

PanACoTA: A modular tool for massive microbial comparative genomics

Amandine PERRIN^{1,2,3} and Eduardo ROCHA¹

¹Microbial Evolutionary Genomics, CNRS, UMR3525, Institut Pasteur, 28, rue Dr Roux, Paris, 75015, France

²Sorbonne Universite, College doctoral, F-75005 Paris, France

³Bioinformatics and Biostatistics Hub, Department of Computational Biology, Institut Pasteur, USR 3756 CNRS, Paris, France

Abstract

The gene repertoires of microbial species, their pangenomes, evolve very fast. Their study facilitates the discrimination between lineages and reveals which genes drive their recent adaptation. It has therefore become a key topic of study in microbial evolution and genomics. Yet, the increase in the number of genomes available to certain species, now reaching many thousands, complicates the establishment of the basic building blocks of comparative genomics. Here, we present PanACoTA, a tool that allows to download all genomes of a species, build a database with those passing quality and redundancy controls, define uniform annotation, and use them to build a pangenome, several variants of core or persistent genomes, their alignments, and a rapid but accurate phylogenetic tree. While many programs have become available in the last few years to build pangenomes, we have focused on a method that tackles all the key steps of the process, from download to phylogenetic inference. This was conceived in a modular way, i.e. while all steps are integrated, they can also be run separately and multiple times to allow rapid and extensive exploration of the space of parameters of interest. The software is built in Python 3 and includes features to facilitate its installation and its future development. We believe PanACoTa is an interesting addition to the current set of bioinformatics software for comparative genomics, since it will accelerate and standardize the more routine parts of the work, allowing microbial genomicists to more quickly tackle their specific questions.

Keywords— Software, bacteria, annotation, core genomes, pangenomes, evolutionary analyses

1 Introduction

Low cost of sequencing and the availability of hundreds of thousands of genomes have made comparative genomics a basic toolkit of many microbiologists, geneticists, and evolutionary biologists. Many bacterial species of interest have now over 100 genomes publicly available in the GenBank RefSeq reference database, and a few have more than ten thousand. This

27 trend will increase with the ever decreasing costs of sequencing, the availability of long-read
28 technologies, and the use of whole-genome sequencing in the clinic for diagnostics and epidemi-
29 ology. As a result, researchers that would like to use available data are faced with extremely
30 large amounts of data to analyze. Comparative genomics has spurred important contributions
31 to the understanding of the organization and evolution of bacterial genomes in the last two
32 decades [1] [2]. It has become a standard tool for epidemiological studies, where the analysis
33 of the genes common to a set of strains - the core or persistent genome - provides unrivalled
34 precision in tracing the expansion of clones of interest [3] [4]. The use of routine sequencing
35 in the clinic will further require rapid and reliable analysis tools to query thousands, and soon
36 possibly millions of genomes from a single species [5]. Population genetics also benefits from
37 this wealth of data because one can now track in detail the origin and fate of mutations of
38 genetic acquisitions to understand what they reveal of adaptive or mutational processes [6].
39 Finally, genome-wide association studies have been recently adapted to bacterial genetics, to
40 account for variants in single nucleotide polymorphism and gene repertoires [7]. They hold the
41 promise of helping biologists to identify the genetic basis of phenotypes of interest. Given the
42 high genetic linkage in bacterial genomes, these studies may require extremely large datasets
43 to detect small effects. More specifically, reverse vaccinology is also a noteworthy application
44 of these pangenomics methods, to identify novel potential antigens among core surface-exposed
45 proteins of a given clade [8].

46 The availability of large genomic datasets puts a heavy burden on researchers, especially
47 those that lack extensive training in bioinformatics, because their analysis implicates the use of
48 automatic processes, efficient tools, extensive standardization, and quality control. Many tools
49 have been recently developed to make rapid searches for sequence similarity with excellent recall
50 rates for highly similar sequences [9] [10] [11].

51 Other tools also provide methods to rapidly cluster large numbers of sequences in families
52 of sequence similarity, to get the families common to a set of genomes, to align them, or
53 to produce their phylogeny, four cornerstones of comparative genomics. A number of recent
54 programs have recently been published that include some of these tools to compute bacterial
55 pangenomes (for a review, see [12]). Many of these programs compute alignments and clusters
56 of families using programs that are very fast. Some use tools that are known to sacrifice

57 accuracy for very high speed, such as DIAMOND [9], USEARCH [13] and CD-HIT [14]. The
58 latter is used, among others by Roary [15], which is currently the most popular tool to compute
59 pangenomes, and Panaroo [16], a very recent tool aiming at reducing the impact of erroneous
60 automated annotation of prokaryotic genomes. BPGA [17], using USEARCH or CD-HIT to
61 cluster proteins, also provides some downstream analyses. PanX [18], which has an outstanding
62 graphical interface, uses DIAMOND to search for similarities among genes.

63 More recently, SonicParanoid introduced the use of the highly efficient and accurate pro-
64 gram mmseqs2 to build pangenomes, and PPanGGOLiN used the same tool to provide a
65 method to statistically class pangenome families in terms of their frequency [19] [20] [21] .
66 Some recent programs also use graph-based approaches to further refine the pangenomes, such
67 as PPanGGOLiN and Panaroo [16]. For that matter, the analysis of a dataset of 319 *Kleb-*
68 *siella pneumoniae* genomes by both tools provide very similar results [16]. Some tools, such
69 as PIRATE [22] have also been recently developed to cluster orthologues between distant
70 genomes. However, all these programs lack some or all of initial and final steps that are essen-
71 tial in comparative genomics, including download, quality control, alignment and phylogenetic
72 inference. This spurred the development of PanACoTA (PANgenome with Annotations, COre
73 identification, Tree and corresponding Alignments). To take advantage of the vast amount of
74 genomic information publicly available, one needs six major blocks of operations. (1) Gather
75 a set of genomes of a clade automatically. This requires some quality control, to avoid drafts
76 with an excessive number of contigs. It is also often convenient to check that the genomes are
77 not too redundant, to minimize computational cost and biases due to pseudo-replication. On
78 the other side, it is important to check that genomes are neither too unrelated, to eliminate
79 genomes that were misclassified in terms of bacterial species (or the taxonomic organisation of
80 relevance). (2) Define *a priori* an uniform nomenclature and annotation, without which the
81 calculation of pangenomes and core genomes becomes unreliable for large datasets. (3) Produce
82 the pangenome, a matrix with the patterns of presence absence of each gene family in the set of
83 genomes, using an accurate, simple, and fast method. (4) Use the pangenome to identify sets
84 of core or persistent genes. (5) Produce multiple alignments of the gene families of the core
85 or persistent genomes. (6) Finally, produce quickly a reasonably accurate phylogeny of the set
86 of core/persistent genes. These four collections of data, pangenome, core genome, alignments,

87 and phylogenetic tree, are the basis of most microbial comparative genomics studies. At the
88 end of this process, the researcher can produce more detailed analyses, specific to the questions
89 of interest, which often lead to changes such as including/excluding taxa, changing the limits
90 of sequence similarity, increasing alignment accuracy, or rebuilding phylogenies using different
91 methods. Such re-definitions can be achieved more efficiently when pipelines are modular and
92 allow to re-start the analyses at several key points in the process.

93 Considering the current availability of pipelines for microbial comparative genomics, we have
94 built one that is modular, easy to setup, uses state-of-the-art tools, and allows simple re-use of
95 intermediate results. The goal was to provide a pipeline that allows to download all genomes
96 from a taxonomic group and make all basic comparative genomics work automatically. The
97 pipeline is entirely built in a single language, Python v3, and uses modern methods to facilitate
98 its future maintenance and to limit unwanted behaviour. PanACoTA is freely available under
99 the open source GNU AGPL license. Here, we describe the method and illustrate it with an
100 analysis of two datasets of 225 complete and 3980 complete or draft genomes of *Klebsiella*
101 *pneumoniae*. This species is interesting for our purposes because there are many genomes
102 available and it has a very open pangenome [23]. The first dataset describes a situation where
103 sequence quality is usually high, and the second illustrates how the method scales-up to a very
104 large dataset where sequence is of lower quality or genes are fragmented due to lack of complete
105 assembly. The procedure is detailed in the Methods section, whereas the illustration of its use,
106 and how it changes in relation to key options in the two datasets, is detailed in the Results
107 section.

108 2 Methods

109 PanACoTa is implemented in 6 sequential modules, described in the six sections below. It was
110 designed to allow the use of a module without requiring the use of previous one(s). This allows
111 to start or stop at any step and re-run an analysis with other parameters (see Figure 1).

112 2.1 Datasets

113 The first module of PanACoTA - `prepare` - allows to fetch all genomes from a given NCBI taxon-
114 omy ID. This uses scripts from `ncbi_genome_download` library (<https://github.com/kbclin/>

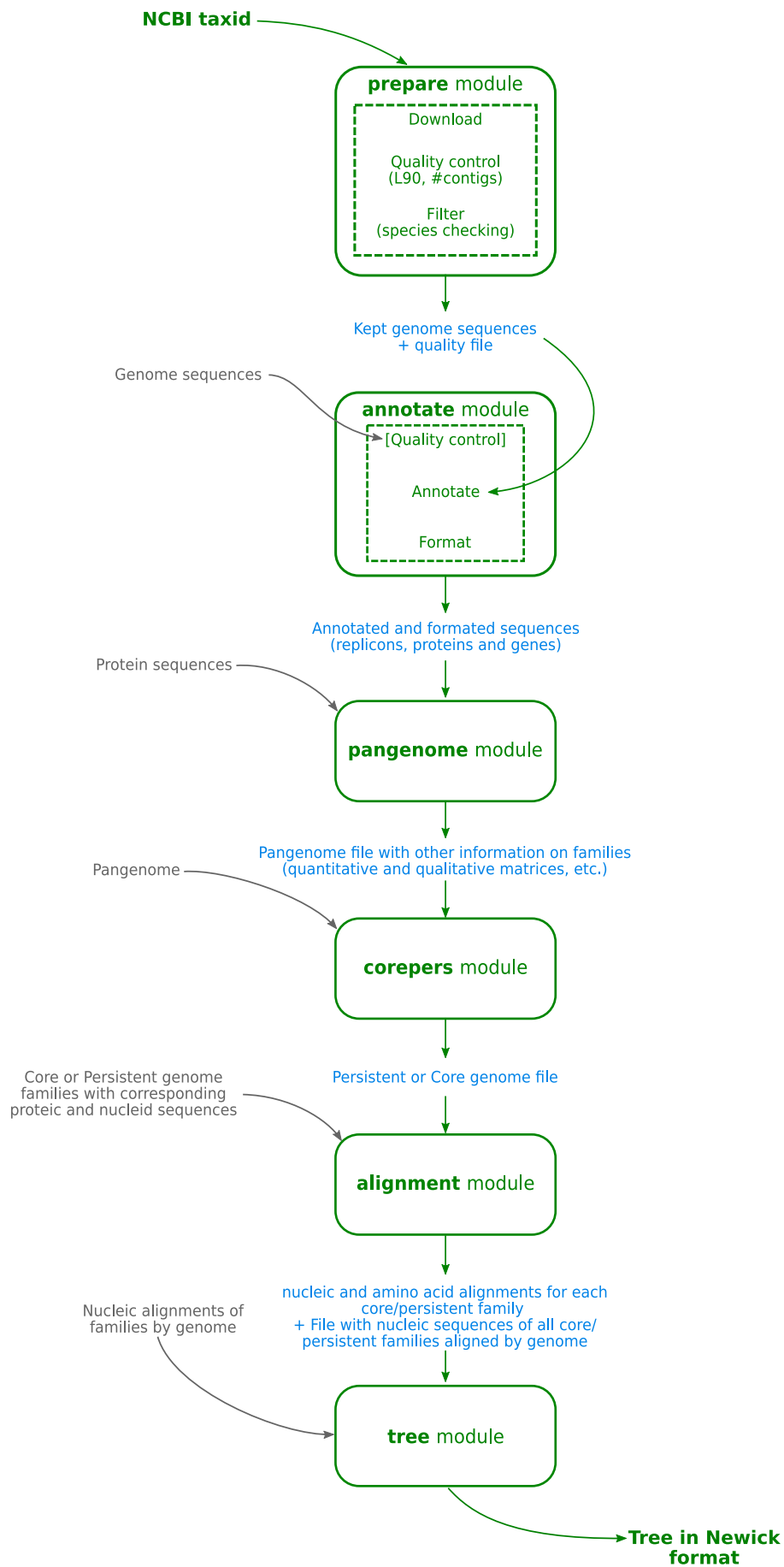


Figure 1: Overview of PanACoTA method

115 [ncbi-genome-download](#)). PanACoTA retrieves the corresponding compressed non-annotated
116 fasta files of the genomes from the NCBI ftp.

117 In this paper, we use two datasets of *Klebsiella pneumoniae* genomes to illustrate how PanA-
118 CoTA functions. The first one, called hereafter DTS1, contains all complete and draft genomes
119 downloaded from the NCBI refseq database on October 2018 10th. The second dataset, DTS2,
120 is the subset of DTS1 containing only the complete genomes (genomes with `assembly_level`
121 = `Complete Genome`, based on the NCBI summary file). This initial database was then given
122 to a quality control procedure, described in the next section. At the end, this module outputs
123 a database with the genomes admissible after this control step: 3980 genomes for DTS1, with
124 a subset of 225 complete genomes for DTS2. The accession numbers of all the genomes are
125 indicated in Table S1.

126 2.2 Quality control procedure

127 The goal of this module is to remove genomes that do not conform with two types of basic
128 requirements in terms of assembly and taxonomy. It is done by the `prepare` module after
129 downloading the genomes, or by the `annotate` module before the annotation step (if the user
130 did not use the `prepare` module). This latter option is useful when the goal is to analyse a pre-
131 defined list of genomes, some of which are eventually not available in GenBank (e.g. in-house
132 sequencing). The module receives as input a set of fasta sequences.

133 The first goal of this control procedure is to filter genomes in terms of sequence quality.
134 Since there is usually no standard description of the quality of the sequence assembly in RefSeq
135 genomes, the program infers it from the sequences. First, it is common usage to put stretches
136 of 'N' to separate contigs in a same fasta sequence. To have a better idea of the sequence
137 quality, and be able to do the analysis more efficiently, we first split sequences at each stretch
138 of at least 5 'N' (this number of 5 can be changed by the user) to get one fasta entry per
139 contig. Assuming that the user is analyzing genomes from the same species, those genomes
140 should have relatively similar characteristics in terms of number of contigs and length. Hence,
141 PanACoTA first calculates two key measures for each genome: the total number of contigs, and
142 the L90 (the minimum number of contigs necessary to get at least 90% of the whole genome).
143 Very high values of these two variables are usually an indication of low quality of sequencing or

144 assembling. They often result in the annotation of numerous truncated genes that spuriously
145 increase the size of pangenomes (because a real gene is split into numerous open reading frames
146 that are classed in distinct families). These poorly assembled genomes also complicate studies
147 of comparative genomics whenever studying genetic linkage is important, because they are
148 dominated by small contigs. The thresholds can be specified by the user. Values by default are
149 set to less than 1000 contigs and L90 lower than 100. Genomes exceeding one of these values
150 are excluded from the rest of the analysis.

151 The second part of the procedure is a filter dedicated to remove redundant and miss-classified
152 genomes. This is done based on the genetic distance between pairs of genomes, as calculated
153 by Mash [24]. We chose Mash for this distance filtering step because it can be computed very
154 fast and is accurate for closely related genomes. Mash reduces each genome sequence to a
155 sketch of representative k-mers, using the MinHash technique [25]. It then compares those
156 sketches, instead of the full sequences. This output Mash distance D strongly correlates with
157 alignment-based measures such as the Average Nucleotide Identity (ANI) which is based on
158 whole-genome sequence comparisons using the blast algorithm [26]: $D \approx 1 - ANI$. For ANI
159 in the range of 90–100%, the correlation with Mash distance is even stronger when increasing
160 the sketch size. Since pangenomes are typically computed for a single bacterial species, we are
161 here using Mash to discriminate genomes having at least 94% identity. A few recent programs
162 have been published showing slightly more accuracy than Mash, but we found them too slow
163 for the use as a systematic filter. For example, using 15 cores, FastANI [27] requires around
164 1h15 to compare all pairs of 200 genomes (40,000 pairwise comparisons), where Mash with a
165 sketch size of 10^6 does the task in less than 3 minutes. Hence, a user requiring a finer grade
166 study of ANI may wish to post-analyse the data from FastANI, instead of running the `prepare`
167 module. However, it is impractical to make all pairwise analyses of very large datasets where
168 one often needs to perform millions of pairwise comparisons.

169 Bacterial species are usually defined as groups of genomes at more than 94% identity [28],
170 and this will typically be used as the upper value for D (`max_mash_dist` is a modifiable param-
171 eter that is fixed by default at 0.06). On the other extreme, genomes with very high similarity
172 (corresponding to low Mash distances) provide very similar information. They can be excluded
173 to lower the computational resources required for the analysis and to diminish eventual over-

174 sampling of certain clades, which could lead to biased results. PanACoTA sets `min_mash_dist`
175 to 10^{-4} by default, but this parameter can also be specified by the user. This distance represents
176 one point change every ten genes on average and may be close to the sequencing and assembling
177 accuracy of many draft genomes.

178 The two procedures, quality control and Mash filtering, are linked together. The information
179 on the number of contigs and L90 is useful to chose the genome that is kept between a pair of
180 very similar genomes. In summary, the control procedure works as follows:

- 181 • Genomes with an excessively high number of contigs or L90 are excluded.
- 182 • Genomes are primarily sorted by increasing L90 value, and secondarily by increasing
183 number of contigs to produce a list ordered in terms of quality.
- 184 • The genomes are compared with Mash. For that, the first genome of the ordered list (the
185 one with best quality) is compared to all the others. The ones which do not obey to the
186 distance thresholds are discarded. The procedure then passes to the subsequent genome
187 in the ordered list (if not rejected before), compares it to all remaining genomes, and
188 discards those not respecting the thresholds. The process continues until the ordered list
189 is exhausted.

190 The output of this module is a database with the genomes that passed the two steps of the
191 quality control procedure. PanACoTA also provides a file listing the discarded genomes and
192 why they were discarded.

193 **2.3 Annotation**

194 The annotation of the genomes is done by the `annotate` module of PanACoTA. The input is a
195 database of fasta sequences, from the `prepare` module or directly provided by the user. If no
196 information is given on the quality control of those genomes (number of contigs and L90), this
197 quality control is done here (see previous section for more information on the quality control
198 step).

199 The goal of this module is to provide a uniform annotation of the gene positions (and
200 functions) across the dataset. PanACoTA annotates all genomes with Prokka [29]. The latter
201 uses Prodigal [30] to identify gene positions. It then adds functional annotations using a series

202 of programs, including BLAST+ [31] to search for homologs in a database of proteins taken from
203 Uniprot and HMMER3 [32] to search for proteins hitting selected profiles from TIGRFAM [33]
204 and PFAM [34]. All annotated sequences are renamed using a standard sequence header format.
205 The header of each gene contains 20 characters and provides human readable information on
206 the genome and contig of the gene, its relative position in the genome, and if it is at the border
207 of a contig (see Figure 2).

Informations contained in name:

Species
Date
Strain number
Contig number
Protein at the border (b) or inside (i) the contig
Protein number

Example:

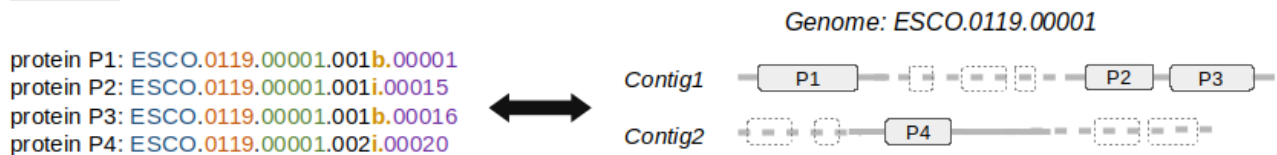


Figure 2: Description of the standard output header format for proteins annotated by PanA-CoTA.

208 If the user does not need the functional annotation, the module gives the possibility of
209 running only the gene finding part, i.e. only running Prodigal. For very large datasets it
210 is much faster to use this option and annotate a posteriori only one gene per family of the
211 pangenome using Prokka or more complete annotation systems like InterProScan [35]. The
212 output of this step consists in five files per genome: the original sequence, the genes, the
213 proteins (all in fasta format), a gff file containing all annotations and a summary information
214 file.

215 2.4 Identification of the pangenome

216 The pangenome is computed with the `pangenome` module of PanACoTA. The input is the set
217 of all proteins from all genomes, e.g. the sequence files of the 'Proteins' folder generated by the
218 `annotate` module.

219 The pangenome is the set of all protein families in the genomes. Its calculation involves com-
220 parisons between all pairs of proteins, i.e. its complexity is to the square of the number of genes
221 (and thus of genomes). To generate a reliable pangenome in a reasonable time, PanACoTA
222 calls the MMseqs2 suite [20]. The `mmseqs search` module has a very good speed/sensitivity
223 trade-off. In order to reduce time, it uses 3 consecutive search stages, with increasing sensitivity
224 and decreasing speed. Everything is highly parallelized and optimized on multiple levels. The
225 first step filters up to 99.9% of the sequences by eliminating high dissimilarities, i.e. sequences
226 not having at least two consecutive kmer matches. The second step filters out another 99% of
227 the remaining sequences using an ungapped alignment. This leaves a small amount of sequences
228 to process with an optimized version of the Smith-Waterman alignment, where only scores are
229 calculated, and not the full alignments.

230 We used the `mmseqs cluster` module included in MMseqs2 suite, with the default `Cascaded`
231 `clustering` option. This module works in two main steps. It first clusters proteins using
232 `linclust` [36], a linear time protein sequence clustering algorithm as a prefilter. Then, the
233 representative sequences of this first step are handled by the `mmseqs search` module, and
234 clustered according to its result. This second step is repeated three times, each time with a
235 higher sensitivity at the `mmseqs search` algorithm module.

236 For the clustering stage, PanACoTA uses the `Connected component` mode, because it has
237 provided results consistent with our previous methods. This mode uses transitive connections
238 to merge pairs of homologous genes: all vertices accessible via a BFS algorithm are members
239 of a cluster. Let's define the graph made from all pairwise comparisons between proteins as
240 follows: each node is a protein and there is an edge between 2 proteins if they are similar
241 (similarity beyond the given threshold). Then, two proteins are in the same family if we can
242 find a path from one to the other in the graph. If desired, the user can choose any of two other
243 clustering modes (`Greedy Set cover`, or `Greedy incremental`) using a dedicated parameter
244 while launching the PanACoTA `pangenome` module. Importantly, the tuning of the options of
245 `mmseqs2` allows the sequence similarity analyses to be exceedingly fast or extremely sensitive
246 [20]. In PanACoTA the user can change the key parameters `-min-seq-id` and `-cluster-mode`,
247 and re-run the `mmseqs cluster` module to explore their effect on the results. More specific
248 `mmseqs2` parameters have, for the time being, to be used with the standalone version of the

249 program.

250 The output of this step is a pangenome file containing one line per family, with the list
251 of all its members. PanACoTA also provides the quantitative and qualitative matrices of the
252 pangenome, as well as a tabular file giving an overview of each family composition.

253 Note that, here, we do not take into account synteny between genes in the genomes, as we
254 think that, for draft genomes, this has a limited interest. By exploring the families generated,
255 we found very few and non-significant differences. However, if the user has very well assembled
256 genomes, or has any particular reason to account for synteny, some tools have been developed,
257 like panOCT [37] [38], SynerClust [39] or PANINI [40].

258 If the user wants to do genome-wide association studies, the output qualitative matrix can
259 be directly used as input for TreeWAS [41].

260 **2.5 Identification of core and persistent genomes**

261 The classification of gene families present in a large number of taxa is done by the `corepers`
262 module of PanACoTA. The input is a pangenome file, like the one generated by the `pangenome`
263 module. In early studies, the pangenome matrix was used to identify the gene families present
264 in all genomes in a single copy: the core genome. However, the increase of the number of
265 genomes in the dataset tends to decrease drastically the size of the core genome. This is because
266 sequencing or annotation errors as well as rare deleterious polymorphism in the populations
267 lead to the rapid decrease of the number of core genes with the increase in the number of input
268 genomes. To overcome this problem, one now commonly identifies the persistent genome. A
269 family is in the persistent genome if it contains members from at least $N\%$ of the genomes. N is
270 defined by the user. The default value is 95% for datasets with more than 1000 genomes. The
271 persistent genome is more robust to rare (true or artifactual) variants. On the other hand, if the
272 goal of computing the persistent genome is to make a phylogenetic tree or analyze population
273 genetics data, one may wish to produce sets with different thresholds. Indeed, a too high value
274 will lead to a small set of persistent genes with few gaps that may be enough to infer a robust
275 phylogeny. On the other side, this will exclude many gene families from other types of analyses,
276 like detection of positive selection or recombination events. The definition of persistent genome
277 may also vary, depending on the subsequent use of the data. PanACoTA defines three types of

278 persistent genomes (see Figure 3):

- 279 • Strict-persistent: a family that contains exactly 1 member in at least $N\%$ genomes (N
280 = 100 means it is a core-family). This definition is particularly practical to reconstruct
281 phylogenies without having to handle the existence of multiple copies per genome.
- 282 • Mixed-persistent: a family where at least $N\%$ of the genomes have exactly 1 member,
283 and other genomes have either 0, either several members in the family. This definition is
284 intermediate between the other two, i.e. it includes the strict-persistent and is included
285 by the multi-persistent.
- 286 • Multi-persistent: a family with at least one member in $N\%$ of the genomes. This definition
287 is interesting to analyse patterns of diversification of nearly ubiquitous protein families.

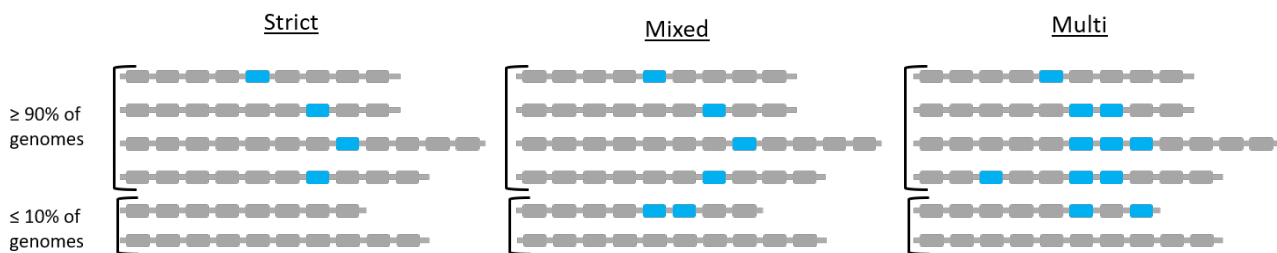


Figure 3: Different types of persistent genomes proposed by PanACoTA, with a threshold of $N = 90\%$.

288 The mixed and multi-persistent definitions are useful to include the gene families with
289 variable numbers of copies in a small (mixed) or large (multi) number of genomes when studying
290 the evolution of nearly ubiquitous gene families. It can be useful when one protein was split
291 in several parts in a few genomes because of sequencing or assembly error(s). This protein
292 family is discarded from the strict-persistent genome, while included in the mixed (and multi)
293 persistent genomes.

294 One should note that the module `corepers` does the re-analysis of the pangenome and
295 therefore it does not use a reference genome whose choice can be questionable. Re-running the
296 module is very fast, because it only requires the re-analysis of the pangenome matrix. Hence,
297 it is easy and fast to re-run the module with different parameters in different analysis to check
298 how they change the final result.

299 The output of this module is a file containing the persistent families of proteins.

300 Note that if the user wants to identify the persistent genome using a statistical approach
301 rather than using fixed thresholds, the gff file generated by `annotate` module is compatible
302 with PPanGGOLiN [21]. This software generates a persistent genome corresponding to our
303 multi-persistent version of the persistent genome (multigenic families are allowed).

304 2.6 Multiple alignments of the persistent gene families

305 The alignment of the persistent gene families is done by the `align` module of PanACoTA. Its
306 input is a persistent genome coming from the `corepers` module, or independently provided by
307 the user. When using the strict-persistent genome, all genes are aligned. When using the other
308 definitions of persistent genomes, some genomes can lack a gene or have it in multiple copies.
309 To produce phylogenetic trees from these alignments, such cases must be handled beforehand.
310 When a genome lacks a member or has more than one member (mixed or multi persistent) of
311 a given gene family, PanACoTA adds a stretch of gaps ('-') of the same length as the other
312 aligned genes. Adding a few "-" has little impact on phylogeny reconstruction. For example,
313 it has been showed that adding up to 60% of missing data in the alignment matrix could
314 still result in informative alignments [42]. In our experience, when this approach is applied to
315 within-species pangenomes, it usually incorporates less than 1% of gaps. The effect of missing
316 data should thus be negligible relative to the advantage of using the phylogenetic signal from
317 many more genes (i.e. in contrast to using the strict-persistent genome). Alignments are more
318 accurate when done at the level of the protein sequence. This has the additional advantage of
319 producing codon-based nucleotide alignments that can be used to study selection pressure on
320 coding sequences. Hence, PanACoTA translates sequences, aligns the corresponding proteins
321 and then back-translates them to DNA to get a nucleotide alignment. This last step constitutes
322 in the replacement of each amino acid by the original codon. Hence, at the end of the process,
323 the aligned sequences are identical to the original sequences.

324 PanACoTA does multiple sequence alignment using MAFFT [10] as it is often benchmarked
325 as one of the most accurate multiple alignment programs available and one of the fastest [43]).
326 It has options that allow to make much faster alignments, at the cost of some accuracy, to
327 handle very large datasets. This loss of accuracy is usually low for very similar sequences as
328 it is the case of orthologous gene families within species, and means that PanACoTA can very

329 rapidly align the persistent genome.

330 This module returns several output files: the concatenate of the alignments of all families
331 to be used for tree inference, and, for each core/persistent genome family, a file with its gene
332 and protein sequences aligned.

333 2.7 Tree reconstruction

334 The phylogenetic inference is done with the `tree` module of PanACoTA. It uses as input the
335 alignments of the `align` module or any other alignments in Fasta format.

336 This is the part that takes most time in the entire pipeline, because the time required for
337 phylogenetic inference grows very fast with the size of the dataset. Even efficient implementa-
338 tions of the maximum likelihood analyses scale with the product of the number of sites and the
339 number of taxa, which is a problem in the case of large datasets (thousands of taxa, with more
340 than ten thousands sites for each one). PanACoTA proposes several different methods to obtain
341 a phylogeny: IQ-TREE [44], FastTreeME [45], fastME [46] and Quicketree [47]. Whatever the
342 software used, the `tree` module takes as input a nucleotide alignment in Fasta format (like, for
343 example, the output of `align` module), and returns a tree in Newick format. According to its
344 needs, the user can choose one of these methods to infer its phylogenetic tree. These trees can
345 be used to build more rigorous phylogenetic inference using methods that are more demanding
346 in computational resources, e.g. by changing the options of IQ-TREE.

347 2.8 Implementation and availability

348 PanACoTA was developed in Python3, trying to follow the best practices for scientific software
349 development [48] [49]. For that, the software is versioned using git, allowing the tracking of
350 all changes in source code during PanACoTA's development. It is freely distributed under the
351 open-source AGPL licence (making it usable by many organizations) and can be downloaded
352 from <https://github.com/gem-pasteur/PanACoTA>.

353 Hosting it on GitHub allows for issue tracking, i.e. users can report bugs, make suggestions
354 or, for developers, participate to the software improvement. To provide a maintainable and
355 reliable software, we set up continuous integration process: each time a modification is pushed,
356 there is an automatic software installation checking, unit tests are done, and, if necessary, an

357 updated version of the documentation is generated, as well as an update of the singularity
 358 image on Singularity Hub.

359 As introduced just before, we also provide a complete documentation, including a step by
 360 step tutorial, based on provided genome examples, so that the user can quickly get started.
 361 It also contains more detailed sections on each module, aiming at helping users to tune all
 362 parameters, in order to adapt the run to more specific needs. This documentation also includes
 363 a 'developer' section, addressed to developers wanting to participate in the project.

364 During its execution, PanACoTA provides logging information, so that user can see real-
 365 time execution progress (a `quiet` parameter is also proposed for users needing empty stdout
 366 and stderr). This also provides log file(s) to keep track on what was ran (command-line used,
 367 time stamp, parameters used etc.).

368 3 Results and discussion

369 All execution times mentioned in this section correspond to wall clock time on 8 cores (except
 370 when the number of cores is given). A summary of all execution times can be found in Table 1.

MODULE	STEP	DTS1 (3980 genomes)	DTS2 (225 genomes)
prepare	<i>downloading</i>	1h (5805 genomes)	3min (266 genomes)
	<i>quality control</i>	<4min	~15sec
	<i>filter</i>	20min	~1min
annotate	<i>with Prokka</i>	5 days	10h
	<i>with Prodigal</i>	6h	30min
pangenome		30min	1min
corepers (1 core)		1min	5sec
align	<i>strict-persistent</i>	3h	10min
	<i>mixed-persistent</i>	7h	11min
Tree (IQ-TREE) (28 cores)	<i>strict-persistent</i>	7h (40GB RAM)	3min10
	<i>mixed-persistent</i>	24h (90GB RAM)	3min30

Table 1: Summary of execution times by (sub)module

371 3.1 Download and preparation of genome sequences

372 The first module of PanACoTA was used to download all genomes of *Klebsiella pneumoniae*
 373 using the TaxID 573. It took approximately 1h to download the 5805 *Klebsiella pneumoniae*
 374 genome sequences (including 266 complete genomes). We used the module `annotate` to make

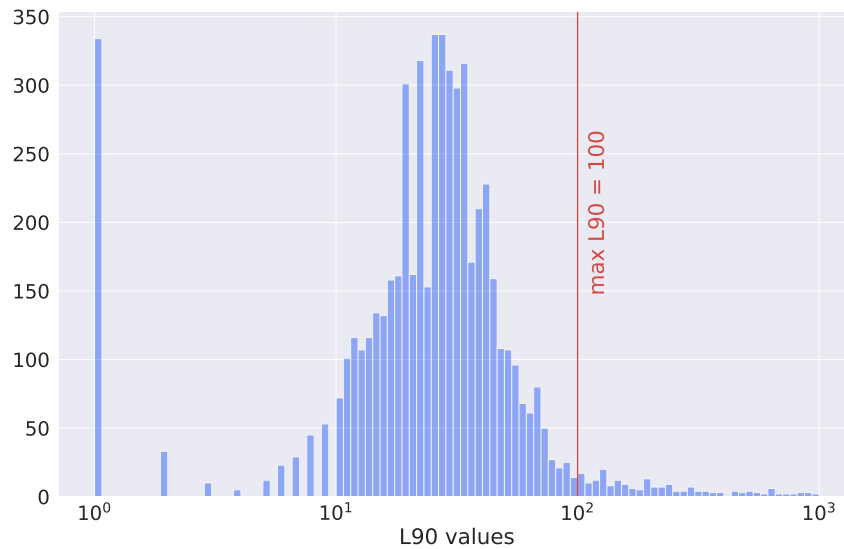
375 the quality control of those 5805 strains. For this study, we used as thresholds: L90<100 and
376 number of contigs<999. The computation of L90, number of contigs and genome size, and
377 the subsequent procedure of discarding genomes according to the thresholds took less than 4
378 minutes. This step discarded 233 draft genomes, leaving 5572 for further analysis (see Figure 4).
379 When the threshold on the number of contigs was decreased by half (number of contigs<500),
380 only 52 more genomes were removed (see Figure 4b). To define the best thresholds to the
381 analysis, the user can preview its dataset quality with a 'dry-run' of the `annotate` module.
382 Then, the user can launch the real analysis, from `prepare` or `annotate` with the adapted
383 thresholds.

384 The analysis of genetic distances across pairs of genomes is performed by Mash (K-mer
385 size of 21 (default), and sketches of at most 10000 non-redundant min-hashed k-mers). A
386 total of 1592 genomes (including 41 complete genomes) did not respect the distance thresholds
387 (`max_mash_dist` = 0.06 and `min_mash_dist` $1e^{-4}$). The vast majority (1448) were too similar
388 to other genomes, whereas 144 were too distant from other strains to be regarded as bona fide
389 *Klebsiella pneumoniae* genomes (figure 5).

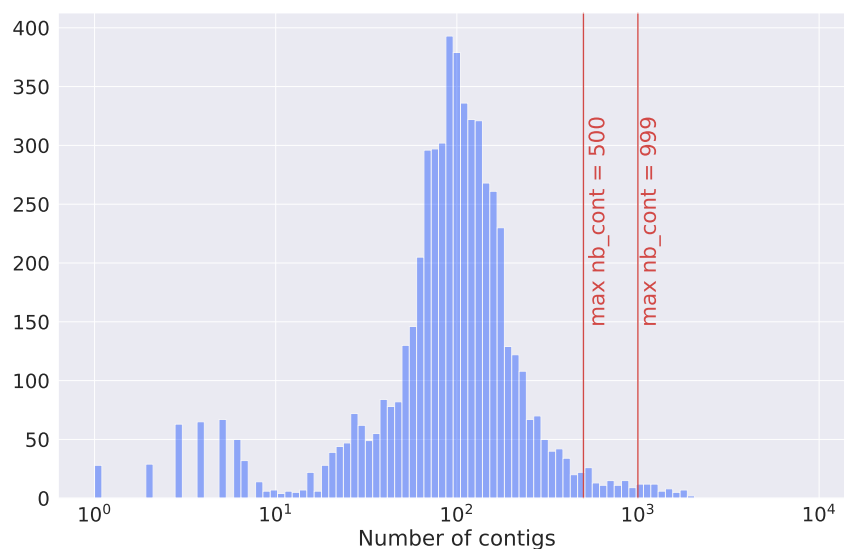
390 Some species can even be defined with narrower ANI values. For example, to identify bona
391 fide *Klebsiella pneumoniae* genomes, Kleborate (<https://github.com/katholt/Kleborate>) uses
392 Mash to compare the given assembly to a curated set of *Klebsiella* assemblies from NCBI. It
393 considers a Mash distance of ≤ 0.01 as a strong species match, and a Mash distance between
394 0.01 and 0.03 as a weak match. With our DTS1, Kleborate would have only removed 22
395 more genomes, that it identifies as *Klebsiella quasipneumoniae subspecies similipneumoniae*.
396 Our method, which is designed for any species, is thus quite consistent with Kleborate results
397 regarding the specific case of *K. pneumoniae* genomes.

398 Three genomes showed an ANI less than 84% identity, meaning they may not even be
399 from the same genus, which emphasizes the necessity of this kind of analysis before computing
400 a pangenome. They were removed from the analysis (GCF_900451665.1, GCF_900493335.1
401 and GCF_900493505.1). Finally, these filters left 3980 genomes in the analysis, with an average
402 of 5307 genes per genome, which will be called the reference database DTS1. Among them,
403 there are 225 complete genomes that form the dataset DTS2 (see Figure 6).

404 The functional annotation part is by far the slowest of the first modules. On a typical 5



(a)



(b)

Figure 4: Histograms describing the features of the 5805 *K. pneumoniae* genomes downloaded from Refseq. (a) Distribution of L90 values. (b) Distribution of the number of contigs per genome.

405 Mb genome, giving 2 cores to Prokka, gene finding takes around 40 seconds and functional
406 annotation around 7 minutes.

407 The annotation of the genomes with Prokka 1.11 took approximately 1min 50s per genome,
408 i.e. around 5 days for the whole dataset. For comparison, we re-did the analysis with the option
409 of restricting the analysis to prodigal 2.60. This analysis took less than 6h in total (annotation

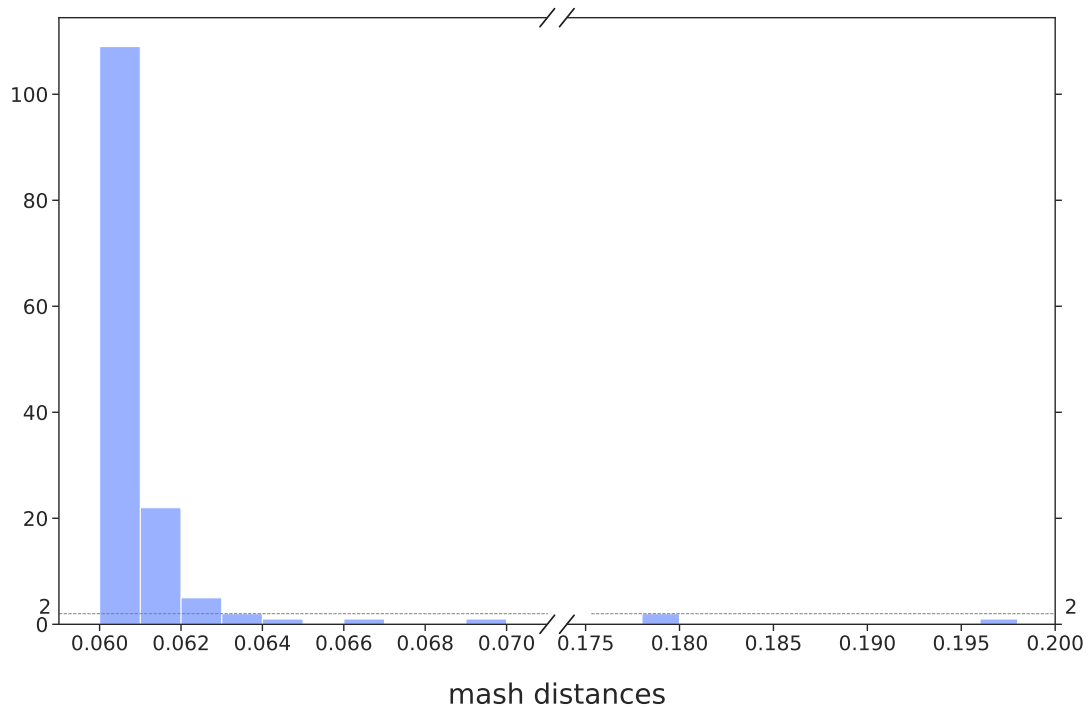


Figure 5: Distribution of Mash distances for the 5572 genomes respecting the L90 and number of contigs thresholds, but having a Mash distance higher than the threshold (0.06).

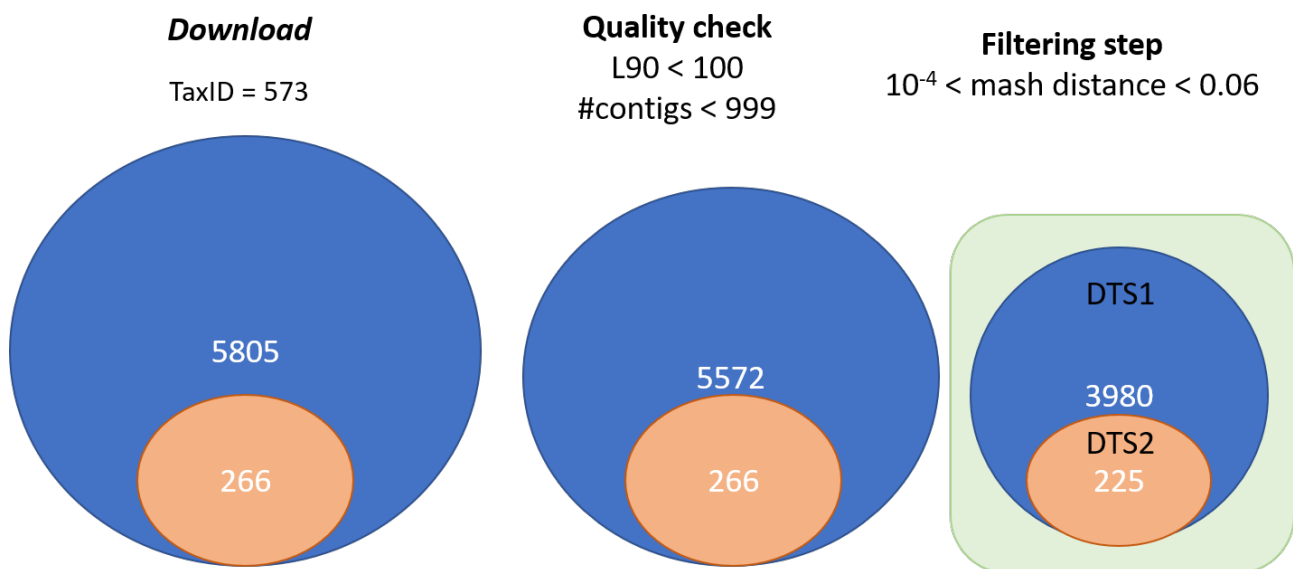


Figure 6: Summary of the procedure to construct DTS1 and DTS2.

410 + formatting of all 3980 genomes), which corresponds to an average of 6 seconds per genome.
 411 In general, this shows that making a simple syntactic annotation leads to a considerable gain
 412 of time. Assuming that genes from the same pangenome family have similar functions, one can
 413 annotate one protein per family at the end of the process and save considerable time.

414 3.2 Building pangenomes

415 We used MMseqs2 Release 11-e1a1c. The 3980 DTS1 genomes contain 20,765,062 proteins.
416 It took less than 30min to create the protein database in the MMseqs2 format, cluster them
417 (with at least 80% identity and 80% coverage of query and target, and other parameters kept as
418 default), and retrieve the pangenome matrices. The pangenome of the smaller DTS2 dataset of
419 225 genomes, 1,190,485 proteins, was computed in less than one minute. The DTS1 pangenome
420 has 86607 families. Among them, 35348 (40%) are singletons (found in a single genome), which
421 is concordant with values observed in *Escherichia coli* [50]. In DTS2 we found 24473 families,
422 including 8975 (37%) singletons.

423 The comparison of these two pangenomes is interesting because it reveals the robustness of
424 the method to changes in sampling size, as summarized in Figure 7. A total of 2147 families
425 contain only members present in both DTS1 and DTS2. This means that, for those families,
426 even with all proteins of DTS1, only proteins from DTS2 were clustered together. Among those
427 families, 2122 are exactly the same in both pangenomes, whereas only 25 families were split
428 in the DTS1 pangenome family relative to the DTS2 pangenome. In that case, they are split,
429 most of the time, into two different families of DTS1. This shows that the clustering procedure
430 is quite robust to the addition of a very large number of genomes.

431 Most important, we observed a total of 22744 families (that is more than 92% of all
432 DTS2 families) that are identical between the independent analysis of the DTS1 and DTS2
433 pangenomes. Identical here means that the DTS2 pangenome gene family is included in a
434 DTS1 pangenome gene family, and the other members of this DTS1 pangenome family are only
435 members of genomes not present in DTS2. Furthermore, around half of the remaining families
436 from the DTS2 pangenome are included in a DTS1 pangenome gene family, which contains a
437 few other proteins from DTS2 genomes. Finally, only 187 gene families of the DTS2 pangenome
438 were split into 2 or 3 different families of DTS1 pangenome. In other words, 24286 families
439 (more than 99%) of DTS2 pangenome are subsets of DTS1 gene families. In conclusion, the
440 construction of pangenome families is robust to large variations in the number of input genomes
441 (see Figure 7).

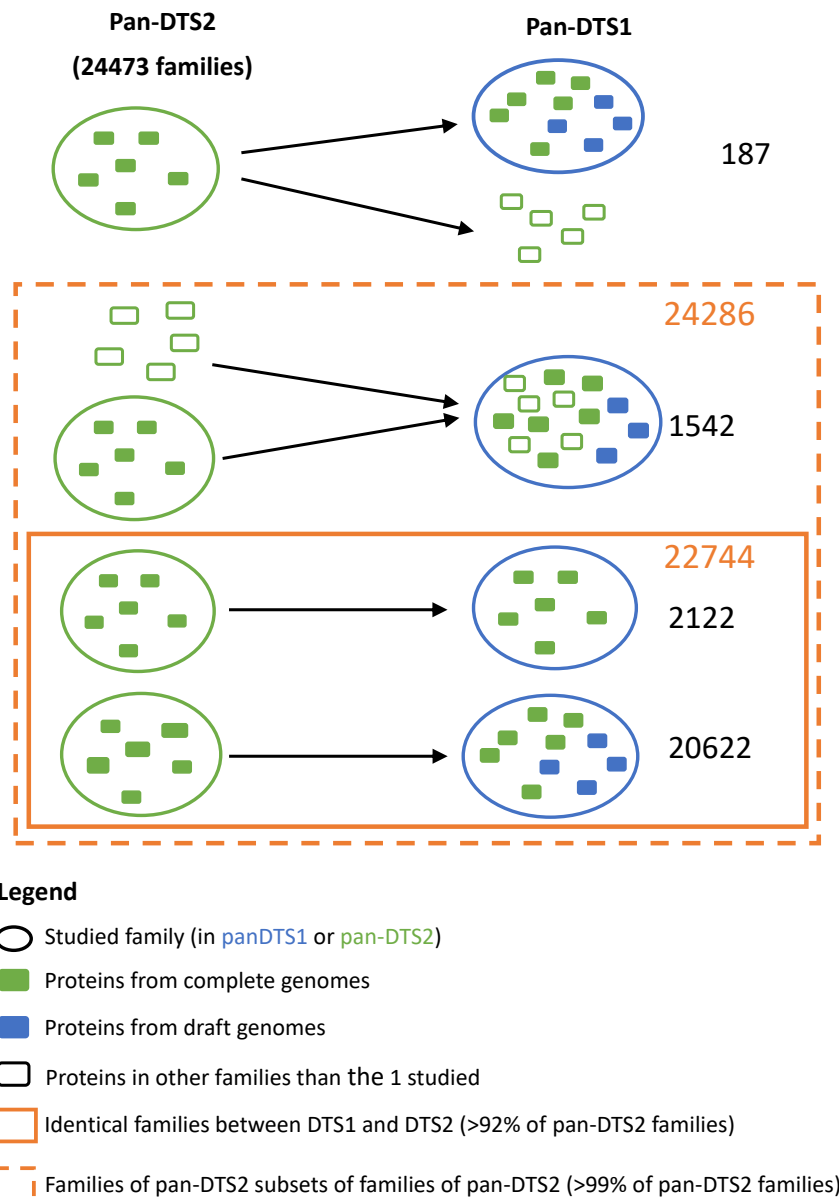


Figure 7: Comparison of the pangenomes generated by PanACoTA for both DTS1 and DTS2.

3.3 Core and persistent genomes

442 This part of the analysis is very fast. Using only 1 core, it took around one minute to generate
 443 a core or persistent genome from DTS1 pangenome. PanACoTA provides a core genome and
 444 three different measures of persistent genome (see Figure 3). The strict-persistent genome
 445 corresponds to cases when the family is present in a single copy in 99% genomes and absent
 446 from the others. In DTS2, the set of complete genomes, the difference between the core and
 447 strict-persistent genome is appreciable (2238 versus 3295 families), i.e. the persistent genome is
 448 50% larger (see Figure 8). The difference becomes huge when the analysis is done on the much
 449

450 larger and less accurate DTS1 dataset, where the two datasets vary by more than one order of
451 magnitude (79 versus 1418 families). In such large datasets of draft genomes the analysis of
452 the core genome is not very useful.

453 The mixed-persistent genome includes the families present in a single copy in 99% genomes
454 and present (potentially in several copies) or absent from the others. It includes the strict-
455 persistent genome and is not much larger than the latter in the small DTS2 dataset. Yet,
456 the difference is much larger in the DTS1 dataset (see Figure 8). While the mixed-persistent
457 genome is 65% percent of the average genome in DTS1, the strict-persistent is only 27% percent
458 in the same dataset. This shows the relevance of using definitions of the core genome adapted
459 to the dataset in order to build robust phylogenetic trees or to analyse patterns of genetic
460 diversification and natural selection.

461 Finally, PanACoTA also computes a multi-persistent genome that includes all gene families
462 present in at least 99% genomes, independently of the copy number. It includes all the other
463 sets and is not much larger than the mixed-persistent genome (see Figure 8). Yet, it includes
464 interesting families. An analysis of these reveals many genes encoding regulators, transporters
465 and enzymes that are nearly ubiquitous, but often present in multiple copies. As a rule, this
466 definition is interesting to study gene families present in most genomes, but present in very
467 different copy number. On the other hand, it is typically not very useful for phylogenetic
468 inference. Since all these sets can be computed very rapidly, it's straightforward to compute
469 them all and use them for different types of analyses.

470 **3.4 Phylogenetic tree inference**

471 PanACoTA ran mafft v.7.467 using `--auto` option to align all families. For DTS1, it selected
472 the FFT-NS-2 method, while for DTS2, it selected FFT-NS-i method. This was done with
473 both the strict-persistent (1418 families) and the mixed-persistent (3441 families). It took 3h
474 (resp. 7h) to align all the families of the strict-persistent (resp. mixed-persistent), giving, for
475 each one, the input file for tree inference.

476 For tree inference, PanACoTA used IQ-TREE multicore version 2.0.6, with `-fast` option.
477 For the tree based on the alignment of the strict-persistent (3980 sequences, 1418 families
478 corresponding to 1,438,179 positions), it took around 7h on 28 cores, requiring 38GB of RAM.

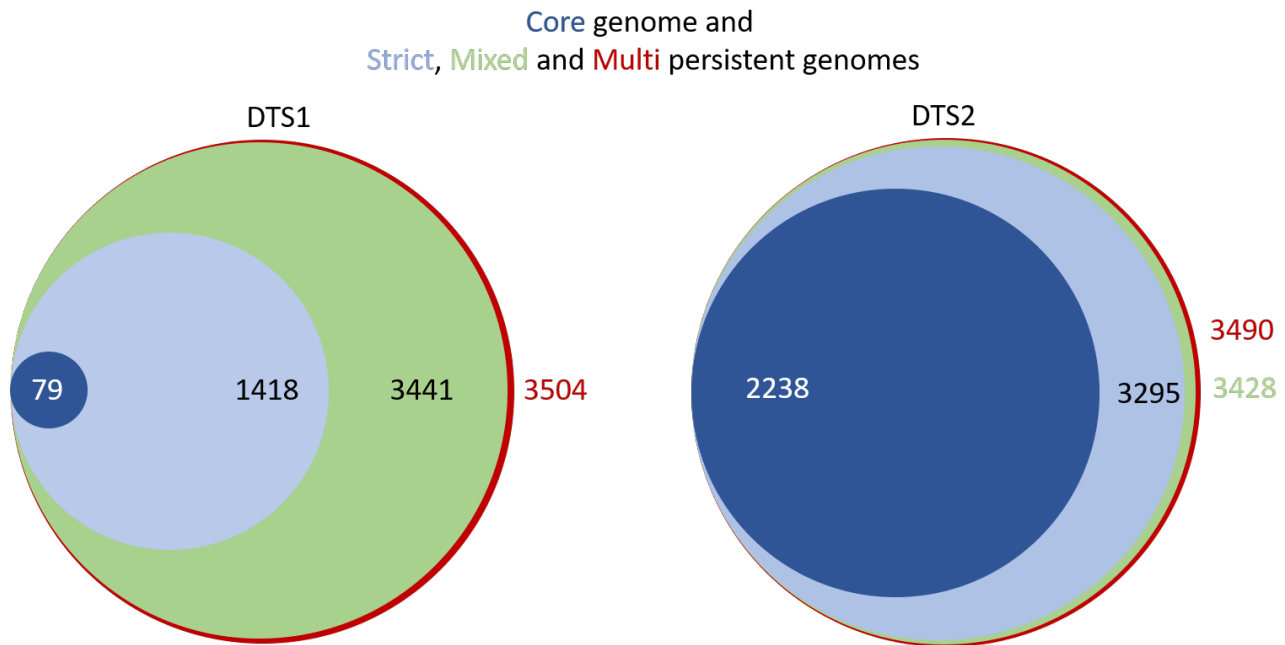


Figure 8: Comparison of the sizes of the core genome and the 3 different types of persistent genomes, for both DTS1 and DTS2. Areas of circles are proportional to the size of the dataset.

479 For the tree based on the alignment of the mixed-persistent (3980 sequences, 3441 families
480 corresponding to 3,393,006 positions), it took 24h using 28 cores, requiring 88GB of RAM.

481 We then wished to understand the differences in phylogenetic inference in terms of the
482 method used to define the persistent genome (strict and mixed persistent). We computed the
483 patristic distance matrix for each tree and a Pearson correlation test showed that they are
484 strongly correlated ($\text{cor} = 0.99138$, $p < 2.2e-16$). This shows that the distances provided by
485 the two methods are very similar. Hence, if the strict persistent is large enough to generate a
486 phylogenetic tree, it provides adequate distances between genomes. Aligning all mixed persistent
487 families would just take much more time, for a similar result. However, if one is interested in
488 having a robust tree topology, one should use the larger (and computationally costlier) dataset.
489 Indeed, the analyses of Robinson-Foulds distance with R `phangorn` package shows a branch-
490 weighted distance of 0.43 and an absolute distance of 2892 [51]. This is because some lineages of
491 *K. pneumoniae* account for a large fraction of the data and these parts of the tree require long
492 informative multiple alignments to produce accurate topologies. Accordingly, the differences
493 in topology between the trees using the DTS2 dataset, which have much larger average branch
494 lengths, show much smaller values of topological distances between the two datasets of persistent
495 genome ($\text{RF}=78$, $\text{wRF}=0.027$).

496 Some researchers use methods to detect recombination in genomes, remove the recombi-

497 nation tracts, and then redo the analyses. This can be done outside PanACoTA by querying
498 the multiple alignments before proceeding to the phylogenetic inference. Yet, previous results
499 have shown that such procedures tend to distort phylogenetic inference at a larger extent than
500 simply using all the information in the multiple alignments [52] [53]), and this explains why we
501 have not included such an option on PanACoTA.

502 4 Conclusion

503 PanACoTA is a pipeline for those wanting to test hypotheses or explore genomic patterns using
504 large scale comparative genomics. We hope that it will be particularly useful for those wishing to
505 use a rapid, accurate and standardized procedure to obtain the basic building blocks of typical
506 analyses of genetic variation at the species level. We built the pipeline having modularity in
507 mind, so that users can produce multiple variants of the analyses at each stage. We also paid
508 particularly care with the portability and evolvability of the software. These two characteristics,
509 modularity and evolvability, will facilitate the implementation of novel procedures in the future.

510 5 Acknowledgements

511 We thank Marie TOUCHON, Matthieu HAUDIQUET and Rémi DENISE for comments,
512 suggestions, and bug reports, and Blaise LI for his help with Singularity. This work was
513 partly supported by the ANR Salmo_Prophages (ANR-16-CE16-0029), the Inception program
514 (PIA/ANR-16-CONV-0005), and the Equipe FRM (EQU201903007835).

515 References

- 516 [1] G. Vernikos, D. Medini, D. R. Riley, and H. Tettelin, “Ten years of pan-genome analyses,”
517 *Current Opinion in Microbiology*, vol. 23, pp. 148–154, feb 2015.
- 518 [2] H. Tettelin and D. Medini, *The pangenome: Diversity, dynamics and evolution of genomes*.
519 Springer International Publishing, jun 2020.
- 520 [3] M. V. Larsen, S. Cosentino, S. Rasmussen, C. Friis, H. Hasman, R. L. Marvig, L. Jels-
521 bak, T. Sicheritz-Pontén, D. W. Ussery, F. M. Aarestrup, and O. Lund, “Multilocus

- 522 sequence typing of total-genome-sequenced bacteria,” *Journal of Clinical Microbiology*,
523 vol. 50, pp. 1355–1361, apr 2012.
- 524 [4] S. Baker, N. Thomson, F. X. Weill, and K. E. Holt, “Genomic insights into the emergence
525 and spread of antimicrobial-resistant bacterial pathogens,” *Science*, vol. 360, pp. 733–738,
526 may 2018.
- 527 [5] T. J. Treangen and M. Pop, “You can’t always sequence your way out of a tight spot,”
528 *EMBO reports*, vol. 19, p. e47036, dec 2018.
- 529 [6] S. K. Sheppard, D. S. Guttman, and J. R. Fitzgerald, “Population genomics of bacterial
530 host adaptation,” *Nature Reviews Genetics*, vol. 19, pp. 549–565, sep 2018.
- 531 [7] D. Falush, “Bacterial genomics: Microbial GWAS coming of age,” *Nature Microbiology*,
532 vol. 1, p. 16059, apr 2016.
- 533 [8] L. B. Zeng, D. Wang, N. Y. Hu, Q. Zhu, K. Chen, K. Dong, Y. Zhang, Y. F. Yao, X. K. Guo,
534 Y. F. Chang, and Y. Z. Zhu, “A novel pan-genome reverse vaccinology approach employing
535 a negative-selection strategy for screening surface-exposed antigens against leptospirosis,”
536 *Frontiers in Microbiology*, vol. 8, p. 396, mar 2017.
- 537 [9] B. Buchfink, C. Xie, and D. H. Huson, “Fast and sensitive protein alignment using DIA-
538 MOND,” *Nature Methods*, vol. 12, pp. 59–60, jan 2014.
- 539 [10] K. Katoh and D. M. Standley, “MAFFT Multiple Sequence Alignment Software Version
540 7: Improvements in performance and usability,” *Molecular Biology and Evolution*, vol. 30,
541 pp. 772–780, apr 2013.
- 542 [11] P. Bradley, H. C. den Bakker, E. P. Rocha, G. McVean, and Z. Iqbal, “Ultrafast search of
543 all deposited bacterial and viral genomic data,” *Nature Biotechnology*, vol. 37, pp. 152–159,
544 feb 2019.
- 545 [12] Y. Kim, C. Gu, H. U. Kim, and S. Y. Lee, “Current status of pan-genome analysis for
546 pathogenic bacteria,” *Current Opinion in Biotechnology*, vol. 63, pp. 54–62, jun 2020.
- 547 [13] R. C. Edgar, “Search and clustering orders of magnitude faster than BLAST,” *Bioinfor-
548 matics*, vol. 26, pp. 2460–2461, aug 2010.

- 549 [14] L. Fu, B. Niu, Z. Zhu, S. Wu, and W. Li, “CD-HIT: Accelerated for clustering the next-
550 generation sequencing data,” *Bioinformatics*, vol. 28, pp. 3150–3152, dec 2012.
- 551 [15] A. J. Page, C. A. Cummins, M. Hunt, V. K. Wong, S. Reuter, M. T. Holden, M. Fookes,
552 D. Falush, J. A. Keane, and J. Parkhill, “Roary: Rapid large-scale prokaryote pan genome
553 analysis,” *Bioinformatics*, vol. 31, pp. 3691–3693, may 2015.
- 554 [16] G. Tonkin-Hill, N. MacAlasdair, C. Ruis, A. Weimann, G. Horesh, J. A. Lees, R. A.
555 Gladstone, S. Lo, C. Beaudoin, R. A. Floto, S. D. Frost, J. Corander, S. D. Bentley,
556 and J. Parkhill, “Producing polished prokaryotic pangenomes with the Panaroo pipeline,”
557 *Genome biology*, vol. 21, p. 180, jul 2020.
- 558 [17] N. M. Chaudhari, V. K. Gupta, and C. Dutta, “BPGA-an ultra-fast pan-genome analysis
559 pipeline,” *Scientific Reports*, vol. 6, pp. 1–10, apr 2016.
- 560 [18] W. Ding, F. Baumdicker, and R. A. Neher, “panX: pan-genome analysis and exploration,”
561 *Nucleic acids research*, vol. 46, p. e5, jan 2018.
- 562 [19] S. Cosentino and W. Iwasaki, “SonicParanoid: Fast, accurate and easy orthology infer-
563 ence,” *Bioinformatics*, vol. 35, pp. 149–151, jan 2019.
- 564 [20] M. Steinegger and J. Söding, “MMseqs2 enables sensitive protein sequence searching for
565 the analysis of massive data sets,” *Nature Biotechnology*, vol. 35, pp. 1026–1028, nov 2017.
- 566 [21] G. Gautreau, A. Bazin, M. Gachet, R. Planel, L. Burlot, M. Dubois, A. Perrin, C. Médigue,
567 A. Calteau, S. Cruveiller, C. Matias, C. Ambroise, E. P. Rocha, and D. Vallenet, “PPanG-
568 GOLiN: Depicting microbial diversity via a partitioned pangenome graph,” *PLoS Computa-
569 tional Biology*, vol. 16, p. e1007732, mar 2020.
- 570 [22] S. C. Bayliss, H. A. Thorpe, N. M. Coyle, S. K. Sheppard, and E. J. Feil, “PIRATE:
571 A fast and scalable pangenomics toolbox for clustering diverged orthologues in bacteria,”
572 *GigaScience*, vol. 8, pp. 1–9, oct 2019.
- 573 [23] K. E. Holt, H. Wertheim, R. N. Zadoks, S. Baker, C. A. Whitehouse, D. Dance, A. Jenney,
574 T. R. Connor, L. Y. Hsu, J. Severin, S. Brisse, H. Cao, J. Wilksch, C. Gorrie, M. B.
575 Schultz, D. J. Edwards, K. Van Nguyen, T. V. Nguyen, T. T. Dao, M. Mensink, V. Le

- 576 Minh, N. T. K. Nhu, C. Schultsz, K. Kuntaman, P. N. Newton, C. E. Moore, R. A.
577 Strugnell, and N. R. Thomson, “Genomic analysis of diversity, population structure, vir-
578 ulence, and antimicrobial resistance in *Klebsiella pneumoniae*, an urgent threat to public
579 health,” *Proceedings of the National Academy of Sciences of the United States of America*,
580 vol. 112, pp. E3574–E3581, jul 2015.
- 581 [24] B. D. Ondov, T. J. Treangen, P. Melsted, A. B. Mallonee, N. H. Bergman, S. Koren,
582 and A. M. Phillippy, “Mash: fast genome and metagenome distance estimation using
583 MinHash.,” *Genome biology*, vol. 17, no. 1, p. 132, 2016.
- 584 [25] A. Z. Broder, “On the resemblance and containment of documents,” in *Proceedings of the*
585 *International Conference on Compression and Complexity of Sequences*, pp. 21–29, 1997.
- 586 [26] S. F. Altschul, W. Gish, W. Miller, E. W. Myers, and D. J. Lipman, “Basic local alignment
587 search tool,” *Journal of Molecular Biology*, vol. 215, no. 3, pp. 403–410, 1990.
- 588 [27] C. Jain, L. M. Rodriguez-R, A. M. Phillippy, K. T. Konstantinidis, and S. Aluru, “High
589 throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries,”
590 *Nature Communications*, vol. 9, pp. 1–8, nov 2018.
- 591 [28] K. T. Konstantinidis and J. M. Tiedje, “Genomic insights that advance the species def-
592 inition for prokaryotes,” *Proceedings of the National Academy of Sciences of the United*
593 *States of America*, vol. 102, pp. 2567–2572, feb 2005.
- 594 [29] T. Seemann, “Prokka: rapid prokaryotic genome annotation.,” *Bioinformatics*, vol. 30,
595 pp. 2068–9, jul 2014.
- 596 [30] D. Hyatt, G. L. Chen, P. F. LoCascio, M. L. Land, F. W. Larimer, and L. J. Hauser,
597 “Prodigal: Prokaryotic gene recognition and translation initiation site identification,” *BMC*
598 *Bioinformatics*, vol. 11, p. 119, mar 2010.
- 599 [31] C. Camacho, G. Coulouris, V. Avagyan, N. Ma, J. Papadopoulos, K. Bealer, and T. L.
600 Madden, “BLAST+: Architecture and applications,” *BMC Bioinformatics*, vol. 10, p. 421,
601 dec 2009.

- 602 [32] S. R. Eddy, “Accelerated profile HMM searches,” *PLoS Computational Biology*, vol. 7,
603 p. e1002195, oct 2011.
- 604 [33] J. D. Selengut, D. H. Haft, T. Davidsen, A. Ganapathy, M. Gwinn-Giglio, W. C. Nelson,
605 A. R. Richter, and O. White, “TIGRFAMs and Genome Properties: Tools for the assign-
606 ment of molecular function and biological process in prokaryotic genomes,” *Nucleic Acids*
607 *Research*, vol. 35, p. D260, jan 2007.
- 608 [34] S. El-Gebali, J. Mistry, A. Bateman, S. R. Eddy, A. Luciani, S. C. Potter, M. Qureshi,
609 L. J. Richardson, G. A. Salazar, A. Smart, E. L. Sonnhammer, L. Hirsh, L. Paladin,
610 D. Piovesan, S. C. Tosatto, and R. D. Finn, “The Pfam protein families database in 2019,”
611 *Nucleic Acids Research*, vol. 47, pp. D427–D432, jan 2019.
- 612 [35] P. Jones, D. Binns, H. Y. Chang, M. Fraser, W. Li, C. McAnulla, H. McWilliam, J. Maslen,
613 A. Mitchell, G. Nuka, S. Pesseat, A. F. Quinn, A. Sangrador-Vegas, M. Scheremetjew,
614 S. Y. Yong, R. Lopez, and S. Hunter, “InterProScan 5: Genome-scale protein function
615 classification,” *Bioinformatics*, vol. 30, pp. 1236–1240, may 2014.
- 616 [36] M. Steinegger and J. Söding, “Clustering huge protein sequence sets in linear time,” *Nature*
617 *Communications*, vol. 9, pp. 1–8, dec 2018.
- 618 [37] D. E. Fouts, L. Brinkac, E. Beck, J. Inman, and G. Sutton, “PanOCT: Automated cluster-
619 ing of orthologs using conserved gene neighborhood for pan-genomic analysis of bacterial
620 strains and closely related species,” *Nucleic Acids Research*, vol. 40, p. e172, dec 2012.
- 621 [38] J. M. Inman, G. G. Sutton, E. Beck, L. M. Brinkac, T. H. Clarke, and D. E. Fouts,
622 “Large-scale comparative analysis of microbial pan-genomes using PanOCT,” *Bioinfor-*
623 *matics*, vol. 35, pp. 1049–1050, mar 2019.
- 624 [39] C. H. Georgescu, A. L. Manson, A. D. Griggs, C. A. Desjardins, A. Pironti, I. Wapinski,
625 T. Abeel, B. J. Haas, and A. M. Earl, “SynerClust: a highly scalable, synteny-aware
626 orthologue clustering tool,” *Microbial genomics*, vol. 4, p. e000231, nov 2018.
- 627 [40] K. Abudahab, J. M. Prada, Z. Yang, S. D. Bentley, N. J. Croucher, J. Corander, and
628 D. M. Aanensen, “PANINI: Pangenome neighbour identification for bacterial populations,”
629 *Microbial Genomics*, vol. 5, p. e000220, apr 2019.

- 630 [41] C. Collins and X. Didelot, “A phylogenetic method to perform genome-wide association
631 studies in microbes that accounts for population structure and recombination,” *PLoS Com-
632 putational Biology*, vol. 14, p. e1005958, feb 2018.
- 633 [42] A. Filipinski, O. Murillo, A. Freydenzon, K. Tamura, and S. Kumar, “Prospects for build-
634 ing large timetrees using molecular data with incomplete gene coverage among species,”
635 *Molecular Biology and Evolution*, vol. 31, pp. 2542–2550, sep 2014.
- 636 [43] J. D. Thompson, B. Linard, O. Lecompte, and O. Poch, “A comprehensive benchmark
637 study of multiple sequence alignment methods: Current challenges and future perspec-
638 tives,” *PLoS ONE*, vol. 6, p. e18093, mar 2011.
- 639 [44] L.-T. Nguyen, H. A. Schmidt, A. von Haeseler, and B. Q. Minh, “IQ-TREE: A fast and
640 effective stochastic algorithm for estimating maximum-likelihood phylogenies,” *Molecular
641 Biology and Evolution*, vol. 32, pp. 268–274, jan 2015.
- 642 [45] M. N. Price, P. S. Dehal, and A. P. Arkin, “FastTree 2 – Approximately Maximum-
643 Likelihood Trees for Large Alignments,” *PLoS ONE*, vol. 5, p. e9490, mar 2010.
- 644 [46] V. Lefort, R. Desper, and O. Gascuel, “FastME 2.0: A Comprehensive, Accurate, and Fast
645 Distance-Based Phylogeny Inference Program,” *Molecular Biology and Evolution*, vol. 32,
646 pp. 2798–2800, oct 2015.
- 647 [47] K. Howe, A. Bateman, and R. Durbin, “QuickTree: Building huge neighbour-joining trees
648 of protein sequences,” *Bioinformatics*, vol. 18, pp. 1546–1547, nov 2002.
- 649 [48] M. List, P. Ebert, and F. Albrecht, “Ten Simple Rules for Developing Usable Software in
650 Computational Biology,” *PLoS Computational Biology*, vol. 13, p. e1005265, jan 2017.
- 651 [49] G. Wilson, D. A. Aruliah, C. T. Brown, N. P. Chue Hong, M. Davis, R. T. Guy, S. H.
652 Haddock, K. D. Huff, I. M. Mitchell, M. D. Plumbley, B. Waugh, E. P. White, and P. Wil-
653 son, “Best Practices for Scientific Computing,” *PLoS Biology*, vol. 12, no. 1, p. e1001745,
654 2014.

- 655 [50] M. Touchon, A. Perrin, J. A. M. De Sousa, B. Vangchhia, S. Burn, C. L. O'Brien, E. Dena-
656 mur, D. Gordon, and E. P. Rocha, "Phylogenetic background and habitat drive the genetic
657 diversification of *Escherichia coli*," *PLoS Genetics*, vol. 16, p. e1008866, jun 2020.
- 658 [51] K. P. Schliep, "phangorn: Phylogenetic analysis in R," *Bioinformatics*, vol. 27, pp. 592–593,
659 feb 2011.
- 660 [52] J. Hedge and D. J. Wilson, "Bacterial phylogenetic reconstruction from whole genomes is
661 robust to recombination but demographic inference is not," *mBio*, vol. 5, pp. e02158–14,
662 nov 2014.
- 663 [53] M. Lapierre, C. Blin, A. Lambert, G. Achaz, and E. P. Rocha, "The Impact of Selection,
664 Gene Conversion, and Biased Sampling on the Assessment of Microbial Demography,"
665 *Molecular biology and evolution*, vol. 33, pp. 1711–1725, jul 2016.