain, P. S., et al. (2014). Complete Genome Sequence of  
887 *Stenotrophomonas maltophilia* Type Strain 810-2 (ATCC 13637). *Genome Announc.* 2, e00974-14-e00974-14. doi:10.1128/genomeA.00974-14.
- 888 Degnan, J. H., and Rosenberg, N. A. (2009). Gene tree discordance, phylogenetic inference and the multispecies coalescent. *Trends Ecol. Evol.* 24, 332–  
889 340. doi:S0169-5347(09)00084-6 [pii]10.1016/j.tree.2009.01.009.
- 890 Deredjian, A., Alliot, N., Blanchard, L., Brothier, E., Anane, M., Cambier, P., et al. (2016). Occurrence of *Stenotrophomonas maltophilia* in agricultural  
891 soils and antibiotic resistance properties. *Res. Microbiol.* 167, 313–324. doi:10.1016/j.resmic.2016.01.001.
- 892 Didelot, X., and Maiden, M. C. J. (2010). Impact of recombination on bacterial evolution. *Trends Microbiol.* 18, 315–322. doi:10.1016/j.tim.2010.04.002.
- 893 Dornburg, A., Townsend, J. P., and Wang, Z. (2017). “Maximizing Power in Phylogenetics and Phylogenomics: A Perspective Illuminated by Fungal Big  
894 Data,” in *Advances in genetics*, 1–47. doi:10.1016/bs.adgen.2017.09.007.
- 895 Felsenstein, J. (1981). Evolutionary trees from DNA sequences: a maximum likelihood approach. *J. Mol. Evol.* 17, 368–76.
- 896 Felsenstein, J. (2004a). *Inferring phylogenies*. Sunderland, MA: Sinauer Associates, INC.
- 897 Felsenstein, J. (2004b). PHYLIP (Phylogeny Inference Package).
- 898 Finkmann, W., Altendorf, K., Stackebrandt, E., and Lipski, A. (2000). Characterization of N(2)O-producing *Xanthomonas*-like isolates from biofilters as  
899 *Stenotrophomonas nitritireducens* sp. nov., *Luteimonas mephitis* gen. nov., sp. nov. and *Pseudoxanthomonas broegbernensis* gen. nov., sp. nov. *Int.*  
900 *J. Syst. Evol. Microbiol.* 50, 273–282. doi:10.1099/00207713-50-1-273.
- 901 Gauthier, J., Vincent, A. T., Charette, S. J., and Derome, N. (2017). Strong Genomic and Phenotypic Heterogeneity in the *Aeromonas sobria* Species  
902 Complex. *Front. Microbiol.* 8, 2434. doi:10.3389/fmicb.2017.02434.
- 903 Ge, S., and Ge, S. C. (2016). Simultaneous Cr(VI) reduction and Zn(II) biosorption by *Stenotrophomonas* sp. and constitutive expression of related genes.  
904 *Biotechnol. Lett.* 38, 877–84. doi:10.1007/s10529-016-2057-8.
- 905 Gomila, M., Busquets, A., Mulet, M., García-Valdés, E., and Lalucat, J. (2017). Clarification of Taxonomic Status within the *Pseudomonas syringae*  
906 Species Group Based on a Phylogenomic Analysis. *Front. Microbiol.* 8, 2422. doi:10.3389/fmicb.2017.02422.
- 907 Gordon, S. P., Contreras-Moreira, B., Woods, D. P., Des Marais, D. L., Burgess, D., Shu, S., et al. (2017). Extensive gene content variation in the  
908 *Brachypodium distachyon* pan-genome correlates with population structure. *Nat. Commun.* 8, 2184. doi:10.1038/s41467-017-02292-8.
- 909 Goris, J., Konstantinidis, K. T., Klappenbach, J. a., Coenye, T., Vandamme, P., and Tiedje, J. M. (2007). DNA-DNA hybridization values and their  
910 relationship to whole-genome sequence similarities. *Int. J. Syst. Evol. Microbiol.* 57, 81–91. doi:10.1099/ijls.0.64483-0.
- 911 Guindon, S., Dufayard, J.-F. F., Lefort, V., Anisimova, M., Hordijk, W., and Gascuel, O. (2010). New algorithms and methods to estimate maximum-  
912 likelihood phylogenies: assessing the performance of PhyML 3.0. *Syst. Biol.* 59, 307–321. doi:syq010 [pii]10.1093/sysbio/syq010.
- 913 Handa, Y., Tazato, N., Nagatsuka, Y., Koide, T., Kigawa, R., Sano, C., et al. (2016). *Stenotrophomonas tumulicola* sp. nov., a major contaminant of the  
914 stone chamber interior in the Takamatsuzuka Tumulus. *Int. J. Syst. Evol. Microbiol.* 66, 1119–1124. doi:10.1099/ijsem.0.000843.
- 915 Heylen, K., Vanparys, B., Peirsegaale, F., Lebbe, L., and De Vos, P. (2007). *Stenotrophomonas terrae* sp. nov. and *Stenotrophomonas humi* sp. nov., two  
916 nitrate-reducing bacteria isolated from soil. *Int. J. Syst. Evol. Microbiol.* 57, 2056–2061. doi:10.1099/ijls.0.65044-0.

- 917 Hoang, D. T., Chernomor, O., von Haeseler, A., Minh, B. Q., Le, S. V., and Vinh, L. S. (2017). UFBoot2: Improving the Ultrafast Bootstrap  
918 Approximation. *Mol. Biol. Evol.* doi:10.1093/molbev/msx281.
- 919 Hughes, G. L., Raygoza Garay, J. A., Koundal, V., Rasgon, J. L., and Mwangi, M. M. (2016). Genome Sequence of *Stenotrophomonas maltophilia* Strain  
920 SmAs1, Isolated From the Asian Malaria Mosquito *Anopheles stephensi*. *Genome Announc.* 4, e00086-16. doi:10.1128/genomeA.00086-16.
- 921 Jeukens, J., Freschi, L., Vincent, A. T., Emond-Rheault, J.-G., Kukavica-Ibrulj, I., Charette, S. J., et al. (2017). A Pan-Genomic Approach to Understand the  
922 Basis of Host Adaptation in *Achromobacter*. *Genome Biol. Evol.* 9, 1030–1046. doi:10.1093/gbe/evx061.
- 923 Junier, T., and Zdobnov, E. M. (2010). The Newick utilities: high-throughput phylogenetic tree processing in the UNIX shell. *Bioinformatics* 26, 1669–  
924 1670. doi:10.1093/bioinformatics/btq243.
- 925 Kaiser, S., Biehler, K., and Jonas, D. (2009). A *Stenotrophomonas maltophilia* multilocus sequence typing scheme for inferring population structure. *J.*  
926 *Bacteriol.* 191, 2934–2943. doi:JB.00892-08 [pii]10.1128/JB.00892-08.
- 927 Kalyaanamoorthy, S., Minh, B. Q., Wong, T. K. F., von Haeseler, A., and Jermini, L. S. (2017). ModelFinder: fast model selection for accurate  
928 phylogenetic estimates. *Nat. Methods* 14, 587–589. doi:10.1038/nmeth.4285.
- 929 Kämpfer, P. (2012). Systematics of prokaryotes: The state of the art. *Antonie van Leeuwenhoek, Int. J. Gen. Mol. Microbiol.* 101, 3–11.  
930 doi:10.1007/s10482-011-9660-4.
- 931 Kaparullina, E., Doronina, N., Chistyakova, T., and Trotsenko, Y. (2009). *Stenotrophomonas chelatiphaga* sp. nov., a new aerobic EDTA-degrading  
932 bacterium. *Syst. Appl. Microbiol.* 32, 157–162. doi:10.1016/j.syapm.2008.12.003.
- 933 Kaufman, L., and Rousseeuw, P. J. eds. (1990). *Finding Groups in Data*. Hoboken, NJ, USA: John Wiley & Sons, Inc. doi:10.1002/9780470316801.
- 934 Kim, H.-B., Srinivasan, S., Sathiyaraj, G., Quan, L.-H., Kim, S.-H., Bui, T. P. N., et al. (2010). *Stenotrophomonas ginsengisoli* sp. nov., isolated from a  
935 ginseng field. *Int. J. Syst. Evol. Microbiol.* 60, 1522–1526. doi:10.1099/ijs.0.014662-0.
- 936 Koeppl, A., Perry, E. B., Sikorski, J., Krizanc, D., Warner, A., Ward, D. M., et al. (2008). Identifying the fundamental units of bacterial diversity: A  
937 paradigm shift to incorporate ecology into bacterial systematics. *Proc. Natl. Acad. Sci.* 105, 2504–2509. doi:10.1073/pnas.0712205105.
- 938 Konstantinidis, K. T., and Tiedje, J. M. (2005). Genomic insights that advance the species definition for prokaryotes. *Proc. Natl. Acad. Sci. U. S. A.* 102,  
939 2567–2572.
- 940 Konstantinidis, K. T., and Tiedje, J. M. (2007). Prokaryotic taxonomy and phylogeny in the genomic era: advancements and challenges ahead. *Curr. Opin.*  
941 *Microbiol.* 10, 504–509.
- 942 Kubatko, L. S., and Degnan, J. H. (2007). Inconsistency of phylogenetic estimates from concatenated data under coalescence. *Syst. Biol.* 56, 17–24.
- 943 Kumar, S., Filipski, A. J., Battistuzzi, F. U., Kosakovsky Pond, S. L., and Tamura, K. (2012). Statistics and Truth in Phylogenomics. *Mol. Biol. Evol.* 29,  
944 457–472. doi:10.1093/molbev/msr202.
- 945 Labat, M., Thierry, S., Macarie, H., Cayol, J.-L., Ouattara, A. S., and Assih, E. A. (2002). *Stenotrophomonas acidaminiphila* sp. nov., a strictly aerobic  
946 bacterium isolated from an upflow anaerobic sludge blanket (UASB) reactor. *Int. J. Syst. Evol. Microbiol.* 52, 559–568. doi:10.1099/00207713-52-  
947 2-559.
- 948 Le, Q., Sievers, F., and Higgins, D. G. (2017). Protein Multiple Sequence Alignment Benchmarking through Secondary Structure Prediction.  
949 *Bioinformatics* 33, btw840. doi:10.1093/bioinformatics/btw840.
- 950 Lee, M., Woo, S.-G., Chae, M., Shin, M.-C., Jung, H.-M., and Ten, L. N. (2011). *Stenotrophomonas daejeonensis* sp. nov., isolated from sewage. *Int. J.*  
951 *Syst. Evol. Microbiol.* 61, 598–604. doi:10.1099/ijs.0.017780-0.
- 952 Lerat, E., Daubin, V., and Moran, N. A. (2003). From gene trees to organismal phylogeny in prokaryotes: the case of the gamma-Proteobacteria. *PLoS Biol.*  
953 1, E19.
- 954 Lewis, P. O. (2001). A likelihood approach to estimating phylogeny from discrete morphological character data. *Syst. Biol.* 50, 913–25.
- 955 Li, H., and Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* 25, 1754–1760.  
956 doi:10.1093/bioinformatics/btp324.
- 957 Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., et al. (2009). The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25,  
958 2078–2079. doi:10.1093/bioinformatics/btp352.
- 959 Lira, F., Berg, G., and Martínez, J. L. (2017). Double-Face Meets the Bacterial World: The Opportunistic Pathogen *Stenotrophomonas maltophilia*. *Front.*  
960 *Microbiol.* 8, 2190. doi:10.3389/fmicb.2017.02190.

- 961 Lira, F., Hernández, A., Belda, E., Sánchez, M. B., Moya, A., Silva, F. J., et al. (2012). Whole-genome sequence of *Stenotrophomonas maltophilia* D457, A  
962 clinical isolate and a model strain. *J. Bacteriol.* 194, 3563–3564. doi:10.1128/JB.00602-12.
- 963 Lukjancenko, O., Wassenaar, T. M., and Ussery, D. W. (2010). Comparison of 61 sequenced *Escherichia coli* genomes. *Microb. Ecol.* 60, 708–720.  
964 doi:10.1007/s00248-010-9717-3.
- 965 Martin, D. P. (2009). Recombination detection and analysis using RDP3. *Methods Mol. Biol.* 537, 185–205. doi:10.1007/978-1-59745-251-9\_9.
- 966 Money, D., and Whelan, S. (2012). Characterizing the Phylogenetic Tree-Search Problem. *Syst. Biol.* 61, 228. doi:10.1093/sysbio/syr097.
- 967 Naas, T., Cuzon, G., Villegas, M. V., Lartigue, M. F., Quinn, J. P., and Nordmann, P. (2008). Genetic structures at the origin of acquisition of the  $\beta$ -  
968 lactamase blaKPC gene. *Antimicrob. Agents Chemother.* 52, 1257–1263. doi:10.1128/AAC.01451-07.
- 969 Nguyen, L.-T., Schmidt, H. A., von Haeseler, A., and Minh, B. Q. (2015). IQ-TREE: A Fast and Effective Stochastic Algorithm for Estimating Maximum-  
970 Likelihood Phylogenies. *Mol. Biol. Evol.* 32, 268–274. doi:10.1093/molbev/msu300.
- 971 Niu, B., Paulson, J. N., Zheng, X., and Kolter, R. (2017). Simplified and representative bacterial community of maize roots. *Proc. Natl. Acad. Sci. U. S. A.*  
972 114, E2450–E2459. doi:10.1073/pnas.1616148114.
- 973 Nourdin-Galindo, G., Sánchez, P., Molina, C. F., Espinoza-Rojas, D. A., Oliver, C., Ruiz, P., et al. (2017). Comparative Pan-Genome Analysis of  
974 *Piscirickettsia salmonis* Reveals Genomic Divergences within Genogroups. *Front. Cell. Infect. Microbiol.* 7, 459. doi:10.3389/fcimb.2017.00459.
- 975 Ochoa-Sánchez, L. E. L. E., and Vinuesa, P. (2017). Evolutionary genetic analysis uncovers multiple species with distinct habitat preferences and antibiotic  
976 resistance phenotypes in the *Stenotrophomonas maltophilia* complex. *Front. Microbiol.* 8, 1548. doi:10.3389/fmicb.2017.01548.
- 977 Ormerod, K. L., George, N. M., Fraser, J. A., Wainwright, C., and Hugenholtz, P. (2015). Comparative genomics of non-pseudomonal bacterial species  
978 colonising paediatric cystic fibrosis patients. *PeerJ* 3, e1223. doi:10.7717/peerj.1223.
- 979 Ouattara, A. S., Le Mer, J., Joseph, M., and Macarie, H. (2017). Transfer of *Pseudomonas pictorum* Gray and Thornton 1928 to genus *Stenotrophomonas*  
980 as *Stenotrophomonas pictorum* comb. nov., and emended description of the genus *Stenotrophomonas*. *Int. J. Syst. Evol. Microbiol.* 67, 1894–1900.  
981 doi:10.1099/ijsem.0.001880.
- 982 Pak, T. R., Altman, D. R., Attie, O., Sebra, R., Hamula, C. L., Lewis, M., et al. (2015). Whole-genome sequencing identifies emergence of a quinolone  
983 resistance mutation in a case of *Stenotrophomonas maltophilia* bacteremia. *Antimicrob. Agents Chemother.* 59, 7117–20. doi:10.1128/AAC.01723-  
984 15.
- 985 Palleroni, N. J. (2005). “Genus IX. *Stenotrophomonas* Palleroni and Bradbury 1993.,” in *Bergey’s Manual of Systematic Bacteriology 2nd Edition*, eds. G.  
986 M. Garrity, D. J. Brenner, N. R. Krieg, and J. T. Staley (New York: Springer), 107–115.
- 987 Palleroni, N. J., and Bradbury, J. F. (1993). *Stenotrophomonas*, a new bacterial genus for *Xanthomonas maltophilia* (Hugh 1980) Swings et al. 1983. *Int. J.*  
988 *Syst. Bacteriol.* 43, 606–9. doi:10.1099/00207713-43-3-606.
- 989 Pan, X., Lin, D., Zheng, Y., Zhang, Q., Yin, Y., Cai, L., et al. (2016). Biodegradation of DDT by *Stenotrophomonas* sp. DDT-1: Characterization and  
990 genome functional analysis. *Sci. Rep.* 6, 21332. doi:10.1038/srep21332.
- 991 Parks, M. B., Wickett, N. J., and Alverson, A. J. (2018). Signal, Uncertainty, and Conflict in Phylogenomic Data for a Diverse Lineage of Microbial  
992 Eukaryotes (Diatoms, Bacillariophyta). *Mol. Biol. Evol.* 35, 80–93.
- 993 Patil, P. P., Midha, S., Kumar, S., and Patil, P. B. (2016). Genome Sequence of Type Strains of Genus *Stenotrophomonas*. *Front. Microbiol.* 7, 309.  
994 doi:10.3389/fmicb.2016.00309.
- 995 Pease, J. B., and Hahn, M. W. (2013). MORE ACCURATE PHYLOGENIES INFERRED FROM LOW-RECOMBINATION REGIONS IN THE  
996 PRESENCE OF INCOMPLETE LINEAGE SORTING. *Evolution (N. Y.)* 67, 2376–2384. doi:10.1111/evo.12118.
- 997 Popescu, A.-A., Huber, K. T., and Paradis, E. (2012). ape 3.0: New tools for distance-based phylogenetics and evolutionary analysis in R. *Bioinformatics*  
998 28, 1536–1537. doi:10.1093/bioinformatics/bts184.
- 999 Posada, D., and Buckley, T. R. (2004). Model selection and model averaging in phylogenetics: advantages of Akaike information criterion and bayesian  
1000 approaches over likelihood ratio tests. *Syst. Biol.* 53, 793–808.
- 1001 Posada, D., and Crandall, K. A. (2002). The effect of recombination on the accuracy of phylogeny estimation. *J. Mol. Evol.* 54, 396–402.
- 1002 Price, M. N., Dehal, P. S., and Arkin, A. P. (2010). FastTree 2--approximately maximum-likelihood trees for large alignments. *PLoS One* 5, e9490.  
1003 doi:10.1371/journal.pone.0009490.

- 1004 Ramos, P. L., Van Trappen, S., Thompson, F. L., Rocha, R. C. S., Barbosa, H. R., De Vos, P., et al. (2011). Screening for endophytic nitrogen-fixing  
1005 bacteria in Brazilian sugar cane varieties used in organic farming and description of *Stenotrophomonas pavanii* sp. nov. *Int. J. Syst. Evol. Microbiol.*  
1006 61, 926–931. doi:10.1099/ijs.0.019372-0.
- 1007 Richards, V. P., Palmer, S. R., Pavinski Bitar, P. D., Qin, X., Weinstock, G. M., Highlander, S. K., et al. (2014). Phylogenomics and the Dynamic Genome  
1008 Evolution of the Genus *Streptococcus*. *Genome Biol. Evol.* 6, 741–753. doi:10.1093/gbe/evu048.
- 1009 Richter, M., and Rossello-Mora, R. (2009). Shifting the genomic gold standard for the prokaryotic species definition. *Proc. Natl. Acad. Sci.* 106, 19126–  
1010 19131. doi:10.1073/pnas.0906412106.
- 1011 Roach, D. J., Burton, J. N., Lee, C., Stackhouse, B., Butler-Wu, S. M., Cookson, B. T., et al. (2015). A Year of Infection in the Intensive Care Unit:  
1012 Prospective Whole Genome Sequencing of Bacterial Clinical Isolates Reveals Cryptic Transmissions and Novel Microbiota. *PLoS Genet.* 11,  
1013 e1005413. doi:10.1371/journal.pgen.1005413.
- 1014 Rosselló-Móra, R., and Amann, R. (2015). Past and future species definitions for Bacteria and Archaea. *Syst. Appl. Microbiol.* 38, 209–216.  
1015 doi:10.1016/j.syapm.2015.02.001.
- 1016 Rousseeuw, P. J. (1987). Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *J. Comput. Appl. Math.* 20, 53–65.  
1017 doi:10.1016/0377-0427(87)90125-7.
- 1018 Ryan, R. P., Monchy, S., Cardinale, M., Taghavi, S., Crossman, L., Avison, M. B., et al. (2009). The versatility and adaptation of bacteria from the genus  
1019 *Stenotrophomonas*. *Nat. Rev. Microbiol.* 7, 514–525. doi:nrmicro2163 [pii]10.1038/nrmicro2163.
- 1020 Salichos, L., and Rokas, A. (2013). Inferring ancient divergences requires genes with strong phylogenetic signals. *Nature* 497, 327–331.  
1021 doi:10.1038/nature12130.
- 1022 Sánchez-Castro, I., Ruiz-Fresneda, M. A., Bakkali, M., Kämpfer, P., Glaeser, S. P., Busse, H. J., et al. (2017). *Stenotrophomonas bentonitica* sp. nov.,  
1023 isolated from bentonite formations. *Int. J. Syst. Evol. Microbiol.* 67, 2779–2786. doi:10.1099/ijsem.0.002016.
- 1024 Sandner-Miranda, L., Vinuesa, P., Cravioto A., and Morales-Espinosa, R. (2018). The genomic basis of intrinsic and acquired antibiotic resistance in the  
1025 genus *Serratia*. *Front. Microbiol.* (in review).
- 1026 Sangal, V., Goodfellow, M., Jones, A. L., Schwalbe, E. C., Blom, J., Hoskisson, P. A., et al. (2016). Next-generation systematics: An innovative approach to  
1027 resolve the structure of complex prokaryotic taxa. *Sci. Rep.* 6, 38392. doi:10.1038/srep38392.
- 1028 Sasser, D., Leardini, I., Villa, L., Comandatore, F., Carta, C., Almeida, A., et al. (2013). Draft Genome Sequence of *Stenotrophomonas maltophilia* Strain  
1029 EPM1, Found in Association with a Culture of the Human Parasite *Giardia duodenalis*. *Genome Announc.* 1, e00182-13-e00182-13.  
1030 doi:10.1128/genomeA.00182-13.
- 1031 Savory, E. A., Fuller, S. L., Weisberg, A. J., Thomas, W. J., Gordon, M. I., Stevens, D. M., et al. (2017). Evolutionary transitions between beneficial and  
1032 phytopathogenic *Rhodococcus* challenge disease management. *Elife* 6. doi:10.7554/eLife.30925.
- 1033 Schierup, M. H., and Hein, J. (2000). Consequences of recombination on traditional phylogenetic analysis. *Genetics* 156, 879–891. Available at:  
1034 [http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Citation&list\\_uids=11014833](http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Citation&list_uids=11014833).
- 1035 Shapiro, B. J., Friedman, J., Cordero, O. X., Preheim, S. P., Timberlake, S. C., Szabo, G., et al. (2012). Population genomics of early events in the  
1036 ecological differentiation of bacteria. *Science (80-. )*. 336, 48–51. doi:336/6077/48 [pii]10.1126/science.1218198.
- 1037 Shapiro, B. J., and Polz, M. F. (2015). Microbial Speciation. *Cold Spring Harb. Perspect. Biol.* 7, a018143. doi:10.1101/cshperspect.a018143.
- 1038 Shen, X.-X., Hittinger, C. T., and Rokas, A. (2017). Contentious relationships in phylogenomic studies can be driven by a handful of genes. *Nat. Ecol.*  
1039 *Evol.* 1, 126. doi:10.1038/s41559-017-0126.
- 1040 Sievers, F., and Higgins, D. G. (2018). Clustal Omega for making accurate alignments of many protein sequences. *Protein Sci.* 27, 135–145.  
1041 doi:10.1002/pro.3290.
- 1042 Sievers, F., Wilm, A., Dineen, D., Gibson, T. J., Karplus, K., Li, W., et al. (2012). Fast, scalable generation of high-quality protein multiple sequence  
1043 alignments using Clustal Omega. *Mol. Syst. Biol.* 7, 539. doi:msb201175 [pii]10.1038/msb.2011.75.
- 1044 Soubrier, J., Steel, M., Lee, M. S. Y., Der Sarkissian, C., Guindon, S., Ho, S. Y. W., et al. (2012). The influence of rate heterogeneity among sites on the  
1045 time dependence of molecular rates. *Mol. Biol. Evol.* 29, 3345–58. doi:10.1093/molbev/mss140.
- 1046 Stackebrandt, E., Frederiksen, W., Garrity, G. M., Grimont, P. A., Kämpfer, P., Maiden, M. C., et al. (2002). Report of the ad hoc committee for the re-  
1047 evaluation of the species definition in bacteriology. *Int. J. Syst. Evol. Microbiol.* 52, 1043–1047. Available at:  
1048 [http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Citation&list\\_uids=12054223](http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Citation&list_uids=12054223).



- 1049 Stackebrandt, E., and Goebel, B. M. (1994). Taxonomic note: A place for DNA-DNA reassociation and 16S rRNA sequence analysis in the present species  
1050 definition in bacteriology. *Int. J. Syst. Bacteriol.* 44, 846–849.
- 1051 Suyama, M., Torrents, D., and Bork, P. (2006). PAL2NAL: robust conversion of protein sequence alignments into the corresponding codon alignments.  
1052 *Nucleic Acids Res.* 34, W609–12. doi:34/suppl\_2/W609 [pii]10.1093/nar/gkl315.
- 1053 Svensson-Stadler, L. a., Mihaylova, S. a., and Moore, E. R. B. (2012). *Stenotrophomonas* interspecies differentiation and identification by *gyrB* sequence  
1054 analysis. *FEMS Microbiol. Lett.* 327, 15–24. doi:10.1111/j.1574-6968.2011.02452.x.
- 1055 Taghavi, S., Garafola, C., Monchy, S., Newman, L., Hoffman, A., Weyens, N., et al. (2009). Genome Survey and Characterization of Endophytic Bacteria  
1056 Exhibiting a Beneficial Effect on Growth and Development of Poplar Trees. *Appl. Environ. Microbiol.* 75, 748–757. doi:10.1128/AEM.02239-08.
- 1057 Tange, O. (2011). GNU Parallel: The Command-Line Power Tool. *USENIX Mag.* 36, 42–47. doi:http://dx.doi.org/10.5281/zenodo.16303}.
- 1058 Tettelin, H., Massignani, V., Cieslewicz, M. J., Donati, C., Medini, D., Ward, N. L., et al. (2005). Genome analysis of multiple pathogenic isolates of  
1059 *Streptococcus agalactiae*: implications for the microbial “pan-genome.” *Proc. Natl. Acad. Sci. U. S. A.* 102, 13950–13955. doi:0506758102  
1060 [pii]10.1073/pnas.0506758102.
- 1061 Timme, R. E., Pettengill, J. B., Allard, M. W., Strain, E., Barrangou, R., Wehnes, C., et al. (2013). Phylogenetic diversity of the enteric pathogen  
1062 *Salmonella enterica* subsp. *enterica* inferred from genome-wide reference-free SNP characters. *Genome Biol. Evol.* 5, 2109–2123.  
1063 doi:10.1093/gbe/evt159.
- 1064 Thompson, C. C., Amaral, G. R., Campeão, M., Edwards, R. A., Polz, M. F., Dutilh, B. E., et al. (2015). Microbial taxonomy in the post-genomic era:  
1065 rebuilding from scratch? *Arch. Microbiol.* 197, 359–70. doi:10.1007/s00203-014-1071-2.
- 1066 Thompson, C. C., Chimetto, L., Edwards, R. a, Swings, J., Stackebrandt, E., and Thompson, F. L. (2013). Microbial genomic taxonomy. *BMC Genomics*  
1067 14, 913. doi:10.1186/1471-2164-14-913.
- 1068 Thompson, C. C., Vicente, A. C., Souza, R. C., Vasconcelos, A. T., Vesth, T., Alves Jr., N., et al. (2009). Genomic taxonomy of *Vibrios*. *BMC Evol Biol* 9,  
1069 258. doi:1471-2148-9-258 [pii]10.1186/1471-2148-9-258.
- 1070 Tu, Q., and Lin, L. (2016). Gene content dissimilarity for subclassification of highly similar microbial strains. *BMC Genomics* 17, 647.  
1071 doi:10.1186/s12864-016-2991-9.
- 1072 Turrientes, M.-C., González-Alba, J.-M., del Campo, R., Baquero, M.-R., Cantón, R., Baquero, F., et al. (2014). Recombination Blurs Phylogenetic Groups  
1073 Routine Assignment in *Escherichia coli*: Setting the Record Straight. *PLoS One* 9, e105395. doi:10.1371/journal.pone.0105395.
- 1074 Valdezate, S., Vindel, A., Martín-Dávila, P., Del Saz, B. S., Baquero, F., and Cantón, R. (2004). High genetic diversity among *Stenotrophomonas*  
1075 *maltophilia* strains despite their originating at a single hospital. *J. Clin. Microbiol.* 42, 693–9.
- 1076 Vandamme, P., and Peeters, C. (2014). Time to revisit polyphasic taxonomy. *Antonie Van Leeuwenhoek* 106, 57–65. doi:10.1007/s10482-014-0148-x.
- 1077 Vandamme, P., Pot, B., Gillis, M., de Vos, P., Kersters, K., and Swings, J. (1996). Polyphasic taxonomy, a consensus approach to bacterial systematics.  
1078 *Microbiol. Rev.* 60, 407–438.
- 1079 Vasileuskaya-Schulz, Z., Kaiser, S., Maier, T., Kostrzewa, M., and Jonas, D. (2011). Delineation of *Stenotrophomonas* spp. by multi-locus sequence  
1080 analysis and MALDI-TOF mass spectrometry. *Syst. Appl. Microbiol.* 34, 35–39. doi:10.1016/j.syapm.2010.11.011.
- 1081 Vinuesa, P. (2010). “Multilocus Sequence Analysis and Bacterial Species Phylogeny Estimation,” in *Molecular Phylogeny of Microorganisms*, eds. A. Oren  
1082 and R. T. Papke (Caister Academic Press), 41–64. Available at: <http://www.horizonpress.com/phylogeny>;
- 1083 Vinuesa, P., and Contreras-Moreira, B. (2015). Robust Identification of Orthologues and Paralogues for Microbial Pan-Genomics Using  
1084 GET\_HOMOLOGUES: A Case Study of pInCA/C Plasmids. *Methods Mol. Biol.* 1231, 203–232. doi:10.1007/978-1-4939-1720-4\_14.
- 1085 Vinuesa, P., and Ochoa-Sánchez, L. E. L. E. (2015). Complete Genome Sequencing of *Stenotrophomonas acidaminiphila* ZAC14D2\_NAIMI4\_2, a  
1086 Multidrug-Resistant Strain Isolated from Sediments of a Polluted River in Mexico, Uncovers New Antibiotic Resistance Genes and a Novel Class-  
1087 II Lasso Peptide Biosynthesis Ge. *Genome Announc.* 10, e01433-15. doi:10.1128/genomeA.01433-15.
- 1088 Vinuesa, P., Rojas-Jimenez, K., Contreras-Moreira, B., Mahna, S. K., Prasad, B. N., Moe, H., et al. (2008). Multilocus sequence analysis for assessment of  
1089 the biogeography and evolutionary genetics of four *Bradyrhizobium* species that nodulate soybeans on the Asiatic continent. *Appl. Environ.*  
1090 *Microbiol.* 74, 6987–6996.
- 1091 Vinuesa, P., Silva, C., Werner, D., and Martínez-Romero, E. (2005). Population genetics and phylogenetic inference in bacterial molecular systematics: the  
1092 roles of migration and recombination in *Bradyrhizobium* species cohesion and delineation. *Mol. Phylogenet. Evol.* 34, 29–54.

- 1093 Walker, B. J., Abeel, T., Shea, T., Priest, M., Abouelliel, A., Sakthikumar, S., et al. (2014). Pilon: an integrated tool for comprehensive microbial variant  
1094 detection and genome assembly improvement. *PLoS One* 9, e112963. doi:10.1371/journal.pone.0112963.
- 1095 Weyenberg, G., Huggins, P. M., Schardl, C. L., Howe, D. K., and Yoshida, R. (2014). kdtrees: Non-parametric estimation of phylogenetic tree  
1096 distributions. *Bioinformatics* 30, 2280–7. doi:10.1093/bioinformatics/btu258.
- 1097 Weyenberg, G., Yoshida, R., and Howe, D. (2017). Normalizing Kernels in the Billera-Holmes-Vogtmann Treespace. *IEEE/ACM Trans. Comput. Biol.*  
1098 *Bioinforma.* 14, 1359–1365. doi:10.1109/TCBB.2016.2565475.
- 1099 Wu, M., and Eisen, J. A. (2008). A simple, fast, and accurate method of phylogenomic inference. *Genome Biol.* 9, R151.
- 1100 Wu, Y., Wang, Y., Li, J., Hu, J., Chen, K., Wei, Y., et al. (2015). Draft Genome Sequence of *Stenotrophomonas maltophilia* Strain B418, a Promising Agent  
1101 for Biocontrol of Plant Pathogens and Root-Knot Nematode. *Genome Announc.* 3, e00015-15. doi:10.1128/genomeA.00015-15.
- 1102 Yang, H.-C., Im, W.-T., Kang, M. S., Shin, D.-Y., and Lee, S.-T. (2006). *Stenotrophomonas koreensis* sp. nov., isolated from compost in South Korea. *Int.*  
1103 *J. Syst. Evol. Microbiol.* 56, 81–84. doi:10.1099/ijs.0.63826-0.
- 1104 Yang, Z. (1995). A space-time process model for the evolution of DNA sequences. *Genetics* 139, 993–1005.
- 1105 Yi, H., Srinivasan, S., and Kim, M. K. (2010). *Stenotrophomonas panacihumi* sp. nov., isolated from soil of a ginseng field. *J. Microbiol.* 48, 30–35.  
1106 doi:10.1007/s12275-010-0006-0.
- 1107 Youenou, B., Favre-Bonté, S., Bodilis, J., Brothier, E., Dubost, A., Muller, D., et al. (2015). Comparative Genomics of Environmental and Clinical  
1108 *Stenotrophomonas maltophilia* Strains with Different Antibiotic Resistance Profiles. *Genome Biol. Evol.* 7, 2484–2505. doi:10.1093/gbe/evv161.
- 1109 Yu, D., Yin, Z., Li, B., Jin, Y., Ren, H., Zhou, J., et al. (2016). Gene flow, recombination, and positive selection in *Stenotrophomonas maltophilia*:  
1110 mechanisms underlying the diversity of the widespread opportunistic pathogen. *Genome* 59, 1063–1075. doi:10.1139/gen-2016-0073.
- 1111 Zhang, L., Morrison, M., O Cuiv, P., Evans, P., and Rickard, C. M. (2013). Genome Sequence of *Stenotrophomonas maltophilia* Strain AU12-09, Isolated  
1112 from an Intravascular Catheter. *Genome Announc.* 1, e00195-13-e00195-13. doi:10.1128/genomeA.00195-13.
- 1113 Zhou, X., Shen, X.-X., Hittinger, C. T., and Rokas, A. (2017). Evaluating Fast Maximum Likelihood-Based Phylogenetic Programs Using Empirical  
1114 Phylogenomic Data Sets. *Mol. Biol. Evol.* doi:10.1093/molbev/msx302.
- 1115 Zhu, B., Liu, H., Tian, W.-X., Fan, X.-Y., Li, B., Zhou, X.-P., et al. (2012). Genome sequence of *Stenotrophomonas maltophilia* RR-10, isolated as an  
1116 endophyte from rice root. *J. Bacteriol.* 194, 1280–1. doi:10.1128/JB.06702-11.
- 1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133  
1134  
1135  
1136  
1137  
1138

1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148

**Table 1.** Overview of key annotation features for the 10 new genome assemblies reported in this study for environmental isolates recovered from Mexican rivers and classified as genospecies 1 (Smc1) and 2 (Smc2) in the study of Ochoa-Sánchez and Vinuesa (2017). Details of their isolation sites and antimicrobial resistance phenotypes can be found therein. All genomes consist of a single gapped chromosome. Supplementary table S1 provides additional information of the assemblies. Their phylogenetic placement within the *Stenotrophomonas maltophilia* complex is shown in Figure 5 (clades Sgn1/Smc1 and Sgn2/Smc2).

Genome	Size_nt	CDSs (coding)	rRNAs	tRNAs	pseudo-genes	RefSeq Acc. num.
<i>Stenotrophomonas</i> genospecies 1 (Smc1; Sgn1) ESTM1D MKCIP4 1	4,475,880	3,904	6	59	67	CP026004
<i>Stenotrophomonas</i> genospecies 1 (Smc1; Sgn1) SAU14A NAIMI4 5	4,570,883	4,020	6	69	66	CP026003
<i>Stenotrophomonas</i> genospecies 1 (Smc1; Sgn1) ZAC14A NAIMI4 1	4,698,328	4,150	7	45	66	CP026002
<i>Stenotrophomonas</i> genospecies 1 (Smc1; Sgn1) ZAC14D1 NAIMI4 1	4,702,461	4,131	6	42	66	CP026001
<i>Stenotrophomonas</i> genospecies 1 (Smc1; Sgn1) ZAC14D1 NAIMI4 6	4,700,343	4,128	6	45	63	CP026000
<i>Stenotrophomonas</i> genospecies 2 (Smc2; Sgn2) SAU14A NAIMI4 8	4,479,100	3,893	5	54	69	CP025999
<i>Stenotrophomonas</i> genospecies 2 (Smc2; Sgn2) YAU14A MKIMI4 1	4,487,007	3,918	7	43	67	CP025998
<i>Stenotrophomonas</i> genospecies 2 (Smc2; Sgn2) YAU14D1 LEIMI4 1	4,319,112	3,819	6	51	66	CP025997
<i>Stenotrophomonas</i> genospecies 2 (Smc2; Sgn2) ZAC14D2 NAIMI4 6	4,431,104	3,882	6	52	66	CP025996
<i>Stenotrophomonas</i> genospecies 2 (Smc2; Sgn2) ZAC14D2 NAIMI4 7	4,468,731	3,918	6	66	62	CP025995

1149  
1150  
1151

1152  
1153  
1154  
1155  
1156  
1157  
1158

**Table 2.** Comparative benchmark analysis of the filtering performance of the GET\_PHYLOMARKERS pipeline when run using the FastTree (FT) and IQ-TREE (IQT) maximum-likelihood algorithms under default and high search-intensity levels. The analyses were started with the stringently defined set of 239 consensus core-genome clusters computed by GET\_HOMOLOGUES for a dataset of 119 genomes (112 *Stenotrophomonas* spp. and 7 *Xanthomonas* spp.).

Test	FTdef	FThigh	IQTdef	IQThigh
Alignments passing the Phi recombination test	127/239 (53.14 %)	125/239 (52.30 %)	125/239 (52.30 %)	127/239 (53.14 %)
Outlier phylogenies (kdetrees test; $k = 1.0$ ) out of the indicated number of non-recombinant alignments	22/127 (17.32 %) passing: 105	18 (14.17 %) passing: 107	19 (14.96 %) passing: 106	22 (17.32%) passing: 105
Alignments passing the phylogenetic signal (mean SH- <i>alrt</i> bipartition support; $m \geq 0.7$ ) test	98/105 (93.33 %)	99/107 (92.52 %)	52/106 (49.05 %)	55/105 (52.38 %)
Concatenated top-scoring markers, <i>lnL</i> score, substitution model and number of independent searches	98 markers var. sites = 36082 <i>lnL</i> = -917444.522 GTR+G searches = 1	99 markers var. sites = 35509 <i>lnL</i> = -899898.614 GTR+G searches = 1	52 markers var. sites = 25383 <i>lnL</i> = -666437.563 GTR+F+ASC+R6 searches = 1	55 markers var. sites = 26988 <i>lnL</i> = -707933.476 GTR+F+ASC+R6 searches = 5
Total wall-clock time of runs on 50 cores	0h:13m:39s	0h:38m:30s	1h:22m:18s	2h:40m:13s

1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172



1173 **Table 3.** RefSeq genome sequences reclassified in this study based on the diverse genomic evidence  
1174 presented herein (see Figures 5, 6 and 7).

DEFINITION (RefSeq classification)	Species/reclassification*	Status	Fragments	BioProject	BioSample	PMID
<i>Stenotrophomonas</i> sp. 69-14	<i>S. acidaminiphila</i>	draft	27	PRJNA279279	SAMN05660631	NA
<i>Stenotrophomonas maltophilia</i> ZBG7B	<i>S. chelatiphaga</i>	draft	145	PRJNA272355	SAMN03280975	26659682
<i>Stenotrophomonas maltophilia</i> AA1	genospecies 2 (Sgn2/Smc2)	complete	1	PRJNA224116	SAMN06130959	28275097
<i>Stenotrophomonas maltophilia</i> 5BA-I-2	genospecies 3	draft	4	PRJNA224116	SAMN02641498	24604648
<i>Stenotrophomonas</i> sp. 92mfc06.1	genospecies 3	draft	11	PRJNA224116	SAMN04488690	NA
<i>Stenotrophomonas maltophilia</i> PierC1	genospecies 3	draft	59	PRJEB8824	SAMEA3309462	26276674
<i>Stenotrophomonas</i> sp. RIT309	genospecies 3	draft	45	PRJNA224116	SAMN02676627	24812212
<i>Stenotrophomonas</i> sp. SC-N050	genospecies 3	draft	24	PRJNA224116	SAMN05720615	NA
<i>Stenotrophomonas maltophilia</i> SeITE02	genospecies 3	draft	63	PRJNA224116	SAMEA3138997	24812214
<i>Stenotrophomonas</i> sp. YR347	genospecies 3	draft	11	PRJNA224116	SAMN05518671	NA
<i>Stenotrophomonas maltophilia</i> B4	genospecies 4	draft	180	PRJNA224116	SAMN03753636	NA
<i>Stenotrophomonas maltophilia</i> Sm41DVV	genospecies 4	draft	26	PRJNA323790	SAMN05188789	27540065
<i>Stenotrophomonas maltophilia</i> SmCVFa1	genospecies 4	draft	30	PRJNA323845	SAMN05190067	27540065
<i>Stenotrophomonas maltophilia</i> 13146	<i>S. bentonitica</i> complex	draft	60	PRJNA224116	SAMN07237143	NA
<i>Stenotrophomonas maltophilia</i> BR12S	<i>S. bentonitica</i>	draft	80	PRJNA224116	SAMN03456145	26472823
<i>Stenotrophomonas</i> sp. HMSC10F07	<i>S. bentonitica</i>	draft	63	PRJNA269850	SAMN03287020	NA
<i>Stenotrophomonas</i> sp. LM091	<i>S. bentonitica</i>	complete	1	PRJNA344031	SAMN05818440	27979933
<i>Stenotrophomonas maltophilia</i> PML168	<i>S. bentonitica</i>	draft	97	PRJNA224116	SAMEA2272452	22887661
<i>Stenotrophomonas</i> sp. Leaf70	<i>S. nitritireducens</i>	draft	11	PRJNA224116	SAMN04151613	26633631
<i>Stenotrophomonas</i> sp. KCTC 12332	<i>S. terrae</i> complex	complete	1	PRJNA310387	SAMN04451766	28689013
<i>Stenotrophomonas nitritireducens</i> 2001	<i>S. terrae</i> complex	complete	1	PRJNA224116	SAMN05428703	NA
<i>Stenotrophomonas maltophilia</i> S028	<i>Stenotrophomonas</i> sp.					
<i>Stenotrophomonas rhizophila</i> QL-P4	<i>Stenotrophomonas</i> sp.	complete	1	PRJNA326321	SAMN05276013	NA

1175 \*The numbered genospecies correspond to novel unnamed species identified by Ochoa-Sánchez and  
1176 Vinuesa (2017) and in this study. Unnamed species classified as members of the *S. maltophilia sensu*  
1177 *lato* clade (Fig. 5) are labeled as *S. maltophilia s. l.* Strains assigned to the *S. terrae* complex most  
1178 likely represent novel species related to *S. terrae*.

1179  
1180  
1181

## 1182 FIGURE LEGENDS

1183

1184 **Figure 1.** Simplified flow-chart of the GET\_PHYLOMARKERS pipeline showing only those parts  
1185 used and described in this work. The left branch, starting at the top of the diagram, is fully under  
1186 control of the master script `run_get_phylomarkes_pipeline.sh`. The names of the worker scripts called  
1187 by the master program are indicated on the relevant points along the flow. Steps involving repetitive  
1188 computational processes, like generating multiple sequence alignments or inferring the corresponding  
1189 gene trees, are run in parallel with the aid of GNU parallel, which is called from  
1190 `run_parallel_cmmnds.pl`. The right-hand branch at the top of the diagram summarizes the analyses that  
1191 can be performed on the pan-genome matrix (PGM). In this work we only present the estimation of  
1192 maximum-likelihood and parsimony pan-genome phylogenies. However, unsupervised learning  
1193 approaches are provided by the `hcluster_pangenome_matrix.sh` script (not shown) for statistical  
1194 analysis of the PGM. In addition, the `plot_matrix_heatmap.sh` script was used to analyze average  
1195 nucleotide identity matrices generated by `get_homologues.pl`. It implements the unsupervised learning  
1196 method described in this work to define the optimal number of clusters in such matrices. The  
1197 `plot_matrix_heatmap.sh` script is distributed with the GET\_HOMOLOGUES suite.

1198

1199

1200 **Figure 2.** Density plots showing the distribution of the number of fragments of the *Stenotrophomonas*  
1201 genomes available in RefSeq as of August 2017, plus the genome of *S. bentonitica* VV6, released in  
1202 January 2018. **A)** Distribution of the number of fragments in the assemblies of 170 annotated  
1203 *Stenotrophomonas* genomes as a function of assembly status (contigs vs. scaffolds) plus 7  
1204 *Xanthomonas* genomes used as outgroup to root the tree. Inset tables provide additional summary  
1205 statistics of the RefSeq assemblies. **B)** Distribution of the number of fragments in the assemblies of the  
1206 119 genomes selected for the analyses presented in this study, which include 102 reference  
1207 *Stenotrophomonas* genomes, 10 new genomes generated for this study, and 7 complete *Xanthomonas*  
1208 spp. genomes.

1209

1210

1211 **Figure 3.** Combined filtering actions performed by GET\_HOMOLOGUES and  
1212 GET\_PHYLOMARKERS to select top-ranking phylogenetic markers to be concatenated for  
1213 phylogenomic analyses, and benchmark results of the performance of the FastTree (FT) and IQ-TREE  
1214 (IQT) maximum-likelihood (ML) phylogeny inference programs. **A)** Venn-diagram indicating the  
1215 number consensus and algorithm-specific core-genome orthologous clusters. **B)** Parallel box-plots  
1216 summarizing the computation time required by FT and IQT when run under “default” (FTdef, IQTdef)  
1217 and thorough (FThigh, IQThigh) search modes (s\_type) on the 239 consensus clusters, as detailed in  
1218 the main text. Statistical significance of differences between treatments were computed with the  
1219 Kruskal-Wallis (robust, non-parametric, ANOVA-like) test. **C)** Distribution of SH-*alrt* branch support  
1220 values of gene-trees found by the FT<sub>high</sub> and IQT<sub>high</sub> searches. Statistical significance of differences  
1221 between the paired samples was computed with the Wilcoxon signed-rank test. This is a non-parametric  
1222 alternative to paired t-test used to compare paired data when they are not normally distributed. **D)**  
1223 Association plot (computed with the `vcd` package) summarizing the results of multi-way Chi-Square  
1224 analyses of the *lnL* score ranks (1 to 4, meaning best to worst) of the ML gene-trees computed from the  
1225 set of 105 codon alignments passing the `kdtrees` filter in the IQT<sub>high</sub> run (**Table 2**) for each search-  
1226 type. The height and color-shading of the bars indicate the magnitude and significance level of the  
1227 Pearson residuals. **E)** Statistical analysis (Kruskal-Wallis test) of the distribution of consensus values  
1228 from majority-rule consensus trees computed from the gene trees passing all the filters, as a function of

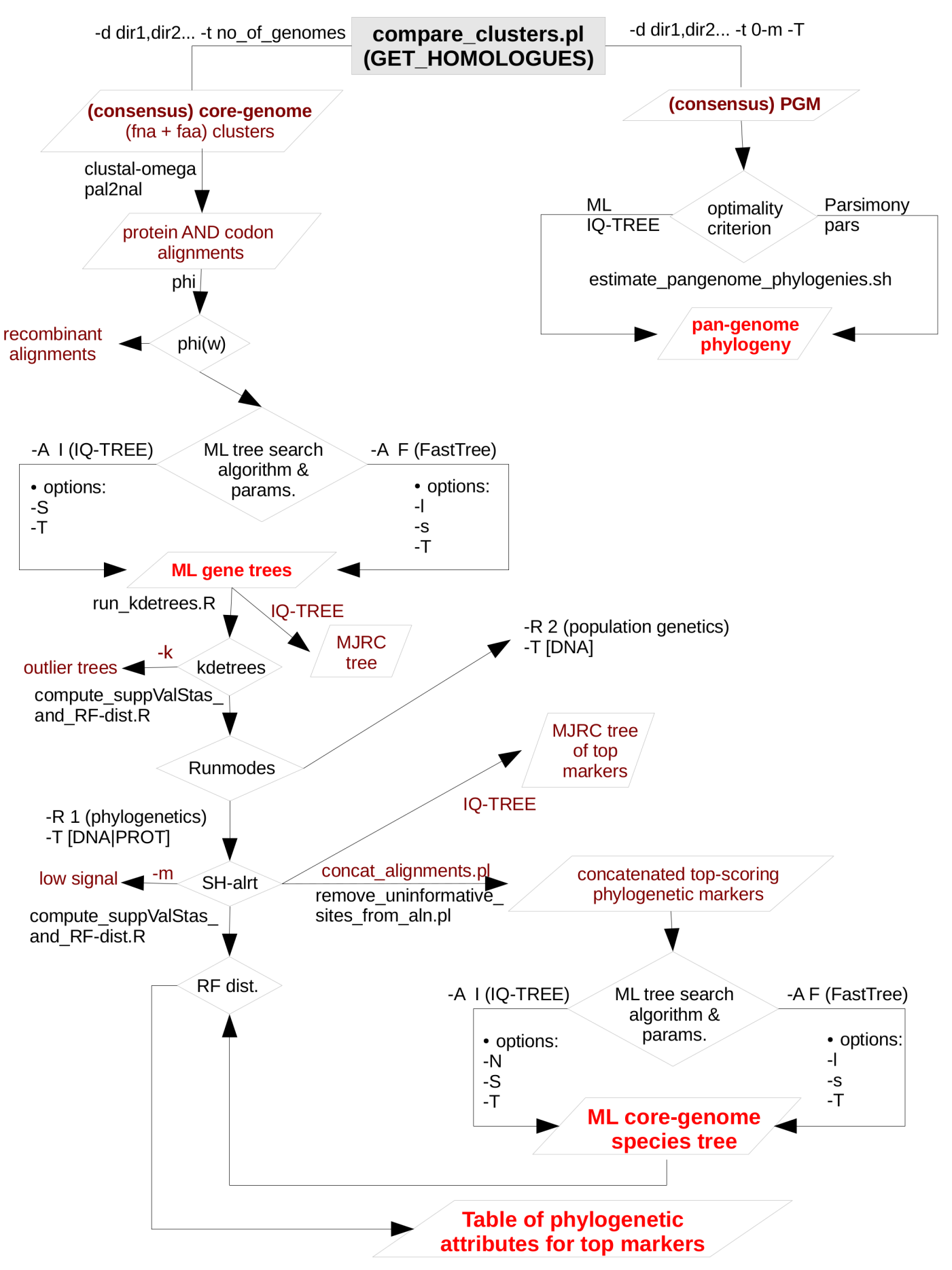
1229 search-type. **D**) Statistical analysis (Kruskal-Wallis test) of the distribution of the edge-lengths of  
1230 species-trees computed from the concatenated top-scoring markers, as a function of search-type.  
1231

1232 **Figure 4.** Comparative analysis of log-likelihood tree search profiles. **A**) Sorted lnL profile of FastTree  
1233 (FT) tree searches launched from 1000 random trees + 1 BioNJ phylogeny, using the “thorough” tree-  
1234 search settings described in the main text and the 55 top-ranking markers (26,988 non-gapped, variable  
1235 sites) selected by the IQThigh run for 119 genomes (**Table 2**). The dashed blue line indicates the score  
1236 of the search initiated from the BioNJ tree. **B**) Sorted lnL profile of 50 independently launched IQ-  
1237 TREE (IQT) searches under the best-fitting model using the same matrix as for the FT search.  
1238

1239 **Figure 5.** Best maximum-likelihood core-genome phylogeny for the genus *Stenotrophomonas* found in  
1240 the IQ-TREE search described in **Fig. 4B**, based on the supermatrix obtained by concatenation of 55  
1241 top-ranking alignments (**Table 2**). The tree was rooted using the *Xanthomonas* spp. sequences as the  
1242 outgroup. Arrows highlight genomes not grouping in the *S. maltophilia sensu lato* clade (Smsl), for  
1243 which we suggest a reclassification, as summarized in **Table 3**. Black arrows indicate misclassified  
1244 strains, while gray ones mark unclassified genomes. The shaded area highlights the strains considered  
1245 as members of the *S. maltophilia* complex (Smc). The genospecies 1 and 2 (Sgn1 = Smc1; Sgn2 =  
1246 Smc2) were previously recognized as separate species-like lineages by Ochoa-Sánchez and Vinuesa  
1247 (2017). Strains grouped in the Smsl clade are collapsed into sub-clades that are perfectly consistent  
1248 with the cluster analysis of core-genome average nucleotide identity (cgANIb) values presented in **Fig.**  
1249 **7** at a cutoff-value of 95.9%. Integers in parentheses correspond to the number of genomes in each  
1250 collapsed clade. **Supplementary Figure S4** displays the same tree in non-collapsed form. Strains from  
1251 genospecies 1, 3 and 5 (Sgn1, Sgn3, Sgn5) marked with an asterisk may represent additional species,  
1252 according to the cgANIb values. Nodes are colored according to the lateral scale, which indicates the  
1253 approximate Bayesian posterior probability values. The scale bar represents the number of expected  
1254 substitutions per site under the best-fitting GTR+ASC+F+R6 model.  
1255

1256 **Figure 6.** Maximum-likelihood pan-genome phylogeny estimated with IQ-TREE from the consensus  
1257 pan-genome displayed in the Venn diagram. Clades of lineages belonging to the *S. maltophilia* complex  
1258 are collapsed and are labeled as in **Fig. 5**. Numbers on the internal nodes represent the approximate  
1259 Bayesian posterior probability/UFBot2 bipartition support values (see methods). The tabular inset  
1260 shows the results of fitting either the binary (GTR2) or morphological (MK) models implemented in  
1261 IQ-TREE, indicating that the former has an overwhelmingly better fit. The scale bar represents the  
1262 number of expected substitutions per site under the binary GTR2+F0+R4 substitution model.  
1263

1264 **Figure 7.** Application of an unsupervised learning approach to the cgANIb distance matrix to identify  
1265 statistically-consistent species-like clusters. The cgANIb matrix was converted to a distance matrix  
1266 (cgANDb) and clustered using the Ward.D2 algorithm. The optimal number of clusters ( $k$ ) was  
1267 determined with the average silhouette-width statistic. The inset shows the statistic's profile, with  $k =$   
1268 19 as the optimal number of clusters. This number corresponds to an cgANIb of 95.5 % (gray dashed  
1269 line). At a cgANDb of 4.1 % (cgANIb = 95.9%) the groups delimited by the clustering approach are  
1270 perfectly consistent with those delimited by the core- and pan-genome ML phylogenies displayed in  
1271 **Figure 5** and **Figure 6**, respectively.

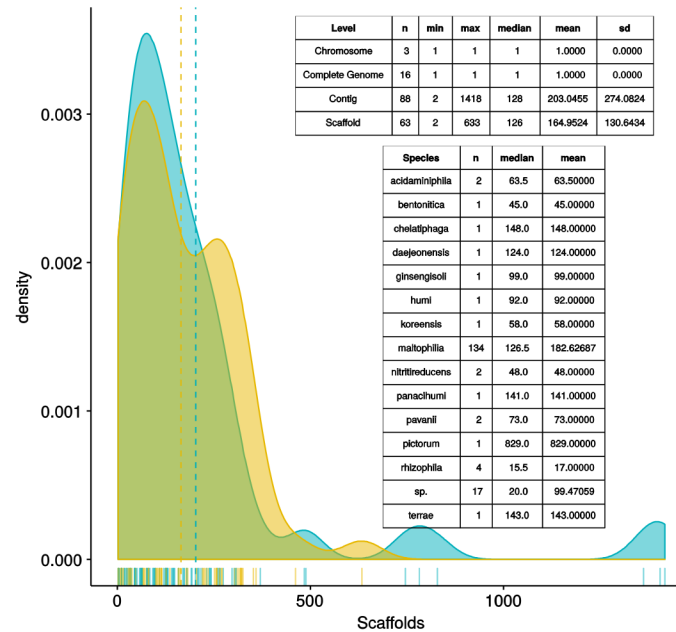




**A**Level  Contig  Scaffold

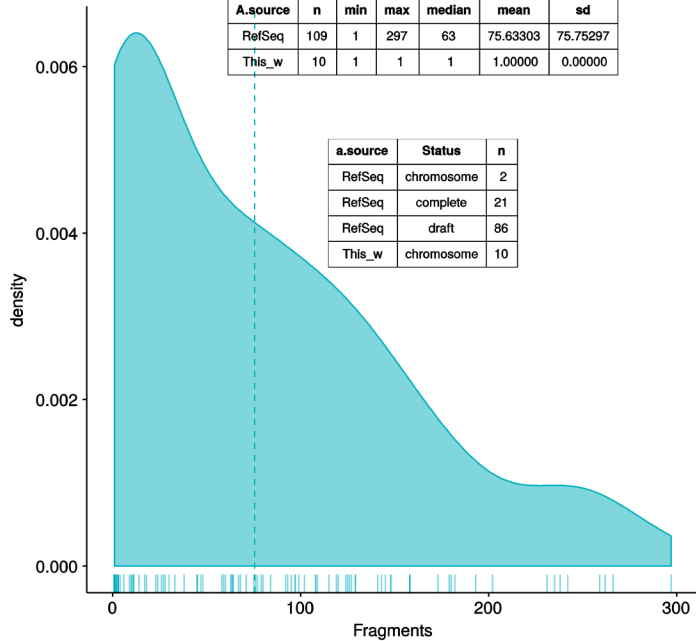
Level	n	min	max	median	mean	sd
Chromosome	3	1	1	1	1.0000	0.0000
Complete Genome	16	1	1	1	1.0000	0.0000
Contig	88	2	1418	128	203.0455	274.0824
Scaffold	63	2	633	126	164.9524	130.6434

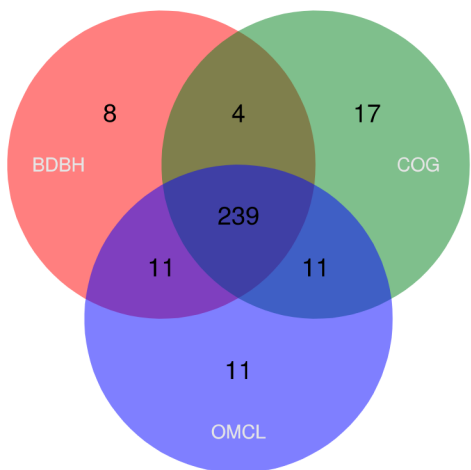
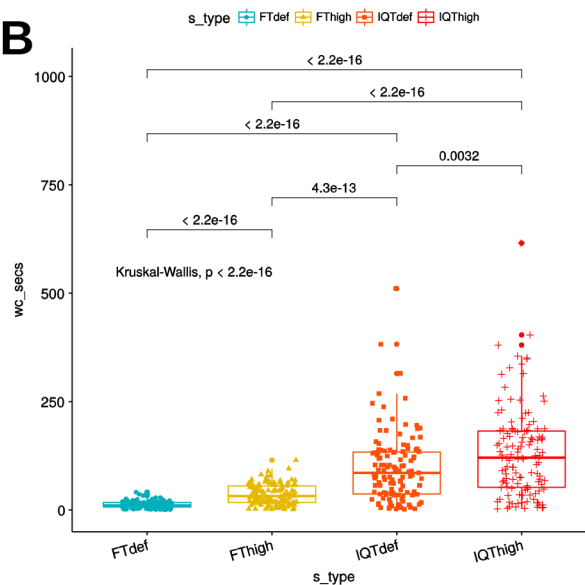
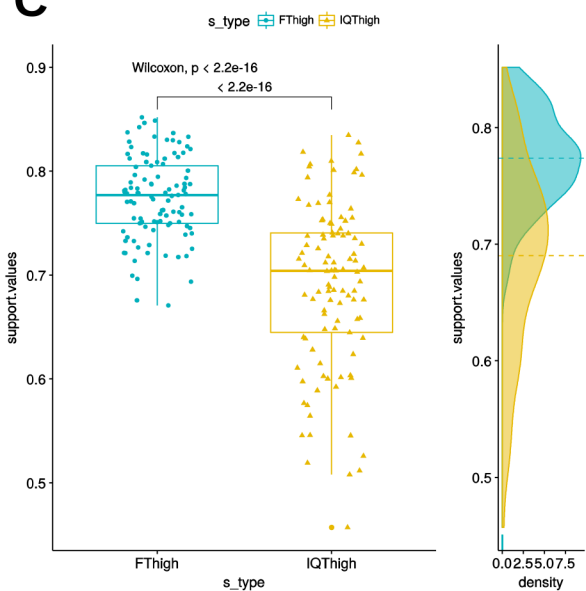
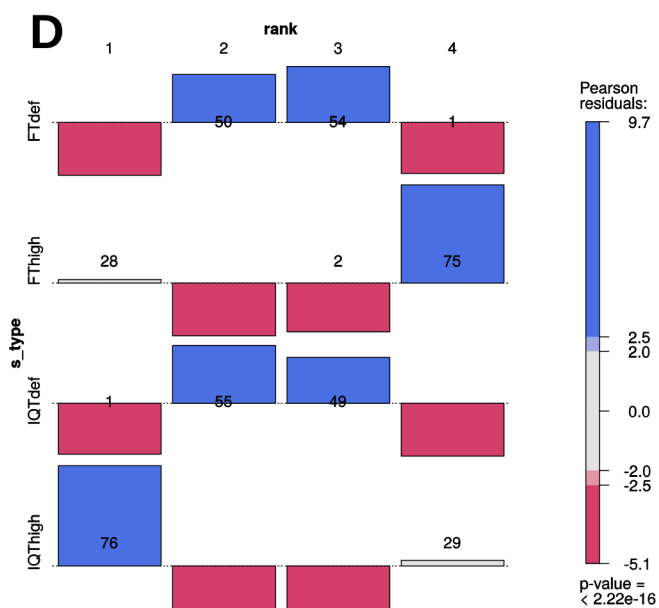
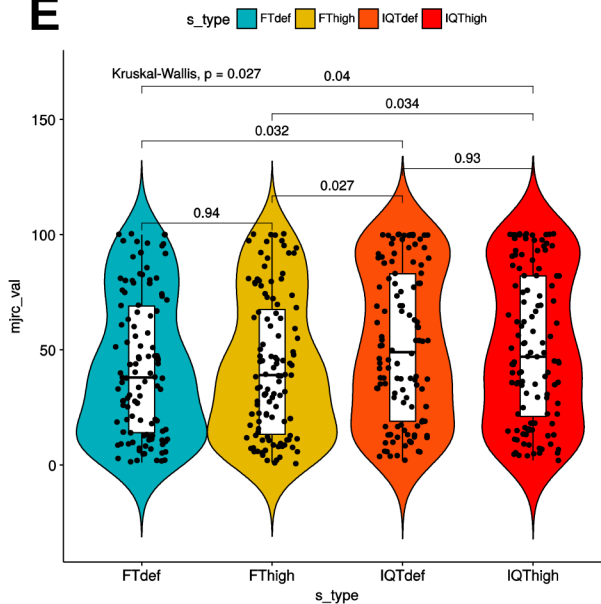
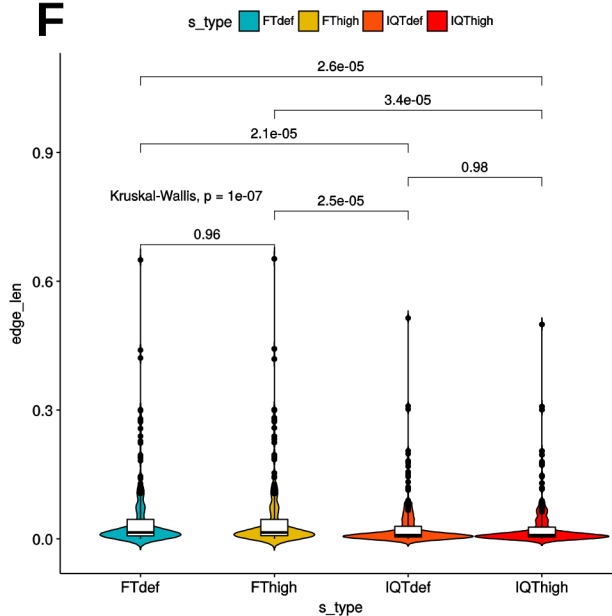
Species	n	median	mean
acidaminiphila	2	63.5	63.50000
bentonitica	1	45.0	45.00000
chelatiphaga	1	148.0	148.00000
daejeonensis	1	124.0	124.00000
ginsengisoli	1	99.0	99.00000
huml	1	92.0	92.00000
koreensis	1	58.0	58.00000
maltophilia	134	126.5	182.62687
nitritireducens	2	48.0	48.00000
panacihuml	1	141.0	141.00000
pavanli	2	73.0	73.00000
pictorum	1	829.0	829.00000
rhizophila	4	15.5	17.00000
sp.	17	20.0	99.47059
terrae	1	143.0	143.00000

**B**a.source  RefSeq

A.source	n	min	max	median	mean	sd
RefSeq	109	1	297	63	75.63303	75.75297
This_w	10	1	1	1	1.00000	0.00000

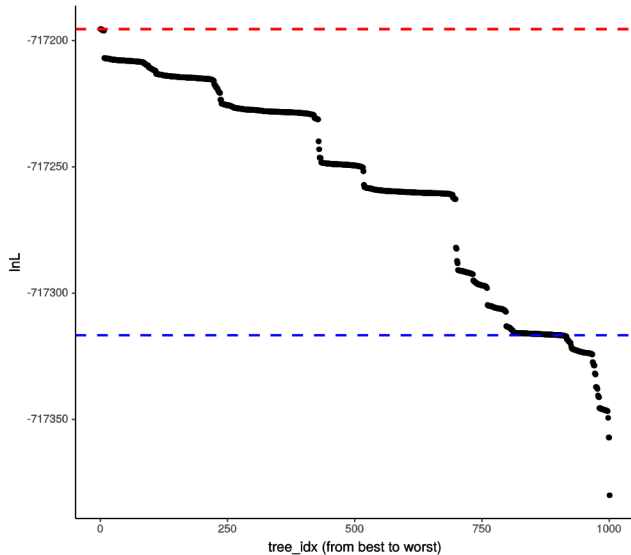
a.source	Status	n
RefSeq	chromosome	2
RefSeq	complete	21
RefSeq	draft	86
This_w	chromosome	10



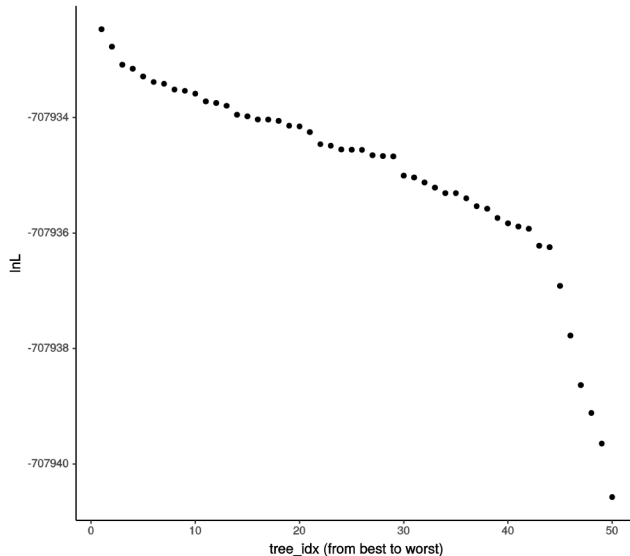
**A****B****C****D****E****F**

**A**

InL profile for 1001 FastTree searches  
max InL = -717195.373 ; min InL = -717380.026  
NJidx = 905; InL = -717316.654  
max-NJ InL difference = 121.281

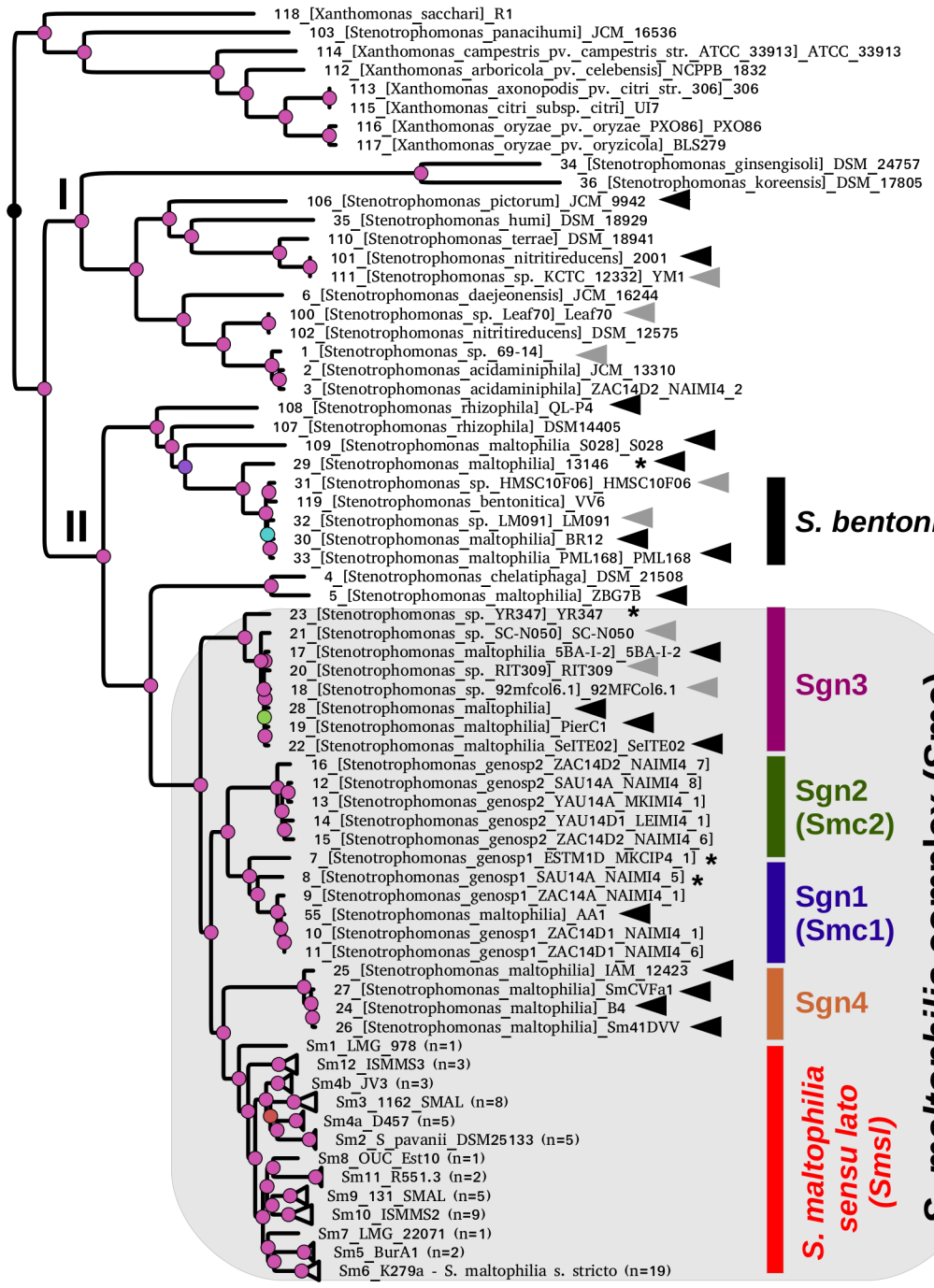
**B**

InL search profile for 50 IQ-TREE searches  
max InL = -707932.468  
min InL = -707940.573  
max-min InL difference = 8.105



aBypp

- /0.953
- /0.986
- /0.989
- /0.994
- /0.997
- /0.998
- /0.999
- /1



**S. bentonitica**

**Sgn3**

**Sgn2 (Smc2)**

**Sgn1 (Smc1)**

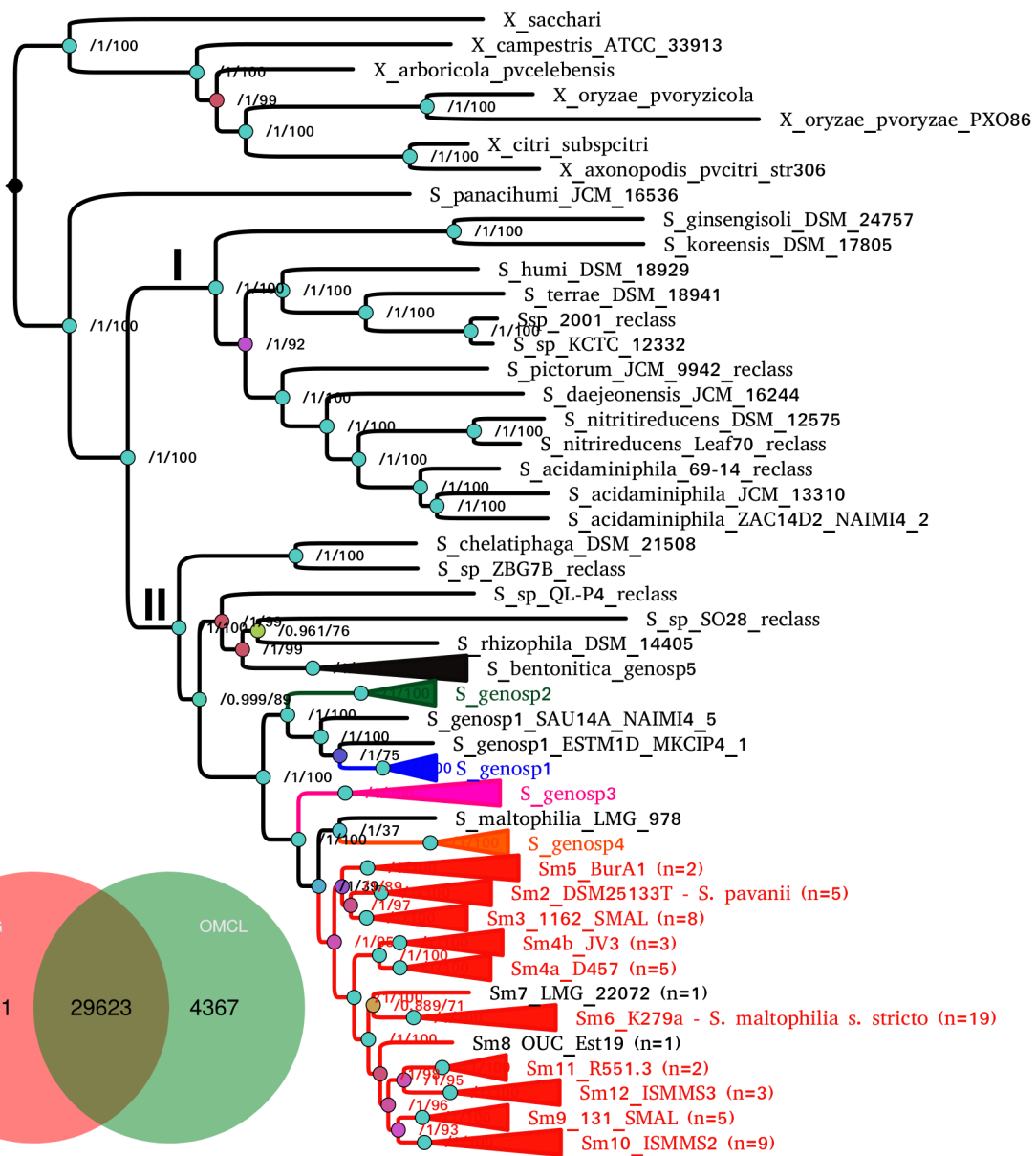
**Sgn4**

**S. maltophilia sensu lato (Sm1)**

**S. maltophilia complex (Smc)**

0.08





Model	LnL	df	AIC	AICc	BIC
GTR2+F0+R4	-322836.2380	242	646156.4760	646160.4791	648164.1821
MK+FQ+R2	-337220.7785	237	674915.5571	674919.3962	676881.7817

