# Large scale automated phylogenomic analysis of bacterial whole-genome isolates and the Evergreen platform — **Source link** ⧉

Judit Szarvas, Johanne Ahrenfeldt, José Cisneros, Martin Christen Frølund Thomsen ...+2 more authors

**Institutions:** Technical University of Denmark

**Published on:** 22 Jun 2019 - bioRxiv (Cold Spring Harbor Laboratory)

**Topics:** Reference genome

Related papers:

- Large scale automated phylogenomic analysis of bacterial isolates and the Evergreen Online platform.

- Multilocus Sequence Typing of Total-Genome-Sequenced Bacteria

- Real-Time Pathogen Detection in the Era of Whole-Genome Sequencing and Big Data: Comparison of k-mer and Site-Based Methods for Inferring the Genetic Distances among Tens of Thousands of Salmonella Samples

- Whole genome single-nucleotide variation profile-based phylogenetic tree building methods for analysis of viral, bacterial and human genomes

- From raw reads to trees: Whole genome SNP phylogenetics across the tree of life

1  Large scale automated phylogenomic analysis of bacterial whole-genome isolates and the

2  Evergreen platform

3

4  Judit Szarvas[1] & Johanne Ahrenfeldt[1], Jose Luis Bellod Cisneros[1], Martin Christen Frølund Thomsen[1], Frank M.

5  Aarestrup[1], Ole Lund[1]

6  [1]Research Group for Genomic Epidemiology, National Food Institute, Technical University of Denmark, Kongens

7  Lyngby, Denmark.

8

9  Abstract

10  Public health authorities whole-genome sequence thousands of pathogenic isolates each month for microbial

11  diagnostics and surveillance of pathogenic bacteria. The computational methods have not kept up with the

12  deluge of data and need for real-time results.

13  We have therefore created a bioinformatics pipeline for rapid subtyping and continuous phylogenomic analysis

14  of bacterial samples, suited for large-scale surveillance. To decrease the computational burden, a two-level

15  clustering strategy is employed. The data is first divided into sets by matching each isolate to a closely related

16  reference genome. The reads then are aligned to the reference to gain a consensus sequence and SNP based

17  genetic distance is calculated between the sequences in each set. Isolates are clustered together with a

18  threshold of 10 SNPs. Finally, phylogenetic trees are inferred from the non-redundant sequences and the

19  clustered isolates are placed on a clade with the cluster representative sequence. The method was

20  benchmarked and found to be accurate in grouping outbreak strains together, while discriminating from non-

21  outbreak strains.

22  The pipeline was applied in Evergreen Online, which processes publicly available sequencing data from

23  foodborne bacterial pathogens on a daily basis, updating the phylogenetic trees as needed. It has so far placed

24  more than 100,000 isolates into phylogenies, and has been able to keep up with the daily release of data. The

25  trees are continuously published on https://cge.cbs.dtu.dk/services/Evergreen

26

27  Keywords

28  Phylogenomics, WGS, subtyping, SNP, automation, epidemiology, outbreak investigation

29

30

31

## Main

Epidemiological typing of bacteria is used by hospitals and public health authorities, as well as animal health authorities, to detect outbreaks of infectious diseases and determine trends over time. Traditionally, that includes culturing and isolating the pathogen, followed by species identification and subtyping using various conventional microbiological and molecular methodologies.

For outbreak investigation, it is necessary to place the infectious agent into a more discriminatory category than species, to establish links between cases and sources. Multi-locus sequence typing (MLST) has been a frequently used molecular subtyping method, where sequence types are assigned to the isolates based on the combinations of alleles for 6-10 housekeeping genes[1].

Whole-genome sequencing (WGS) has opened a new chapter in microbial diagnostics and epidemiological typing. WGS data can be used to determine, amongst other properties, both MLST types and serotype of several bacterial species[2,3]. Several studies for multiple bacterial species have shown the value of WGS for elucidating the bacterial evolution and phylogeny, and identifying outbreaks[4–6].

The use of WGS has enabled the unbiased comparison of samples processed in different laboratories, boosting surveillance and outbreak detection, but the methods for sharing and comparing a large number of samples have not been established yet[7,8]. Therefore, a number of national, regional and international initiatives have been launched with the aim of facilitating the sharing, analyses and comparison of WGS data[9–11].

Since 2012, the US Food and Drug Administration (FDA) is leading a network of public health and university laboratories, called GenomeTrakr. These laboratories sequence bacterial isolates from food and environmental samples and upload the data to the National Center for Biotechnology Information (NCBI). GenomeTrakr is restricted to foodborne pathogens and currently includes data from seven such bacterial species.[12] All raw WGS data are publicly shared through NCBI, facilitating the collaboration between laboratories. Furthermore, the raw data are picked up by the NCBI Pathogen Detection pipeline[13], that assembles the samples into draft genomes to predict the nearest neighbors and construct phylogenetic trees for each within-50-SNPs cluster using an exact maximum compatibility algorithm[14]. This approach requires access to all of the raw data or assembled genomes, and very extensive computational resources for larger databases, like *Salmonella enterica*. In addition, no sub-species taxonomical classification has so far been implemented in the pipeline.

Focusing on the same bacterial species as GenomeTrakr, PulseNet USA has also established procedures for use of WGS data for outbreak detection. In their vision, an extension of the highly successful MLST approach into a core-genome (cgMLST) or whole-genome (wgMLST) scheme, with genes in the order of thousand, would allow for sharing information under a common nomenclature. Meanwhile, all of the raw data could be kept locally. Only data from individual strains would have to be shared when further confirmation of an outbreak is required.[11] MLST schemes are offered from several databases[15–17], and a number of, at times conflicting, cg- and wgMLST schemes have recently been proposed for a limited number of bacterial species[16,18–24]. Moreover, few of the proposed schemes provide a definitive nomenclature of sequence types to go with the allele profiles. The existing schemes do not cover all of the potential allelic variation: a recent study showed, that for *Campylobacter jejuni*, that has maintained MLST schemes, only approximately 53% of the strains of animal origin could be assigned to an existing unique allelic profile[25]. Continuous curation of the hundreds of relevant bacterial species, that are known human, animal and plant pathogens, would require great effort. A centralized database for the distribution of the allele profiles and sequences would be also necessary. Furthermore, for comparable results, and surveillance, the same analysis pipeline or software should be used for the prediction of the allelic profiles. Single-linkage cgMLST clusters can be generated of public and private uploaded data on
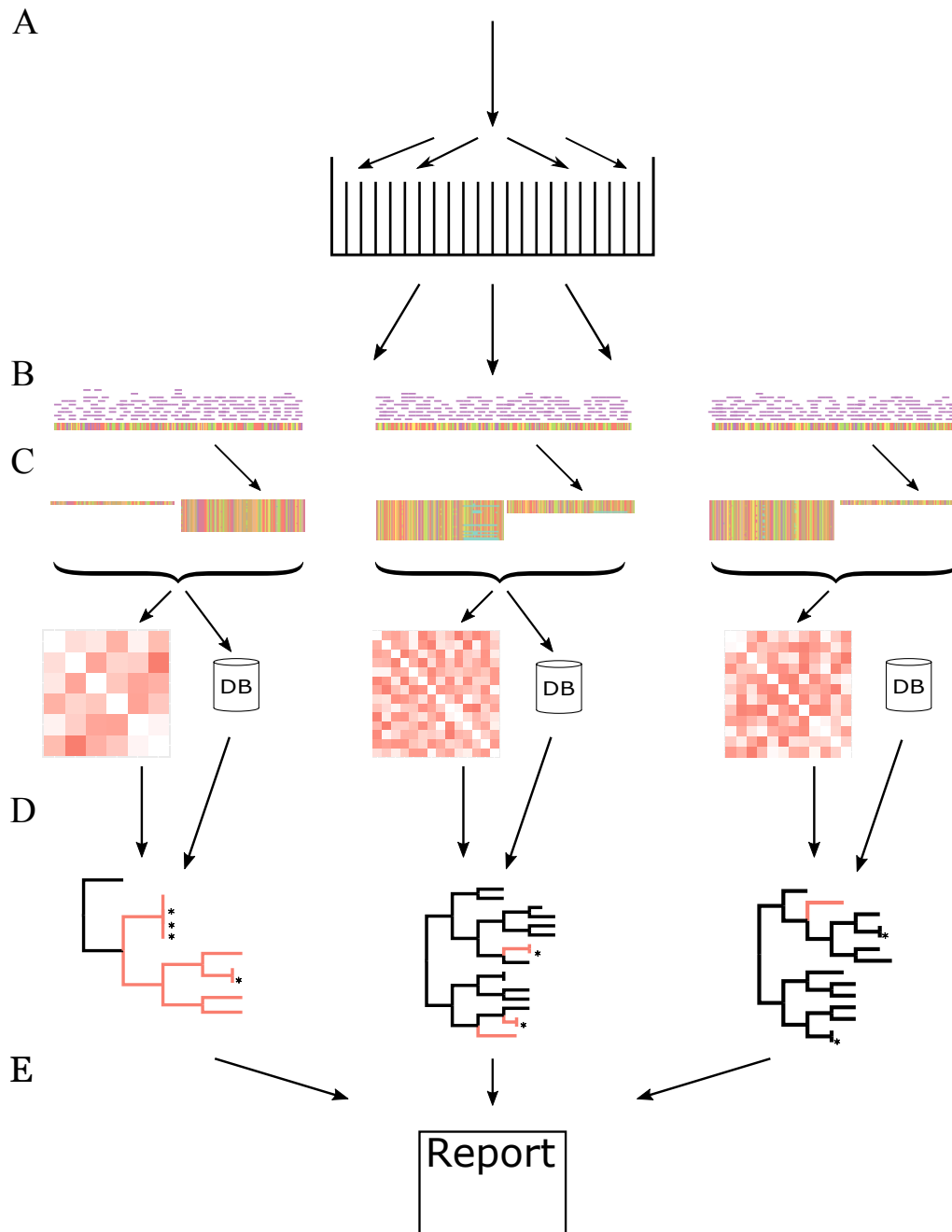
74    EnteroBase[15], and up to 1000 sequences on Pathogenwatch[26], by manual selection of strains to be included in
75    the analysis.

76    The approaches mentioned above yield preliminary results for outbreak detection, as they often lack the
77    necessary resolution, thus, in most cases, selected WGS data are further analyzed using single nucleotide
78    profiling. Here, genomic variants (single nucleotide polymorphisms (SNPs), insertions and deletions) are
79    derived by aligning WGS reads to a reference genome. For each bacterial species, custom single nucleotide
80    profiling (SNP validation, cluster threshold determination, etc.) is necessary in order to achieve results that are
81    biologically relevant and informative. The samples (of current interest and historical) included in the analysis
82    and the reference genome are chosen on a case-by-case basis, usually based on subtyping results. EnteroBase
83    offers SNP analysis of user selected strains based on the genotypes, but Alikhan et al. dismiss the feasibility of
84    large-scale SNP analyses[15]. Various SNP analysis pipelines are used by laboratories and research groups for
85    inferring phylogenetic trees for isolates of interest[27–32]. For example, Public Health England developed and uses
86    SnapperDB for outbreak detection without initial cluster analysis by cg- or wg-MLST. SnapperDB consists of
87    tools to create a database of SNPs compared to a given reference sequence, and assign each isolate a SNP
88    address based on single linkage clustering.[33]

89    We present here a whole-genome, single nucleotide-based method for subtyping and preliminary
90    phylogenomic analysis, that circumvent the known limitations of current gene- and SNP-based approaches.
91    PAPABAC carries out rapid and automated subtyping of bacterial whole-genome sequenced isolates and
92    generates continuously updated phylogenetic trees based on nucleotide differences. We demonstrate two
93    applications, a standalone version for local monitoring of bacterial isolates, and Evergreen Online, for global
94    surveillance of foodborne bacterial pathogens. We also suggest a stable naming scheme for each isolate,
95    making the results from the pipeline easier to communicate to others. To the best of our knowledge, no such
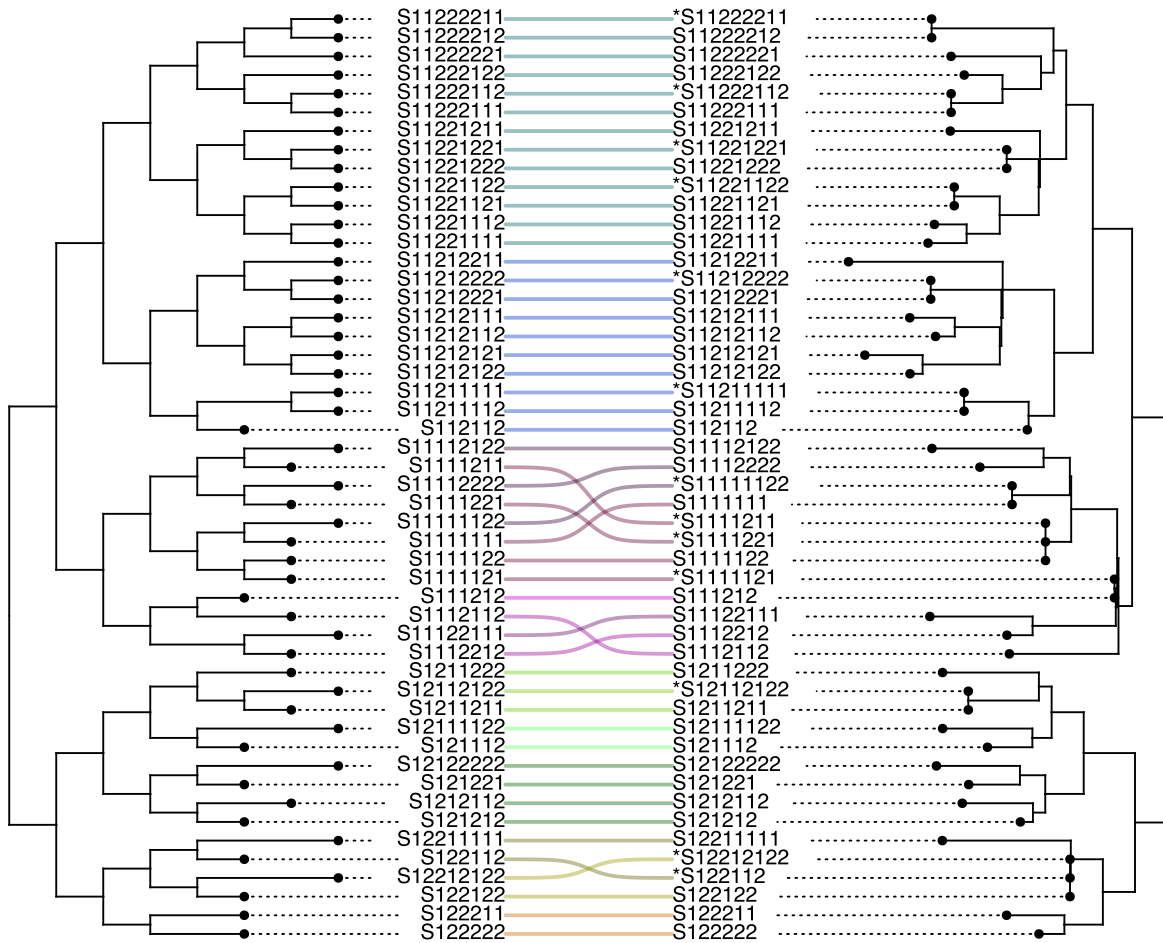96    tool exists at the moment.

97

3

98



99

*Figure 1 Overview of PAPABAC. (A) The input raw read files are classified into sets based on k-mer similarity to NCBI RefSeq complete prokaryotic chromosomal genomes. (B) The raw reads are mapped to the reference genome and a consensus sequence is generated via strict statistical evaluation ($p < 0.05$) of the mapped bases in each position. (C) The resulting consensus sequences are of equal length in each template set. The new isolates in each set are clustered to the non-redundant isolates already in the set if the pairwise nucleotide difference based genetic distance is less than 10. The remaining new isolates undergo the same clustering process. (D) Pairwise genetic distance between all non-redundant isolate in the set is used as input for neighbor-joining algorithm. If there are less than 600 non-redundant isolates in a set, an approximately maximum likelihood phylogenetic tree is also inferred based on the consensus sequences (red: new isolates). The clustered isolates are placed back onto the trees with 0 distance to the cluster representative (marked with an asterisk). (E) The information about the acquired isolates, the sets, the clusters, and the phylogenetic trees is stored in SQLite databases, which are queried once all sets with new isolates are processed to output the results to the users.*

111

Figure 2 Comparison of the ideal tree (left) to the PAPABAC maximum likelihood tree made of the in vitro experiment dataset[34]Taxa with an asterisk were clustered together with the taxa in the same clade.

114 Results

115 **Pipeline for automated phylogenomic analysis of bacterial whole-genome sequences (PAPABAC)**

116 We developed PAPABAC (Figure 1), a phylogenomic pipeline for the automated analysis of bacterial isolates,
117 that needs no additional input besides WGS data (fastq files) and generates clusters of closely related isolates.
118 PAPABAC first matches the isolates to complete bacterial chromosomal genome reference sequences with
119 greater than 99.0% sequence identity and a minimum average depth of 11. These reference sequences serve as
120 templates for the alignment of the raw reads. The aligned bases at each position are statistically evaluated to
121 determine the consensus sequence, as previously described for a nucleotide difference method[35]. Positions
122 that do not fulfil the significance criteria remain ambiguous, get assigned "N", and are disregarded during the
123 pairwise genetic distance calculation. These steps ensure that there is high confidence in the consensus
124 sequence that is the basis of the genetic distance estimation.

125 The pipeline retains analysis results in such a manner that input is added to the previously processed data. The
126 phylogenomic analysis is carried out on the current input and the previously found non-redundant isolates
127 (singletons and cluster representatives). The genetic distance is estimated in a pairwise manner, comparing the
128 given two sequences for all non-ambiguous positions, i.e. positions where none of the two sequences have an
129 "N" assigned. The distances between the previously processed runs are stored on disk, saving computational
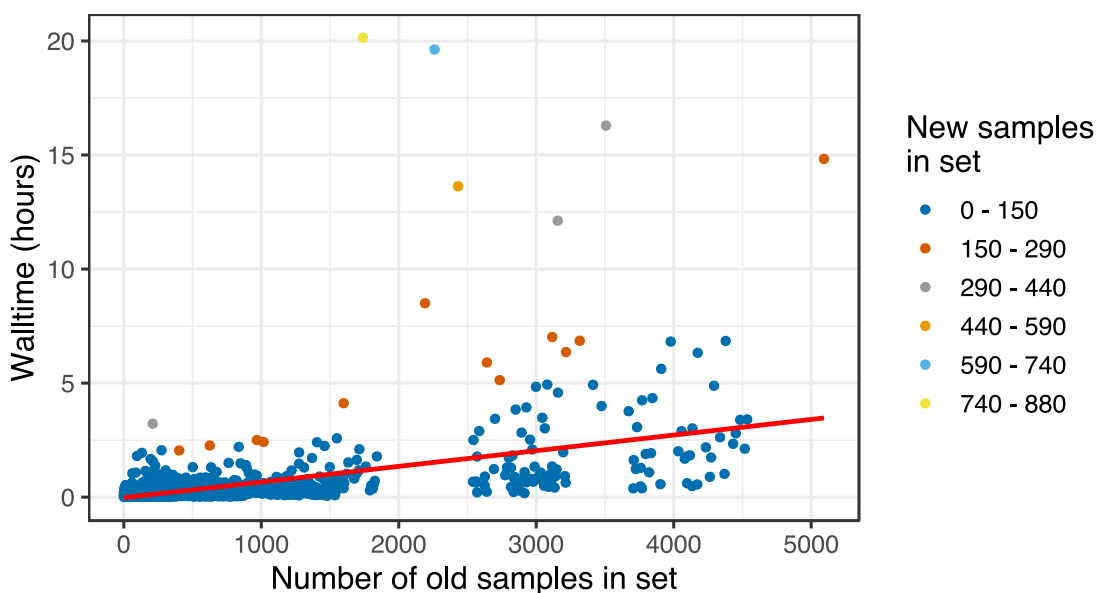130 time, and only the distances to the new isolates are computed in a given run.

131 A clustering step during the genetic distance calculation forms clusters of closely related isolates and reduces
132 the number of similar sequences in each set, and thereby also reduce the computation time. After identifying a
133 non-redundant isolate and a closely related isolate to it, the one previously deemed non-redundant will be the
134 cluster representative and kept, while the clustered one is omitted from the subsequent runs of the pipeline.
135 However, the information about the clustering is added to a database and the clustered isolate will be placed
136 on the inferred phylogenetic tree. The cluster representatives remain constant through the subsequent runs of
137 the pipeline, and the clusters only increase in size if new isolates are clustered with the representative.
138 Therefore, each cluster is stable in the sense that an isolate will newer change which representative it is
139 associated with and each cluster can be reliably identified by the template name and the identifier of its cluster
140 representative.

141 The pipeline can be run on a computer with 8 Gb RAM and Unix system. The computational time is reduced
142 compared to re-running the whole analysis each time new samples are added, even without parallelisation
143 (Figure S1).

144 PAPABAC was benchmarked against three SNP pipeline benchmarking datasets. An *Escherichia coli in vitro*
145 evolution experiment dataset[34] provided 50 closely related samples on a short temporal scale with less than
146 100 nucleotide differences across the dataset. The algorithm clustered together 7 out of 10 samples with the
147 same ancestor that were taken on the same day and presumably had less than 10 nucleotide differences
148 between them. The PAPABAC maximum likelihood (Figure 2) and neighbor-joining (Figure S2) trees with the
149 clustered isolates pruned to resolve the polytomies were comparable to the ideal phylogeny of the *in vitro*
150 experiment dataset: the normalized Robinson-Foulds distances were 0.18 and 0.12, respectfully.
151 Benchmarking against the *Campylobacter jejuni* (Figure S3A) and the *Listeria monocytogenes* (Figure S3B)
152 datasets from Timme et al.[36], PAPABAC correctly clustered the related outbreak strains (colored) and the
153 outgroups, where the genetic distance was below the clustering threshold. The topologies of the maximum
154 likelihood phylogenetic trees closely resembled the tree topologies given.

155

6

Figure 3 Time requirement of the phylogenomic analysis for given number of non-redundant and new strains, on 20 CPUs.
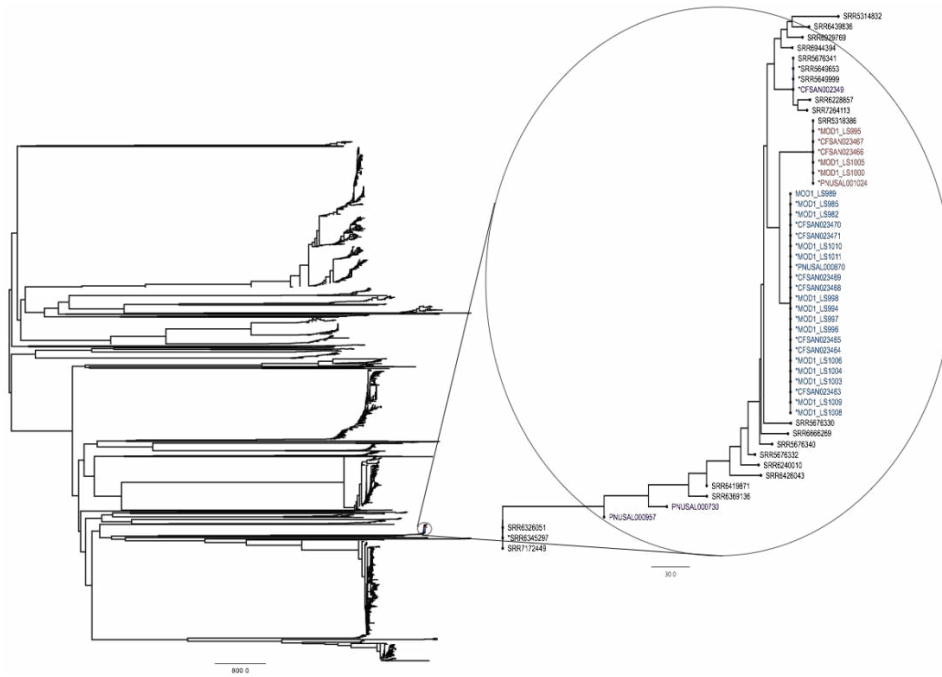
**Evergreen Online for surveillance of foodborne bacterial pathogens**

Evergreen Online was built on PAPABAC. Raw WGS data files of five major foodborne pathogens (*C. jejuni*, *E. coli*, *L. monocytogenes*, *Salmonella enterica*, and *Shigella* spp.) are downloaded daily from public repositories with the aim of global surveillance of potential outbreaks worldwide. The inferred phylogenetic trees and information about all of the isolates in the system are available and searchable on the website (http://cge.cbs.dtu.dk/services/Evergreen).

The platform has been available since October 1st 2017, with logs reliably saved since October 28th 2017. The number of raw read files downloaded fluctuates with the work week of the public health laboratories. On busier days, more than 800 isolates are downloaded. The average number of isolates downloaded per day is 418. Downloading and mapping to the reference genomes take 130 minutes on average, with the majority of the time spent on downloading. Alignment of the raw reads and the generation of the consensus sequences takes on average 9 minutes per isolate. The computing time for the template sets is dependent on the number of non-redundant and new sequences in each set, but in most cases even the slowest is finalized within five hours (Figure 3).

As of June 26th 2018, the pipeline downloaded 82,043 isolates. Out of these, 63,276 isolates have been mapped to references with at least 99.0% identity and average depth of 11 (Figure S4A). The majority of the isolates were typed as *Salmonella enterica* (59.1%), followed by *Escherichia coli* (19.4%) (Figure S4B). The two largest template sets are *S.* Dublin and *S.* Typhimurium serovars, with both close to 9,500 isolates in total. After the homology reduction there were 3,216 and 5,093 non-redundant sequences in these sets, respectively. On average, 67% of the sequences are non-redundant in the template sets, while the E. coli template sets are the most diverse and the *Listeria monocytogenes* ones are the least diverse (Figure S4C). There were 122 isolates predicted to have a type not specified by the query (Table S1). Of these, 14 isolates were mixed samples, composed of both the queried and the non-queried organisms.
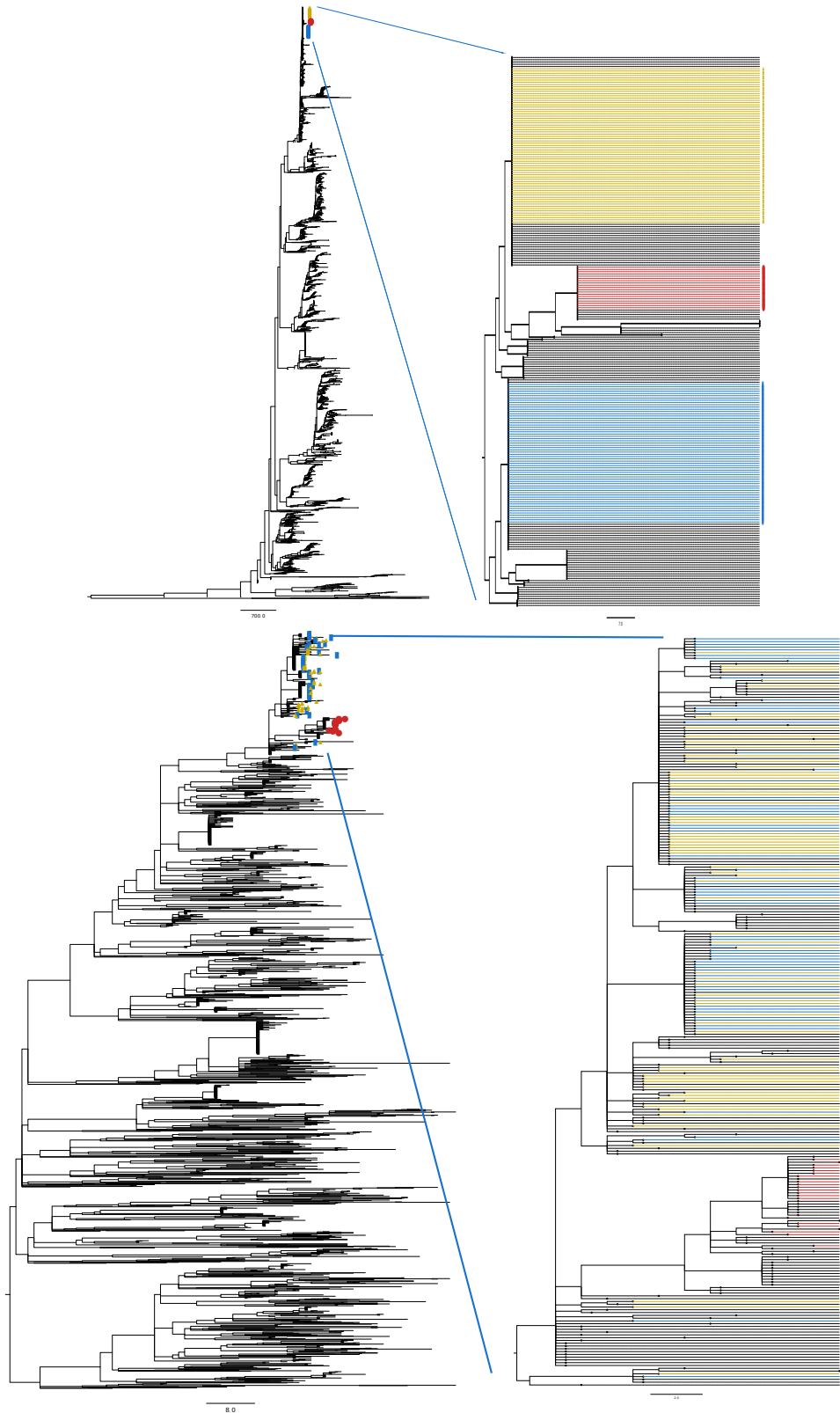
7

182

*Figure 4 Neighbor-joining tree for the Listeria_monocytogenes_07PF0776_NC_017728_1 set after the samples of the L. monocytogenes dataset were added. Isolates colored in concordance with Figure S3B*

185 The *L. monocytogenes* SNP pipeline benchmarking dataset[36] was added to the template set
186 (Listeria_monocytogenes_07PF0776_NC_017728_1) of the corresponding reference genome in Evergreen
187 Online, to test the sensitivity and accuracy of the clustering in large datasets. This template set at that moment
188 contained more than 2400 isolates, of which 1398 were non-redundant. The isolates were placed onto a clade
189 of a clonal lineage. The outbreak and outgroup isolates were separated in concordance with the ideal
190 phylogeny (Figure 4). The smaller clade of outbreak samples clustered to a sample (SRR538386) of an
191 environmental swab in 2014, from California, USA.

192 Isolates that were presumed to be from an *E. coli* O157:H7 outbreak were selected for the comparison of
193 Evergreen Online and the NCBI Pathogen Detection platform (NCBI-PD). They were located on the
194 Escherichia_coli_O157_H7_str_Sakai_chromosome_NC_002695_1 neighbor-joining (NJ) tree from Evergreen
195 Online and the PDS000000952.271 SNP cluster tree from NCBI-PD. The labelled isolates appeared in three
196 clusters on the NJ tree. There were 19.9 nucleotide differences between the yellow and the red cluster
197 representatives and 12.6 nucleotide differences between the yellow and the blue cluster representative. On
198 the PD tree, the isolates marked with red circles were on the same clade, while the ones marked with blue and
199 yellow were intermixing on clades that were, at most, 15 compatible characters apart (Figure 5).

8

200



201

*Figure 5 Selected isolates in the Escherichia_coli_O157_H7_str_Sakai_chromosome_NC_002695_1 NJ tree (top) and on the PDS000000952.271 SNP cluster maximum compatibility tree (bottom). The three largest clusters of the selected samples on the NJ tree are labelled with yellow, red and blue dots. These isolates were marked with the same labels on the NCBI-PD tree. The red labelled ones are on a single clade on the PD tree, while the blue and yellow isolates are mixing on two other clades.*

9

## Discussion

206    Whole-genome sequencing, performed alongside the traditional methods in routine microbiology, yields
207    hundreds to thousands of WGS isolates yearly in hospital, public health and food safety laboratories. This
208    amount of data is overwhelming for many, and there is a lack of methods to generate a quick overview and
209    help prioritize resources. The timely analysis of the sequencing data would allow the detection of more
210    bacterial outbreaks and aid the prevention of further spread. However, lack of human and computational
211    resources for this demanding task often hampers the prompt analysis of the data. Automating the initial
212    subtyping phase would facilitate the start of an outbreak investigation. PAPABAC offers rapid subtyping for a
213    wide range of prokaryotic organisms: the supplied database covers all bacterial subtypes with complete
214    genomes present in NCBI RefSeq. Further reference genomes could be added to increase the covered sequence
215    space, but the active curation of the reference database is not required for routine use. The selection of the
216    reference sequence for the phylogenomic analysis is fast and robust. It is independent of pre-assumptions
217    about the isolates. Misclassification during previous analysis does not introduce errors into the downstream
218    analysis. Contamination from another species is discarded during the consensus sequence generation. The
219    subtyping step via k-mer based mapping to a close reference also serves as a sequencing quality control
220    measure, because low-quality sequencing runs will typically result in isolates with low identity to any reference
221    and/or low depth. These isolates do not progress further to the phylogenomic analysis, as they would not yield
222    reliable results.
223    

224    The phylogenomic analysis performed on the template sets has higher discriminatory power than cg- or wg-
225    MLST. The underlying nucleotide difference method was validated in five different studies[6,34,35,37,38]. By using all
226    positions in the consensus sequences for estimating the genetic distance, instead of considering only selected
227    loci, we ensure a high level of sensitivity, as we also include mutations that occur between genes.

228    The clustering step during the genetic distance calculation was introduced in order to reduce the homology in
229    the template sets and thus reduce the computational burden as the template sets increase in size. However,
230    the clustering threshold of 10 nucleotide differences also constructs informative clusters of highly similar
231    isolates. Benchmarking with the *E. coli in vitro* evolution experiment dataset (Figure 2) showed that the
232    algorithm was capable of correctly clustering isolates that were derived from the same ancestor, while
233    distinguishing them from other closely related strains. The same sensitivity was demonstrated on empirical
234    outbreak datasets (Figure S3), where the pipeline clustered the outbreak-related strains and separated them
235    from the outgroup strains. Both the maximum likelihood inferred and the neighbor-joining trees placed the
236    outbreak strains correctly in the phylogeny. These results show, that PAPABAC provides quick and reliable
237    information about the close relatives of an outbreak strain to provide candidates to perform a more thorough
238    analysis on.

239    The design of PAPABAC means that once an isolate passed the homology reduction step, it will be present in
240    the subsequent runs of the pipeline. When an incoming isolate is highly similar to a non-redundant one, the
241    more recent will be the one that is clustered, added to the database and removed from further runs. Hence,
242    the cluster representatives and clusters are robust to the addition of new data to the analysis. Therefore,
243    PAPABAC yields a stable and communicable name for the clusters, comprised of the template name and the
244    cluster representative. This is an advantage over cg- and wg-MLST, where allelic profiles don't necessarily have
245    communicable names, and the clusters could merge.

246    Evergreen Online has been steadily processing WGS data of foodborne bacterial pathogen isolates collected
247    worldwide in real time (Figure S4A). It has been able to keep pace with the flow of the generated data that
248    mainly came from public health and food safety laboratories. Excluding the download time and the optional
249    maximum likelihood based phylogenetic inference, the whole analysis is done in less than a day, even for

250 template sets with thousands of isolates (Figure 3). This turnover time facilitates quick response in a potential
251 outbreak scenario.

252 The isolates are not distributed equally across the templates in the system (Figure S4B). Out of the five queried
253 species, *S. enterica* isolates are disproportionally represented. Sequences in the *S.* Dublin and the *S.*
254 Typhimurium LT2 template sets comprise in total approximately half of the *S. enterica* isolates. The sequence
255 diversity in the template sets is varied, but the homology reduction on the template sets reduces the number
256 of sequences approximately by a third, significantly decreasing the computational time. The *L. monocytogenes*
257 template sets were the least diverse, which could be due to sampling bias: bacteria that are present in the
258 environment are routinely sampled from food production sites multiple times, producing highly similar
259 sequences, that are then removed from the ongoing analysis. We also tested how a large number of sequences
260 already present in a template set would affect the ability of the pipeline to discriminate between samples
261 (Figure 4). The template set that corresponded to the stone fruit *L. monocytogenes* outbreak dataset reference
262 had more than 1,000 non-redundant isolates, which was ideal for the test analysis. The isolates that were part
263 of the same outbreak clustered together and formed the two expected outbreak clusters, despite the
264 confounding presence of the sequences already in the template set. The smaller clade, however, had a
265 different cluster representative when using all data for the template set, compared with analysis of the
266 outbreak data alone: an environmental sample, that could be related to the outbreak, as it was sampled from
267 the same US state and year (California, 2014) as the samples in the outbreak dataset. These findings indicate
268 that the pipeline is capable of identifying closely related samples, however it is necessary to conduct
269 epidemiological analysis and apply other knowledge when interpreting the results.

270 Evergreen Online allows for automated selection of closely related isolates out of thousands, which is also the
271 objective of NCBI-PD. *E. coli* isolates, situated on three clusters in Evergreen Online and supposedly from an
272 outbreak, were located in NCBI-PD and their placement in the SNP cluster tree was compared to the Evergreen
273 Online tree (Figure 5). One cluster (red) was in agreement between the two platforms, and samples from the
274 other two (yellow and blue) clusters were intermixing on a clade on the NCBI-PD tree. The nucleotide
275 difference counts between these samples are low and the differences between the phylogenomic methods
276 could lead to differences in the finer details of the inferred phylogenies. The homology reducing clustering in
277 Evergreen Online means that any sample in the cluster is less than 10 nucleotide differences from the cluster
278 representative, however, the differences between the samples could amount to 18 nucleotides. The
279 compatible character distances on the NCBI-PD tree between the mixed samples are less than 18 characters.
280 Taking this into account, the observed distribution of the yellow and blue labeled samples is concordant with
281 our results.

282 *Table 1 Comparison of pipelines for large-scale surveillance for pathogenic bacteria*

|  | SnapperDB | NCBI-PD | PAPABAC |
|---|---|---|---|
| For a wide range of bacterial species | x | x | x |
| Only sequencing data is required input | - | x | x |
| Whole-genome based | x | x | x |
| Assembly-free | x | - | x |
| Quality control steps | x | x | x |
| Automated phylogenomic analysis | - | x | x |
| Stable clustering of samples across runs | - | - | x |
| Communicable nomenclature for subtype and cluster | X | - | x |
| Open source | x | - | x |

11

283

284   In summary, we developed PAPABAC with the aim of rapid subtyping and continuous phylogenomic analysis on
285   a growing number of bacterial samples. PAPABAC overcomes limitations of cg- and wg-MLST approaches by
286   tolerating genomic variation during subtyping, but providing greater sensitivity during the phylogenomic
287   analysis. It was benchmarked on datasets created for testing SNP-based pipelines, and was proved to be
288   accurate in discriminating between outbreak related and non-related samples. The software is open source and
289   fulfills expectations put to WGS-based surveillance pipelines (*Table 1*). Evergreen Online, an application made
290   for the global surveillance of foodborne bacterial pathogens, demonstrates the accuracy, speed, stability and
291   practicality of PAPABAC on thousands of samples via an on-going analysis, where the results are published
292   online.

293

294   References

295   1.     Maiden, M. C. J. Multilocus Sequence Typing of Bacteria. *Annu. Rev. Microbiol.* **60,** 561–588 (2006).

296   2.     Larsen, M. V. *et al.* Multilocus Sequence Typing of Total-Genome-Sequenced Bacteria. *J. Clin. Microbiol.*
297         **50,** 1355–1361 (2012).

298   3.     Joensen, K. G., Tetzschner, A. M. M., Iguchi, A., Aarestrup, F. M. & Scheutz, F. Rapid and easy in silico
299         serotyping of Escherichia coli isolates by use of whole-genome sequencing data. *J. Clin. Microbiol.* **53,**
300         2410–2426 (2015).

301   4.     Köser, C. U. *et al.* Rapid whole-genome sequencing for investigation of a neonatal MRSA outbreak. *N.*
302         *Engl. J. Med.* **366,** 2267–75 (2012).

303   5.     Mellmann, A. *et al.* Prospective Genomic Characterization of the German Enterohemorrhagic
304         Escherichia coli O104:H4 Outbreak by Rapid Next Generation Sequencing Technology. *PLoS One* **6,**
305         e22751 (2011).

306   6.     Joensen, K. G. *et al.* Real-time whole-genome sequencing for routine typing, surveillance, and outbreak
307         detection of verotoxigenic Escherichia coli. *J. Clin. Microbiol.* **52,** 1501–1510 (2014).

308   7.     WHO. *Whole genome sequencing for foodborne disease surveillance: landscape paper.* (2018).

309   8.     Deng, X., den Bakker, H. C. & Hendriksen, R. S. Genomic Epidemiology: Whole-Genome-Sequencing-
310         Powered Surveillance and Outbreak Investigation of Foodborne Bacterial Pathogens. *Annu Rev Food Sci*
311         *Technol* **7,** 1–22 (2016).

312   9.     GenomeTrakr Network. Available at:
313         https://www.fda.gov/food/foodscienceresearch/wholegenomesequencingprogramwgs/ucm363134.ht
314         m. (Accessed: 27th June 2018)

315   10.    COMPARE Europe. Available at: http://www.compare-europe.eu.

316   11.    Nadon, C. *et al.* PulseNet International: Vision for the implementation of whole genome sequencing
317         (WGS) for global food-borne disease surveillance. *Euro Surveill.* **22,** (2017).

318   12.    Timme, R. E., Sanchez Leon, M. & Allard, M. W. Utilizing the Public GenomeTrakr Database for
319         Foodborne Pathogen Traceback. in 201–212 (Humana Press, New York, NY, 2019). doi:10.1007/978-1-
320         4939-9000-9_17

321   13.    Pathogen Detection - NCBI. Available at: https://www.ncbi.nlm.nih.gov/pathogens/. (Accessed: 27th

322    June 2018)

323    14.    Cherry, J. L. A practical exact maximum compatibility algorithm for reconstruction of recent evolutionary
324          history. *BMC Bioinformatics* **18,** 127 (2017).

325    15.    Alikhan, N.-F., Zhou, Z., Sergeant, M. J. & Achtman, M. A genomic overview of the population structure
326          of Salmonella. *PLOS Genet.* **14,** e1007261 (2018).

327    16.    Cody, A. J., Bray, J. E., Jolley, K. A., McCarthy, N. D. & Maiden, M. C. J. Core Genome Multilocus
328          Sequence Typing Scheme for Stable, Comparative Analyses of Campylobacter jejuni and C. coli Human
329          Disease Isolates. *J. Clin. Microbiol.* **55,** 2086–2097 (2017).

330    17.    Institut Pasteur MLST databases and software. Available at: https://bigsdb.pasteur.fr/. (Accessed: 28th
331          May 2019)

332    18.    Ghanem, M. & El-Gazzar, M. Development of Mycoplasma synoviae (MS) core genome multilocus
333          sequence typing (cgMLST) scheme. *Vet. Microbiol.* **218,** 84–89 (2018).

334    19.    Higgins, P. G., Prior, K., Harmsen, D. & Seifert, H. Development and evaluation of a core genome
335          multilocus typing scheme for whole-genome sequence-based typing of Acinetobacter baumannii. *PLoS*
336          *One* **12,** e0179228 (2017).

337    20.    Ghanem, M. *et al.* Core Genome Multilocus Sequence Typing: a Standardized Approach for Molecular
338          Typing of *Mycoplasma gallisepticum*. *J. Clin. Microbiol.* **56,** (2017).

339    21.    Bletz, S., Janezic, S., Harmsen, D., Rupnik, M. & Mellmann, A. Defining and Evaluating a Core Genome
340          Multilocus Sequence Typing Scheme for Genome-Wide Typing of *Clostridium difficile*. *J. Clin. Microbiol.*
341          **56,** (2018).

342    22.    Zhou, H., Liu, W., Qin, T., Liu, C. & Ren, H. Defining and Evaluating a Core Genome Multilocus Sequence
343          Typing Scheme for Whole-Genome Sequence-Based Typing of Klebsiella pneumoniae. *Front. Microbiol.*
344          **8,** (2017).

345    23.    Kohl, T. A. *et al.* Whole-Genome-Based Mycobacterium tuberculosis Surveillance: a Standardized,
346          Portable, and Expandable Approach. *J. Clin. Microbiol.* **52,** 2479–2486 (2014).

347    24.    Moran-Gilad, J. *et al.* Design and application of a core genome multilocus sequence typing scheme for
348          investigation of Legionnaires' disease incidents. *Euro Surveill.* **20,** (2015).

349    25.    Leekitcharoenphon, P. *et al.* Comparative genomics of quinolone-resistant and susceptible
350          Campylobacter jejuni of poultry origin from major poultry producing European countries (GENCAMP).
351          *EFSA Support. Publ.* **15,** (2018).

352    26.    Pathogenwatch | A Global Platform for Genomic Surveillance. Available at: https://pathogen.watch/.
353          (Accessed: 28th May 2019)

354    27.    Kvistholm Jensen, A. *et al.* Whole-genome Sequencing Used to Investigate a Nationwide Outbreak of
355          Listeriosis Caused by Ready-to-eat Delicatessen Meat, Denmark, 2014. *Clin. Infect. Dis.* **63,** 64–70 (2016).

356    28.    Schjørring, S. *et al.* Cross-border outbreak of listeriosis caused by cold-smoked salmon, revealed by
357          integrated surveillance and whole genome sequencing (WGS), Denmark and France, 2015 to 2017. *Euro*
358          *Surveill.* **22,** (2017).

359    29.    Ford, L. *et al.* Incorporating Whole-Genome Sequencing into Public Health Surveillance: Lessons from
360          Prospective Sequencing of Salmonella Typhimurium in Australia. *Foodborne Pathog. Dis.* **15,** 161–167
361          (2018).

13

30. Holmes, A., Dallman, T. J., Shabaan, S., Hanson, M. & Allison, L. Validation of Whole-Genome Sequencing for Identification and Characterization of Shiga Toxin-Producing Escherichia coli To Produce Standardized Data To Enable Data Sharing. *J. Clin. Microbiol.* **56,** (2018).

31. Woksepp, H., Ryberg, A., Berglind, L., Schön, T. & Söderman, J. Epidemiological characterization of a nosocomial outbreak of extended spectrum β-lactamase Escherichia coli ST-131 confirms the clinical value of core genome multilocus sequence typing. *APMIS* **125,** 1117–1124 (2017).

32. Davis, S. *et al.* CFSAN SNP Pipeline: an automated method for constructing SNP matrices from next-generation sequence data. *PeerJ Comput. Sci.* **1,** e20 (2015).

33. Dallman, T. *et al.* SnapperDB: a database solution for routine sequencing analysis of bacterial isolates. *Bioinformatics* **81,** 3946–3952 (2018).

34. Ahrenfeldt, J. *et al.* Bacterial whole genome-based phylogeny: construction of a new benchmarking dataset and assessment of some existing methods. *BMC Genomics* **18,** 19 (2017).

35. Leekitcharoenphon, P., Nielsen, E. M., Kaas, R. S., Lund, O. & Aarestrup, F. M. Evaluation of Whole Genome Sequencing for Outbreak Detection of Salmonella enterica. *PLoS One* **9,** e87991 (2014).

36. Timme, R. E. *et al.* Benchmark datasets for phylogenomic pipeline validation, applications for foodborne pathogen surveillance. *PeerJ* **5,** e3893 (2017).

37. Kaas, R. S., Leekitcharoenphon, P., Aarestrup, F. M. & Lund, O. Solving the problem of comparing whole bacterial genomes across different sequencing platforms. *PLoS One* **9,** e104984 (2014).

38. Joensen, K. G. *et al.* Evaluating next-generation sequencing for direct clinical diagnostics in diarrhoeal disease. *Eur. J. Clin. Microbiol. Infect. Dis.* **36,** 1325–1338 (2017).

39. Clausen, P. T. L. C., Aarestrup, F. M. & Lund, O. Rapid and precise alignment of raw reads against redundant databases with KMA. *BMC Bioinformatics* **19,** 307 (2018).

40. Hobohm, U., Scharf, M., Schneider, R. & Sander, C. Selection of representative protein data sets. *Protein Sci.* **1,** 409–417 (1992).

41. Saitou, N. & Nei, M. The neighbor-joining method: a new method for reconstructing phylogenetic trees. *Mol. Biol. Evol.* **4,** 406–25 (1987).

42. Studier, J. & Keppler, K. A note on the neighbor-joining algorithm of Saitou and Nei. *Mol. Biol. Evol.* **5,** 729–731 (1988).

43. Nguyen, L.-T., Schmidt, H. A., von Haeseler, A. & Minh, B. Q. IQ-TREE: a fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol. Evol.* **32,** 268–74 (2015).

44. Argimón, S. *et al.* Microreact: visualizing and sharing data for genomic epidemiology and phylogeography. *Microb. Genomics* **2,** (2016).

45. Huerta-Cepas, J. *et al.* ETE 3: Reconstruction, Analysis, and Visualization of Phylogenomic Data. *Mol. Biol. Evol.* **33,** 1635–1638 (2016).

46. Revell, L. J. phytools: an R package for phylogenetic comparative biology (and other things). *Methods Ecol. Evol.* **3,** 217–223 (2012).

47. Schliep, K. P. phangorn: phylogenetic analysis in R. *Bioinformatics* **27,** 592–593 (2011).

48. CDC. Multistate Outbreak of E. coli O157:H7 Infections Linked to Romaine Lettuce (Final Update) | Investigation Notice: Multistate Outbreak of E. coli O157:H7 Infections April 2018 | E. coli | CDC.

14

401       Available at: https://www.cdc.gov/ecoli/2018/o157h7-04-18/index.html. (Accessed: 7th August 2018)

402

403

## Methods

**Bioinformatics pipeline: PAPABAC**

406  The pipeline takes raw whole-sequencing reads (fastq files) as input. Matching reference sequences
407  (templates) in our reference database, that have greater than 99.0% identity and a minimum average depth of
408  11, are identified for the isolates using 16-mers via KMA[39] in sparse mode. Multiple templates are accepted, if
409  they meet the criteria, allowing for the procession of mixed samples. Information about the runs and their
410  templates are inserted into the main SQLite database. The isolates are grouped into sets according to the
411  matched templates. The next steps are performed in these sets in parallel. The isolate reads are mapped to the
412  template using the mapping algorithm of NDtree[35], yielding equal-length consensus sequences. The Z-score
413  threshold for accepting a base is set to 1.96, and the majority base have to be present in 90% of the mapped
414  reads.

415  Genetic distance based on nucleotide difference is calculated pairwise between the previous, non-redundant
416  isolates and the new isolates. Positions with ambiguous bases are discarded. The new isolates are clustered to
417  the non-redundant ones with a threshold of 10, in order to reduce the homology in each set and form
418  informative clusters. In this step, the non-redundant isolate is prioritized over the new isolate and becomes the
419  cluster representative. After the clustering, the remaining new isolates are clustered together with the
420  Hobohm 1 algorithm[40]. In this case, the cluster representative is the one that has already passed the
421  redundancy threshold. The information about new or extended clusters is saved to the main SQLite database. A
422  distance matrix is constructed for all non-redundant isolates and saved to disk for use in the next run. A
423  distance-based phylogenetic tree is inferred by neighbor-joining[41,42]. If there are less than 600 non-redundant
424  isolates in the set, then a whole-genome based approximate maximum likelihood phylogenetic tree is also
425  inferred using IQ-tree[43], where the neighbor-joining tree is the starting tree and the GTR nucleotide
426  substitution model is used. The clustered isolates are placed back onto the clades with zero distances to the
427  cluster representative. Their tip labels start with an asterisk. The information about the trees is saved to the
428  main SQLite database.

429  When all the phylogenetic trees with new isolates have been inferred, then the main SQLite database is
430  queried for the list of all isolates, their templates, cluster representatives (if there is any) and the latest
431  phylogenetic tree they are on. This information is printed to a tab-separated file.

432  Scripts and installation instructions are available on Bitbucket:
433  https://bitbucket.org/genomicepidemiology/evergreen

**Online Evergreen platform**

435  A query is made to the National Center for Biotechnology Information (NCBI) Sequencing Read Archive (SRA)
436  for the newly published Illumina paired-end sequenced isolates of *Campylobacter jejuni*, *Escherichia coli*,
437  *Listeria monocytogenes*, *Salmonella enterica*, and *Shigella spp.* on a daily basis. Fastq files of raw sequencing
438  reads and the corresponding metadata (collection date, location, institute, source, etc.) are acquired either
439  from SRA or from the European Nucleotide Archive (ENA). The sample inclusion criteria is known metadata for
440  collection date and location, and in addition, samples are included from the following institutions: Unites States
441  Center for Disease Control, United States Food and Drug Administration, Food Safety and Inspection Service,

15

442  Public Health England, University of Aberdeen, University Hospital Galway, Statens Serum Institut, Norwegian
443  Institute of Public Health. The downloaded isolates are the input to PAPABAC. The metadata are saved in the
444  main SQLite database, and added to the tip labels on the phylogenetic trees.

445  Individual subtrees are inferred from isolates with less than 20 SNPs distance from each cluster-representative,
446  considering only the positions in the sequences where there is no missing data. No tree is inferred, if no genetic
447  difference is found. The subtrees are inserted into an SQLite database.

448  Once all instances of the second wrapper script have finished, then the SQLite databases are queried for the list
449  of available phylogenetic trees (the maximum likelihood trees preferred over neighbor-joining ones), changes
450  in the clusters and the list of all isolates in the system, which is then used to update the website. For
451  visualization in external programs, such as Microreact[44], the phylogenetic trees can be downloaded as newick
452  files and the corresponding metadata as tab separated files.

453  **Architecture**

454  The pipeline is written in Python 2.7 and Bash in Unix environment. In addition to the standard Anaconda
455  Python 2.7 packages, it also requires ETE Toolkit v3.0[45] and Joblib v0.11 (https://pythonhosted.org/joblib)
456  packages to be installed. Neighbor program from the PHYLIP package v3.697
457  (http://evolution.genetics.washington.edu/phylip.html) and IQ-tree v1.6.4[43] are used for the phylogenetic tree
458  inference. The SQL database management is performed with SQLite v3.20.1 (https://www.sqlite.org).

459  The two main parts of the pipeline have their own wrapper scripts. PAPABAC can be run on a personal
460  computer with as few as four cores.
461  Evergreen Online is running on a high-performance computing cluster, utilizing the Torque (Adaptive
462  Computing Inc., USA) job scheduler. The first wrapper is run in one instance on 20 cores, meanwhile the second
463  wrapper is run once on 20 cores for each template that has at least one new run, in a parallel fashion. When all
464  of these instances are finished running, a Bash script is launched to collect the information from the SQL
465  database, the website is updated and the job for the next day is scheduled.

466  **Reference database**

467  The reference sequences are complete prokaryotic chromosomal genomes from the NCBI RefSeq database.
468  Homology reduction was performed at a 99.0% sequence identity threshold with the Hobohm 1 algorithm. The
469  curated NCBI prokaryotic reference genomes were given priority in the process. The reference sequences and
470  the classification database could be downloaded via ftp
471  (ftp://ftp.cbs.dtu.dk/public/CGE/databases/Evergreen/).

472  **Website**

473  The phylogenetic trees are interactively visualized on the website (https://cge.cbs.dtu.dk/services/Evergreen/)
474  using the Phylocanvas API (http://phylocanvas.org). The isolates and clusters can be searched by SRA run ID,
475  which allows the quick localization of the clusters that increased in size via their cluster representative.

476  **Computational time comparison of continued phylogenomic analysis**

477  101 samples from the Escherichia coli in vitro evolution experiment dataset by Ahrenfeldt et al. were batched
478  according to their sampling time. The parallelization in PAPABAC was disabled. The traditional method meant
479  that the analysis was carried out on all the samples up to the given batch, starting anew each time, but using
480  the same scripts as PAPABAC.

16

**481 Benchmarking of PAPABAC with the Escherichia coli in vitro evolution experiment dataset by Ahrenfeldt et**
**482 al.**

483 The last samples in each lineage were selected for the benchmarking. Therefore, the benchmarking dataset
484 constituted 50 tips on the ideal phylogeny. These samples were batched according to their sampling time (6th,
485 7th and 8th day). The batches were processed by PAPABAC chronologically. The pipeline was run with the
486 default parameters. Both maximum likelihood and neighbor-joining trees were inferred.

487 The phylogenetic trees inferred on all 50 isolates were trimmed for the reference sequence and compared with
488 the ideal phylogeny using the phytools R package (v0.6-60)[46]. The normalized Robinson-Foulds distance was
489 calculated between the ideal and the maximum likelihood, and the ideal and the neighbor-joining trees, after
490 the clustered isolates are removed from each pair of trees. The RF.dist function was utilized from the phangorn
491 R package (v2.4.0)[47].

**492 Benchmarking of PAPABAC with datasets from Timme et al.**

493 Each dataset was downloaded with the provided script into a distinct directory. The pipeline was run
494 individually on the datasets with default parameters. If the isolates were mapped to more than one template,
495 the phylogenetic trees of the template set with the highest number of isolates were evaluated. The maximum
496 likelihood trees were visually compared to the ideal phylogenies and checked for the distribution of the isolates
497 amongst the clades.

**498 Comparison with the NCBI Pathogen Detection platform**

499 *Escherichia coli* isolates were queried from the SQL database of Evergreen Online (EO) for the period of 2018-
500 03-15 and 2018-06-01, corresponding to a multistate outbreak of *E.coli* O157:H7 in the USA[48]. These samples
501 were subtyped using traditional MLST[2], as it was assumed, that the sequence type with the most isolates would
502 also include the outbreak samples. Sequence type 11, which is commonly corresponds to the O157:H7
503 serotype, was selected for further analysis. The corresponding samples and their SNP clusters were found in
504 the NCBI-PD platform. The phylogenetic tree for the SNP cluster with the most samples (PDS000000952.271)
505 was downloaded. The common samples were marked on both the NCBI-PD and the EO phylogenetic tree
506 (Escherichia_coli_O157_H7_str_Sakai_chromosome_NC_002695_1). The marked samples on the three biggest
507 clusters on the EO tree were labeled, and their placement on the NCBI-PD tree was visually inspected.
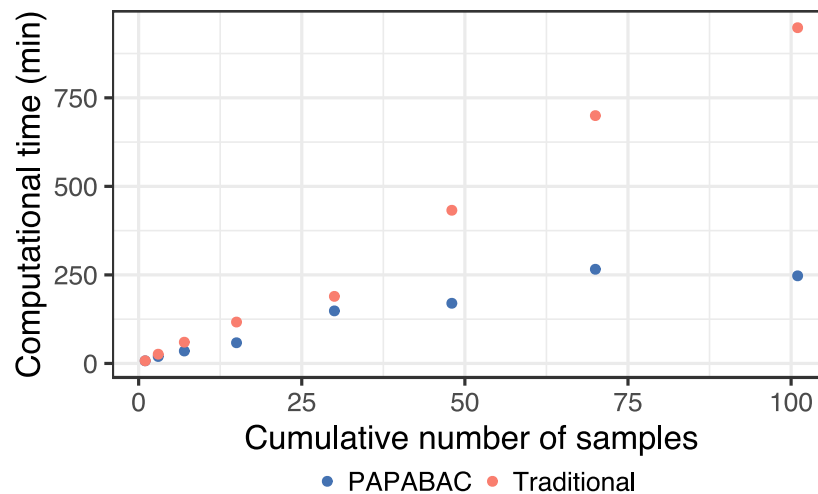
508    Supplementary material

509

510    *Table S1 Non-queried species, due to mislabelled or mixed samples*

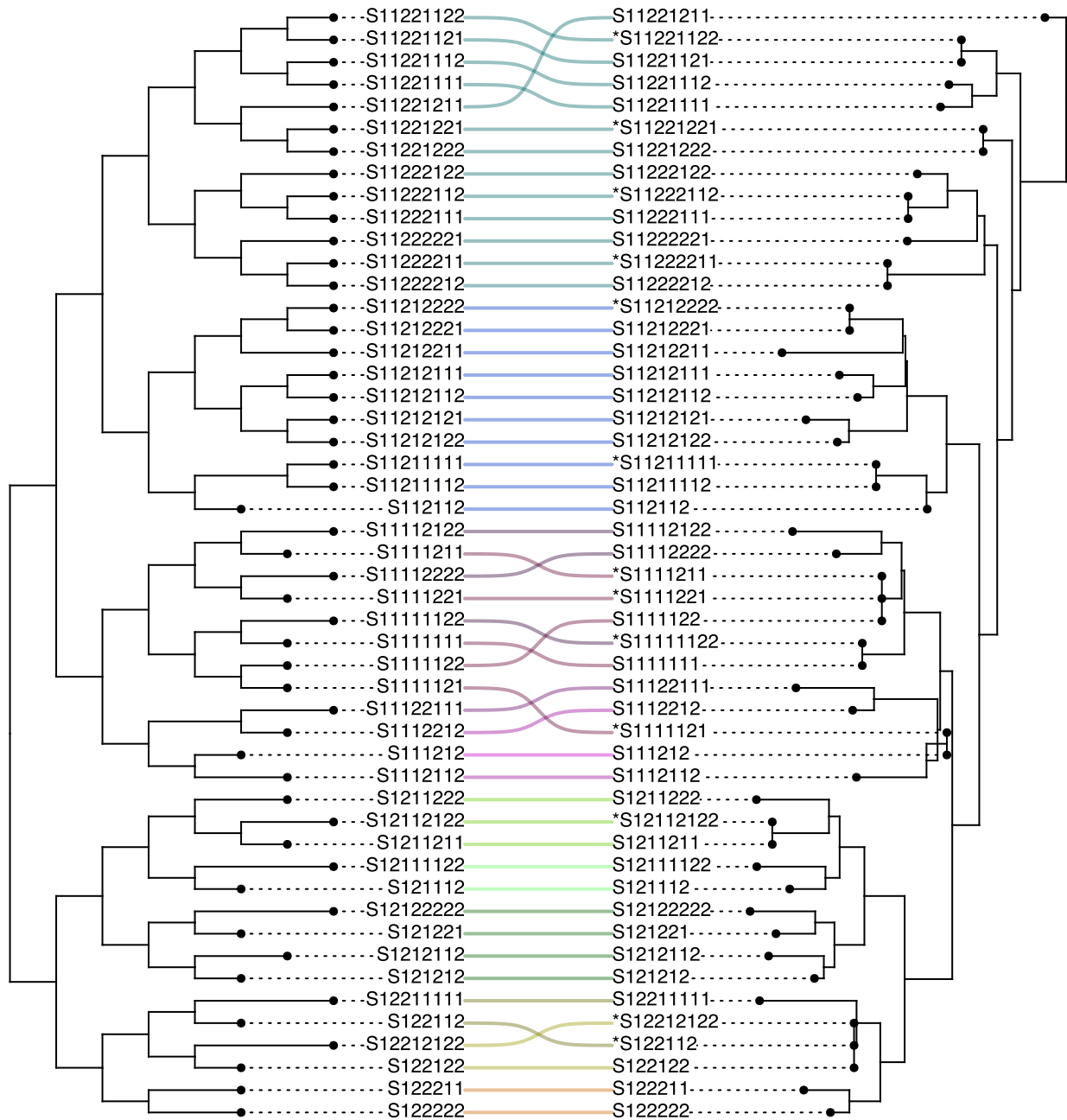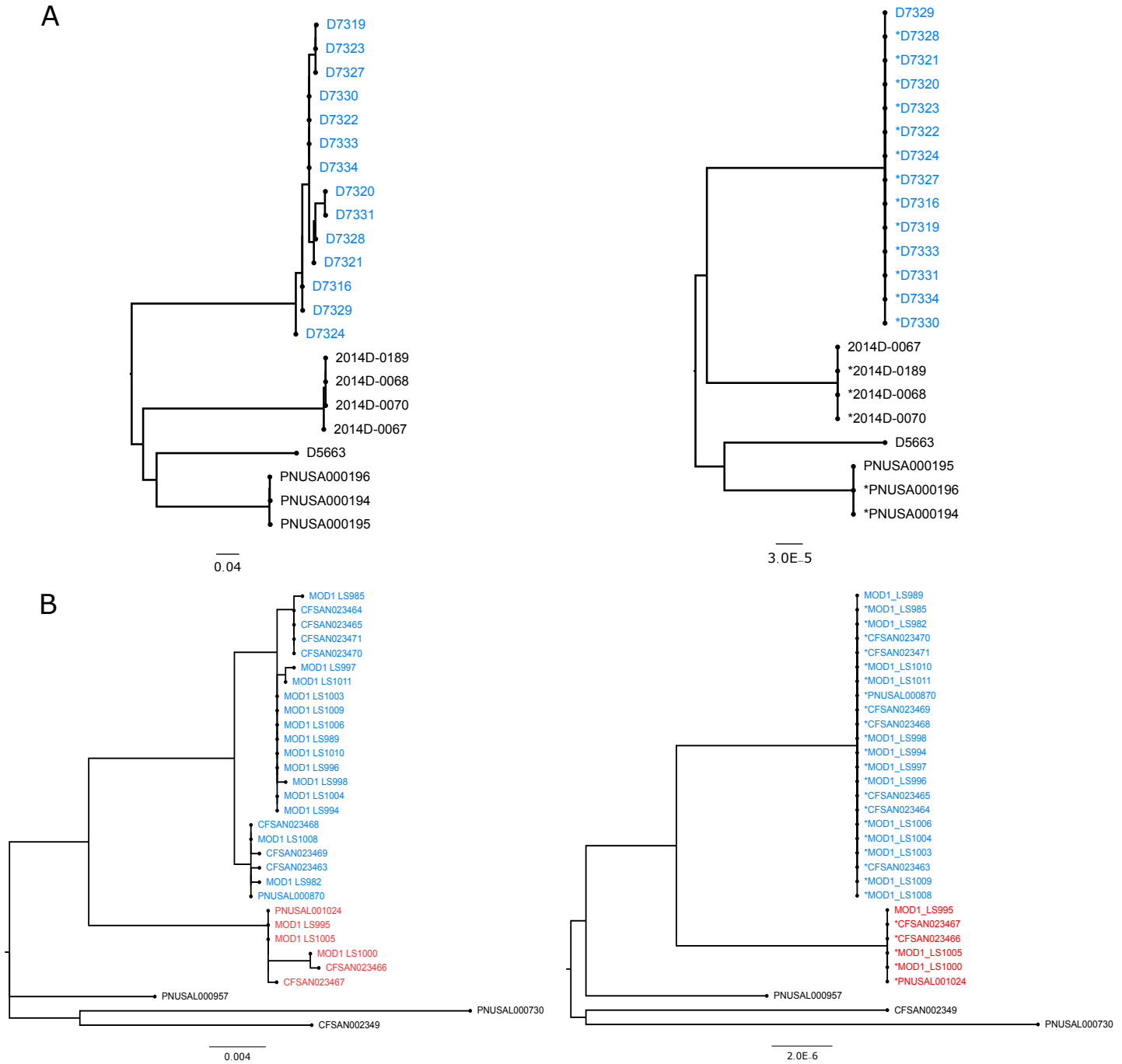| Genus | Species | Isolate |
|---|---|---|
| *Bacillus* | *subtilis* | 3 |
| *Bacillus* | *pumilus* | 2 |
| *Campylobacter* | *coli* | 58 |
| *Campylobacter* | *fetus* | 1 |
| *Citrobacter* | *amalonaticus* | 1 |
| *Enterobacter* | *cloacae* | 2 |
| *Enterococcus* | *faecalis* | 1 |
| *Escherichia* | *albertii* | 5 |
| *Hafnia* | *alvei* | 3 |
| *Klebsiella* | *pneumoniae* | 7 |
| *Listeria* | *ivanovii* | 1 |
| *Morganella* | *morganii* | 7 |
| *Peptoclostridium* | *difficile* | 1 |
| *Proteus* | *mirabilis* | 7 |
| *Providencia* | *stuartii* | 2 |
| *Pseudomonas* | *aeruginosa* | 6 |
| *Raoultella* | *ornithinolytica* | 1 |
| *Salmonella* | *bongori* | 11 |
| *Staphylococcus* | *epidermidis* | 1 |
| *Streptococcus* | *agalactiae* | 1 |

511

512



513

514    *Figure S1 Computational time of the Escherichia coli in vitro evolution dataset where the samples were added in batches based on the*
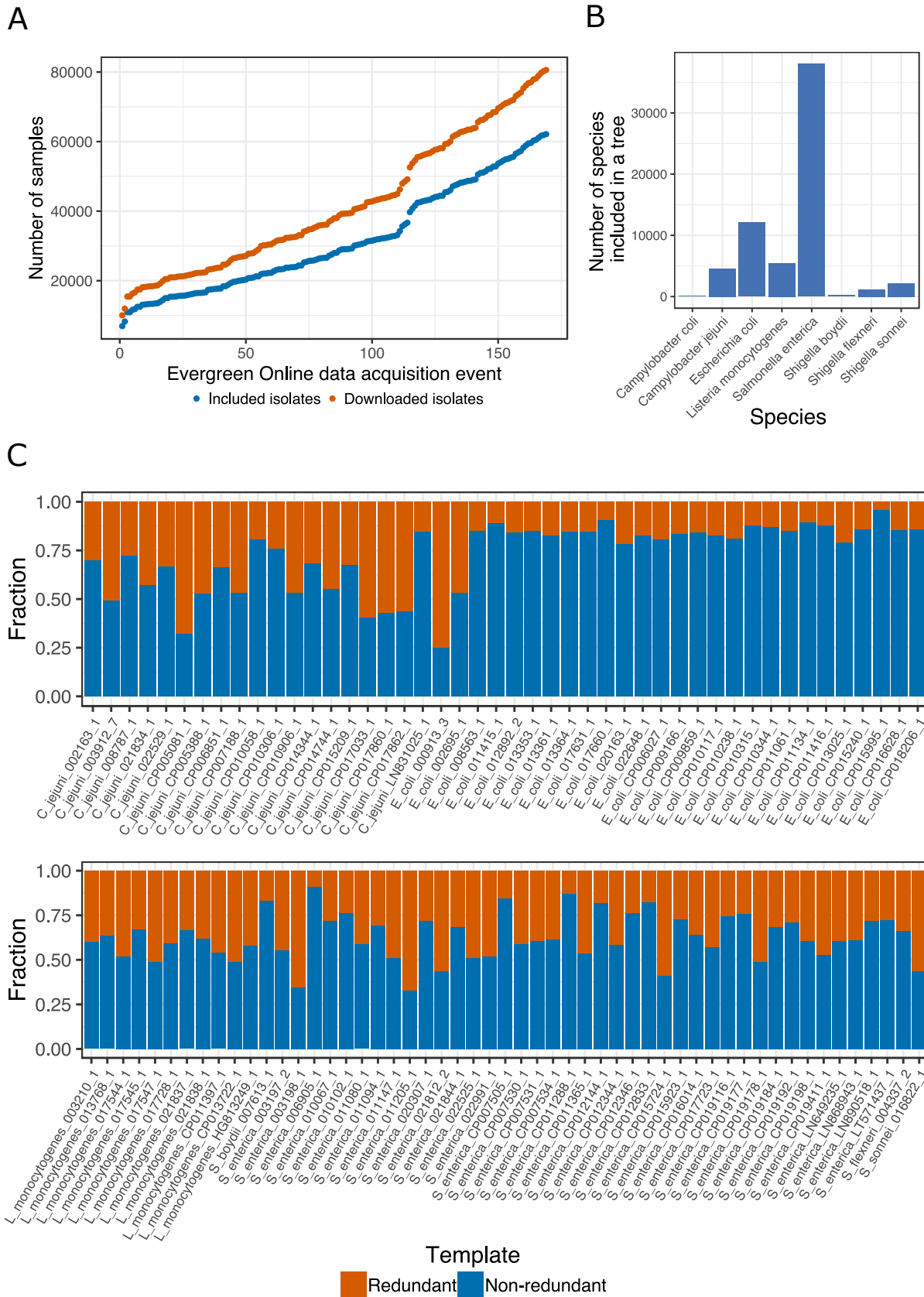515    *sampling time.*

516

18

517

518  *Figure S2 Comparison of the ideal tree (left) to the PAPABAC neighbor-joining tree made of the in vitro experiment dataset[34]Taxa with an*
519  *asterisk were clustered together with the taxa in the same clade.*

520

Figure S3 Maximum likelihood trees of (A) Campylobacter jejuni and (B) Listeria monocytogenes SNP pipeline benchmarking datasets. The trees on the left are the "ideal" phylogenies by Timme et al. The colored (blue, red) clades contain the outbreak strains, while the black ones are non-related isolates. The reference sequences were trimmed from the trees.

524

525  *Figure S4 A) Number of downloaded and included isolates as function of data acquisition events B) Number of isolates for the species we*
526  *query for C) Fraction of non-redundant isolates in template sets larger than 100 isolates*