

1 **RAREFAN: a webservice to identify REPINs and RAYTs in bacterial genomes**

2 Carsten Fortmann-Grote, Julia Balk and Frederic Bertels

3

4 Max-Planck-Institute for Evolutionary Biology, Department of Microbial Population Biology

5

6

7 Corresponding author: Frederic Bertels, August-Thienemann-Straße 2, 24306 Plön, Germany,

8 bertels@evolbio.mpg.de.

9

10

11

12 Running title: REPIN/RAYT Finder and ANalyzer

13

14 Keywords: sequence analysis – mobile genetic elements – bacterial genomes –

15 *Stenotrophomonas maltophilia**maltophilia*

16

17 **Abstract**

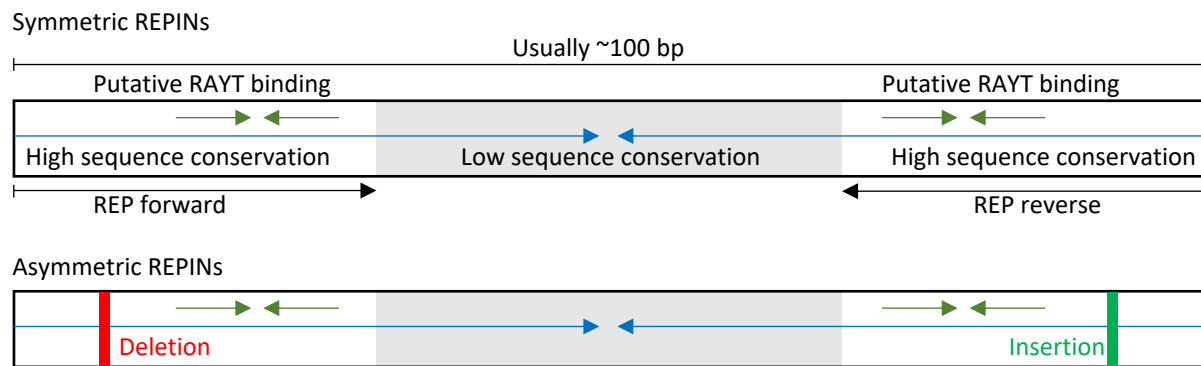
18 Compared to eukaryotes, ~~mobile genetic elements~~repetitive sequences are rare in bacterial
19 genomes and usually do not persist for long ~~in the genome~~. Yet, there is at least one class of
20 persistent prokaryotic mobile genetic elements: REPINs. REPINs are non-autonomous
21 transposable elements replicated by single-copy transposases called RAYTs. REPIN-RAYT
22 systems are mostly vertically inherited and have persisted in individual bacterial lineages for
23 millions of years. Discovering and analyzing REPIN populations and their corresponding RAYT
24 transposases in bacterial species can be rather laborious, hampering progress in understanding
25 REPIN-RAYT biology and evolution. Here we present RAREFAN, a webservice that identifies
26 REPIN populations and their corresponding RAYT transposase in a given set of bacterial
27 genomes. We demonstrate RAREFAN's capabilities by analyzing a set of 49 *Stenotrophomonas*
28 *maltophilia* genomes, containing nine different REPIN-RAYT systems. We guide the reader
29 through the process of identifying and analyzing REPIN-RAYT systems across *S. maltophilia*,
30 highlighting erroneous associations between REPIN and RAYTs, and ~~providing~~provide solutions
31 on how to find ~~the~~ correct associations. RAREFAN enables rapid, large-scale detection of
32 REPINs and RAYTs, and provides insight into the fascinating world of intragenomic sequence
33 populations in bacterial genomes.

34

35 **Introduction**

36 Repetitive sequences in bacteria are rare compared to most eukaryotic genomes. In eukaryotic
37 genomes, repetitive sequences are the result of the activities of persistent parasitic transposable
38 elements. In bacteria, in contrast, parasitic transposable elements cannot persist for long periods
39 of time (Park *et al.* 2021; van Dijk *et al.* 2022). (Park *et al.* 2021; van Dijk *et al.* 2022). To persist in
40 the gene pool, transposable elements must have to constantly infect novel hosts (Sawyer *et al.*
41 1987; Lawrence *et al.* 1992; Bichsel *et al.* 2010; Rankin *et al.* 2010; Wu *et al.* 2015; Park *et al.*
42 2021). Yet, there is at least one exception: a class of transposable elements called REPINs.

43



44

45 **Figure 1. The structure of symmetric and asymmetric REPINs.** A typical REPIN consists of two
46 highly conserved regions at the 5' and 3' end of the REPIN (white), separated by a spacer region
47 of lower sequence conservation (grey). The entire REPIN is a palindrome (blue arrows), which
48 means it can form hairpin structures in single stranded DNA or RNA. Each 5' and 3' region contains
49 a nested imperfect palindrome, which is referred to as REP (repetitive extragenic palindromic)
50 sequence and has first been described in *Escherichia coli* (Higgins *et al.* 1982). REPINs can be
51 either symmetric or asymmetric. Asymmetric REPINs have a deletion and a corresponding
52 insertion in the highly conserved 5' or 3' end, which leads to "bubbles" in the hairpin structure.
53 REPINs in *E. coli* are asymmetric, which makes analyses with RAREFAN more challenging. REPINs
54 (**REP-doublet forming hairPINS**) are bacterial repetitive sequences that occur in extragenic spaces
55 (Bertels, Rainey 2011). REPINs are non-autonomous mobile genetic elements that are duplicated
56 by a domesticated, single-copy RAYT transposase (Nunvar *et al.* 2010; Bertels, Rainey 2011; Ton-

57 ~~Hoang *et al.* 2012). In contrast to typical bacterial mobile genetic elements, REPINs have persisted~~
58 ~~for at least 100 million years in various species (Bertels, Gallie, *et al.* 2017; Park *et al.* 2021;~~
59 ~~Bertels, Rainey 2022), in the absence of horizontal transfer of RAYT transposases or REPIN~~
60 ~~populations (Bertels, Gallie, *et al.* 2017; Park *et al.* 2021; Bertels, Rainey 2022). These~~
61 ~~evolutionary characteristics are consistent with REPIN-RAYT systems providing a benefit to the~~
62 ~~host (Bertels, Rainey 2022).~~

The study Figure adapted from (Bertels, Rainey 2022).

63

64 REPINs are short (~100 bp) nested palindromic sequences (Figure 1) that consist of two inverted
65 REP (repetitive extragenic palindromic (Higgins *et al.* 1982)) sequences that can be present
66 hundreds of times per genome (Bertels, Rainey 2011a). Most REPINs are symmetric where the 5'
67 REP sequences is identical to the 3' REP sequences, with the occasional substitution (Bertels,
68 Rainey 2011a; b). However, there are also asymmetric REPINs where the 5' REP sequence differs
69 from the 3' REP sequence by a point deletion or insertion (Bertels, Rainey 2011a, 2022), which
70 makes the analysis and detection significantly more difficult (e.g., *Escherichia coli* REPINs).
71 Isolated REP sequences, REP singlets can also be found in the genome. These sequences are
72 decaying remnants of REPINs that are not mobile anymore (Bertels, Rainey 2011a). REPINs are
73 non-autonomous mobile genetic elements, which means they require a RAYT (REP Associated
74 tYrosine Transposase) transposase gene (also referred to as $tnpA_{REP}$) to replicate inside the
75 genome (Nunvar *et al.* 2010; Bertels, Rainey 2011a; Ton-Hoang *et al.* 2012).

76

77 Within a genome, each REPIN population is usually only associated to a single RAYT gene. Hence,
78 RAYT genes occur only in single copies per genome and do not copy themselves, unlike for
79 example insertion sequences where often multiple identical sequences are present inside the
80 genome. Unlike insertion sequences RAYT genes are also only inherited vertically, meaning they
81 are host-beneficial transposases that are coopted by the host (Bertels, Gallie, *et al.* 2017; Bertels,
82 Rainey 2022). The fact that REPINs and their corresponding RAYT genes are confined to a single

83 bacterial lineage makes them very special, in comparison to all other parasitic mobile genetic
84 elements in bacterial genomes (Bertels, Rainey 2022).

85

86 Of a total of five different RAYT families, there are only two RAYT families that are associated
87 with REPINs: Group 2 and Group 3 RAYTs (Bertels, Gallie, *et al.* 2017). Group 2 RAYTs are present
88 in most Enterobacteria and usually occur only once per genome associated with a single REPIN
89 population. In contrast, Group 3 RAYTs are found in most *Pseudomonas* species and are usually
90 present in multiple divergent copies per genome, each copy associated with a specific REPIN
91 population (Bertels, Gallie, *et al.* 2017).

92

93 REPINs and their corresponding RAYT genes occur exclusively in bacterial genomes and are
94 absent in eukaryotic or archaeal genomes (Bertels, Gallie, *et al.* 2017; Bertels, Rainey 2022).
95 Within bacterial genomes REPINs and RAYTs have been evolving in single bacterial lineages for
96 millions maybe even for a billion years (Bertels, Gallie, *et al.* 2017). The long term persistence of
97 REPINs in single bacterial lineages can also be observed when analyzing REPIN populations
98 (Bertels, Gokhale, *et al.* 2017; Bertels, Rainey 2022).

99

100 Parasitic insertion sequences usually occur in identical copies in bacterial genomes reflecting the
101 fact that insertion sequences persist only briefly in the genome before they are eradicated from
102 the genome or kill their host (Park *et al.* 2021). REPINs in contrast are only conserved at the ends
103 of the sequence (presumably due to selection for function), the rest of the sequence is highly
104 variable and only the hairpin structure is conserved (Bertels, Rainey 2011a). The sequence
105 variability of REPINs within the same genome reflects their long-term persistence in single
106 bacterial lineages (Bertels, Rainey 2022). REPINs cannot simply reinfect another bacterial lineage
107 since they rely for mobility on their corresponding RAYT, which itself is immobile.

108

109 RAYTs and REPINs are distinct from typical parasitic insertion sequences, yet we know very little
110 about their evolution or biology. Currently, it is completely unclear what kind of beneficial
111 function maintains REPINs and RAYTs for millions of years in the genome. The reason for our lack
112 of knowledge is not because REPINs and RAYTs are rare. They are ubiquitously found in many
113 important and well-studied model bacteria such as Enterobacteria, Pseudomonads, Neisseriads,
114 Xanthomonads. Microbial molecular biologists presumably encounter REPINs quite frequently.
115 However, connecting the presence or absence of REPINs/RAYTs with phenotypes is difficult if we
116 do not know when it is a REPIN that is present close to a gene of interest or a different type or
117 repeat sequence. Even if the scientist knows about the presence of a REPIN it is probably also
118 important to know whether a corresponding RAYT is present, since the function of REPINs largely
119 depends on the function of the presence of a corresponding RAYT gene (Bertels, Rainey 2022).

120

121 Yet, the identification of REPIN populations and their corresponding RAYTs can be rather
122 cumbersome.~~To facilitate REPIN studies, we have developed a webservice called~~ if done from
123 scratch. This is particularly true if the microbial molecular biologist is not aware of all the ins and
124 outs of REPIN and RAYT biology. Identifying REPINs starts with an analysis of short repetitive
125 sequences in the genome. If there are excessively abundant short sequences present in the
126 genome, the distribution of these sequences is analyzed next. If they are exclusively identical
127 tandem repeats without sequence variation and present in only one or two loci in the genome
128 then these sequences are probably part of a CRISPR array and not REPINs.

129

130 Here, we present RAREFAN (RAYT/REPIN Finder and ~~ANalyzer~~-Analyzer), a webservice that
131 automates the identification of REPINs and their corresponding RAYTs. RAREFAN is publicly
132 accessible at <http://rarefan.evolbio.mpg.de> and identifies REPIN populations and RAYTs inside a
133 set of bacterial genomes. RAREFAN also generates graphs to visualize the population dynamics
134 of REPINs, and assigns RAYT genes to their corresponding REPIN groups. Here we will
135 demonstrate RAREFAN's functionality by analyzing REPIN-RAYT systems in the bacterial species
136 *Stenotrophomonas maltophilia*.

A) Input (genomes)



B) Identifying sequence groups in reference

- 1) Determine all 21bp long sequences above a certain threshold
- 2) Group 1 Group 2 Group 3 Group sequences by vicinity (<30bp) in reference
- 3) Group 1 Group 2 Group 3 Determine most common sequence in each group

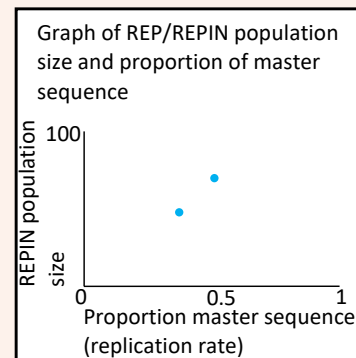
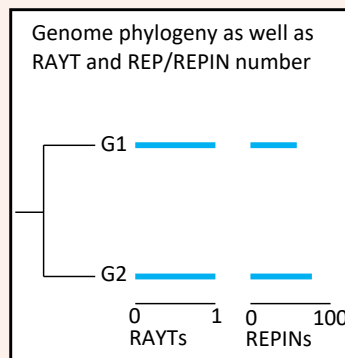
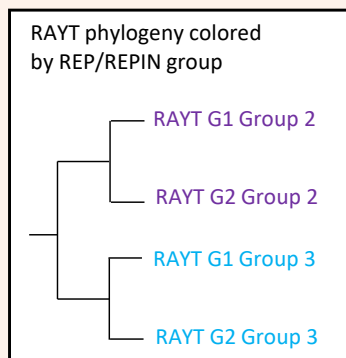
C) Identify REPINS/REPs in all genomes

- 1) Group 1 Group 3 Group 2 Generate all possible single mutants of each for each group sequence
 - 2) Retain sequences that are present in the genome, and repeat until no more new sequences are identified in the genome
- 2) < 200 bp REPIN
- If two sequences are found at a distance of less than 200 bp and in inverted orientation then they are designated a REPIN

D) Identify RAYTs and corresponding REP/REPIN group

- 1) RAYT
↓ TBLASTN
Input genome
Use TBLASTN of either Type 2 (e.g. from *E. coli*) or Type 3 (e.g. from *P. fluorescens* SBW25) RAYTs to identify RAYT occurrences in the input genomes
- 2) Group 2 REP/REPIN Group 3 REP/REPIN
Group 2 RAYT Group 3 RAYT
The RAYT gene's group becomes the group of the REP/REPIN in its vicinity (<200bp), if there is no REP/REPIN the RAYT remains unassociated

E) Output



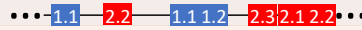
And other files, for example: REP/REPIN and RAYT positions and sequences, REP/REPIN frequency, REP/REPIN conservation.

A) Input (genomes)



B) Identifying REP sequence groups in reference genome

- 1) — 232 — 128 — 122 — 65 — 56 Determine all 21bp long sequences that occur more than 55 times
- 2) Group 1 Group 2 Group sequences by vicinity (< 15 bp) in reference genome
- 3) Group 1 Group 2 Group 3 Determine most common sequence in each group



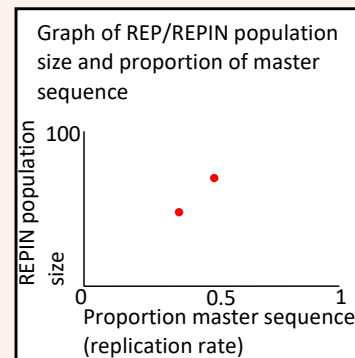
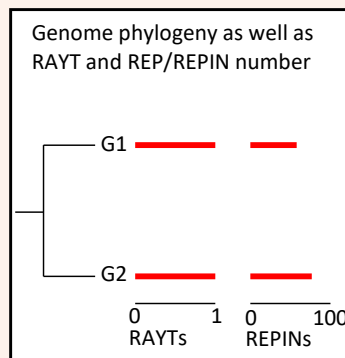
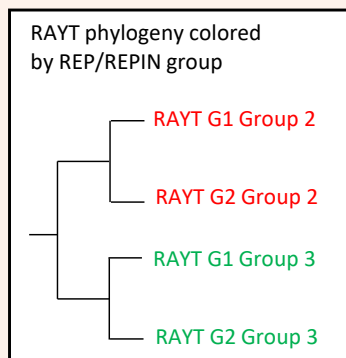
C) Identify REPINS/REPs in all genomes

- 1) Group 1 Group 2 Group 3 Generate all possible single mutants for each sequence
 - 2) < 200 bp REPIN
- Retain sequences that are present in the genome, and repeat until no more new sequences are identified in the genome
- If two sequences of the same group are found at a distance of less than 200 bp and in inverted orientation then they are designated a REPIN

D) Identify RAYTs and corresponding REP/REPIN group

- 1) RAYT Use TBLASTN of either Type 2 (e.g., from *E. coli*) or Type 3 (e.g., from *P. fluorescens* SBW25) RAYTs to identify RAYT occurrences in the input genomes
 - 2) Group 2 REP/REPIN Group 3 REP/REPIN
- The RAYT gene's group becomes the group of the REP/REPIN in its vicinity (< 200 bp), if there is no REP/REPIN the RAYT remains unassociated

E) Output



And other files, for example: REP/REPIN and RAYT positions and sequences, REP/REPIN frequency, REP/REPIN conservation.

Figure 4-2. RAREFAN workflow. (a) By default, RAREFAN requires the user to supply input sequences containing RAYTs and REPINs. These are, ideally, fully sequenced and complete genomes. (b) RAREFAN then identifies seed sequence groups (potential REP sequences) in the reference genome by first isolating all 21 bp (adjustable parameter) long sequences that occur more than 55 times (adjustable parameter) in the reference genome. It is likely that a large number of these sequences belong to the same REPIN sequence type, since the conserved part of REPINs is longer than 20bp. Hence, we grouped all sequences together that occur within 15 bp (adjustable parameter) of each other, anywhere in the genome. For example, if 'sequence 1' occurs 55 times and 'sequence 2' occurs 42 times then only one of these occurrences of 'sequence 1' needs to be within 15 bp of 'sequence 2' in order to be sorted into the same sequence group. All further analyses are performed only with the most common sequence in each sequence group. This sequence will be called the seed sequence. (c) The occurrences of the seed and mutated seed sequences are identified in all submitted genomes. If a mutated seed sequence is identified in a genome, then all single mutants of that seed sequence are searched

recursively in the same genome. All identified sequences that occur within 130 bp in inverted orientation of each other are designated REPINs. All other identified seed sequences and mutated seed sequences are REP singlets. **(d)** TBLASTN is used to identify RAYT homologs ([e-Value < 1e-30, adjustable parameter](#)) of either *E. coli* (Group 2 RAYT) or from *Pseudomonas fluorescens* SBW25 (Group 3 RAYT) across all submitted genomes. If a RAYT homolog is in the vicinity ([default 200 bp, adjustable parameter](#)) of a previously identified REPIN or REP singlet, then this RAYT is designated as associated with this REPIN group. ~~(e) Finally, RAREFAN plots three different summary graphs and~~ **(e)** The first graph contains a RAYT phylogeny computed from a nucleotide alignment of all identified RAYT genes. The alignment is calculated with MUSCLE (Edgar 2004) and a phylogeny with PHYML (Guindon *et al.* 2010). The RAYT phylogeny indicates what RAYTs are associated with what REPIN populations (largest sequence cluster calculated with MCL) via colour coding. In a second graph the abundance of each REPIN population and RAYT copy number are displayed on a genome phylogeny. If no genome phylogeny is supplied RAREFAN calculates a whole genome phylogeny of the submitted genomes using *andi* (Haubold *et al.* 2015). In the last graph REPIN population sizes are plotted in relation to the proportion of master sequences. Master sequences are the most abundant REPIN in each population. The REPIN population is the largest sequence cluster that is formed by REPIN sequences (REP sequences are excluded). The largest sequence cluster is identified by applying MCL with an inflation parameter of 1.2 to a sequence matrix where only sequences are connected that differ in exactly one position (Van Dongen 2000). RAREFAN also generates various files containing, for example, REP, REPIN, or RAYT sequences and their positions in the query genomes.

137 **Methods**

138 *Implementation*

139 RAREFAN is a modular webservice. It consists of a web frontend written in the python
140 programming language (Van Rossum, Drake Jr 1995) using the flask framework (Grinberg 2018),
141 a java (Arnold *et al.* 2005) backend for genomic sequence analysis and an R (~~[R Core Team 2016](#)~~)
142 ~~[shiny app \(RStudio, Inc 2013\) for data visualization. The software is developed and tested under \(R](#)~~
143 ~~[Core Team 2016\) shiny app \(RStudio, Inc 2013\) for data visualization. The software is developed](#)~~
144 ~~[and tested on](#)~~ the Debian GNU/Linux operating system (Kleinmann *et al.* 2021). All components
145 are released under the MIT opensource license (Initiative 2021) and can be obtained from our
146 public GitHub repository at [https://github.com/mpievolbio-](https://github.com/mpievolbio-scicom/rarefan)
147 ~~[scicom/rarefan](https://github.com/mpievolbio-scicom/rarefan)~~[https://github.com/mpievolbio-](https://github.com/mpievolbio-scicom/rarefan)
148 ~~[scicom/rarefan](https://github.com/mpievolbio-scicom/rarefan)~~[https://github.com/mpievolbio-scicom/rarefan.](https://github.com/mpievolbio-scicom/rarefan)

149 [The public RAREFAN instance at http://rarefan.evolbio.mpg.de](http://rarefan.evolbio.mpg.de) runs on a virtual cloud server with
150 [4 single-threaded CPUs and 16GB of shared memory provided and maintained by the Gesellschaft](#)
151 [für Wissenschaftliche Datenverarbeitung Göttingen \(GWDG\) and running the Debian GNU/Linux](#)
152 [Operating System \(Kleinmann *et al.* 2021\).](#)

153 The java backend drives the sequence analysis. It makes system calls to TBLASTN (Altschul *et al.*
154 1990) to identify RAYT homologs and to MCL (Van Dongen 2000) for clustering REPIN sequences-
155 [in order to determine REPIN populations.](#)

156 Jobs submitted through the web server are queued and executed as soon as the required
157 resources become available. Users are informed about the status of their jobs. After job
158 completion, the user can trigger the R shiny app to visualize the results.

159 The java backend can also be run locally *via* the command line interface (available for download
160 at <https://github.com/mpievolbio-scicomp/rarefan/releases>).

161 *Usage of the webservice*

162 The front page of our webservice allows users to upload their bacterial genomes in FASTA (.fas)
163 format (**Figure 1A2A**). Optionally, users may also provide RAYT protein FASTA sequences (.faa) or
164 phylogenies in NEWICK (.nwk) format. After successful completion of the upload process, the
165 user fills out a web form to specify the parameters of the algorithm:

166 • Reference sequence: Which of the uploaded genome sequences will be designated as
167 reference genome (see below for explanations). Defaults to the first uploaded filename
168 in alphabetical order.

169 • [Query RAYT: The RAYT gene that is used to identify homologous RAYTs in the query](#)
170 [genomes.](#)

171 • [Tree file: A phylogenetic tree of the reference genomes that can be provided by the user,](#)
172 [otherwise the tree will be calculated using andi \(Haubold *et al.* 2015\).](#)

173 • ~~Seed sequence length: The seed sequence length (in base pairs) is used to identify REPIN~~
174 ~~candidates from the input genomes. Default is 21 bp.~~

- 175 • Minimum seed sequence frequency: Lower limit on seed sequence frequency in the
176 reference genome to be considered as a REP candidate. Default is 55.
- 177 • Seed sequence length: The seed sequence length (in base pairs) is used to identify REPIN
178 candidates from the input genomes. Default is 21 bp.
- 179 • Distance group seeds: The maximum distance between a single occurrence of short
180 repetitive sequences to still be sorted into the same sequence group.
- 181 • Association distance REPIN-RAYT: The maximum distance at which a REP sequence can
182 be located from a RAYT gene to be linked to that RAYT gene.
- 183 • e-value cut-off: Alignment e-value cut-off for identifying RAYT homologs with TBLASTN.
184 Default is 1e-30.
- 185 • Analyse REPINs: Ticked REPINs will be analysed (two inverted REP sequences found at a
186 distance of less than 130 bp), if not ticked only short repetitive 21 bp long sequence will
187 be analysed.
- 188 • User email (optional): If provided, then the user will be notified by email upon run
189 completion.

190 The job is then ready for submission to the job queue. Upon job completion, links to browse and
191 to download the results, as well as a link to a visualization dashboard are provided. If a job runs
192 for a long time then users may also come back to RAREFAN at a later time, query their job status
193 and eventually retrieve their results by entering the run ID into the search field at
194 <http://rarefan.evolbio.mpg.de/results>. Relevant links and the run ID are communicated either on
195 the status site or by email if the user provided their email address during run configuration. Runs
196 are automatically deleted from the server after six months.

197 *Identification of REPs and REPINs*

198 ~~The algorithm to determine REP sequence groups has been described in previous papers and is~~
199 ~~slightly improved in our implementation (Bertels, Rainey 2011, 2022; Bertels, Gokhale, et al.~~
200 ~~2017). First, all N bp (21 bp by default) long seed sequences that occur more than M times (55 by~~
201 ~~default) are extracted from the reference genome. N and M are the seed sequence length and~~

202 ~~minimum seed sequence frequency, respectively (Figure 1B). All sequences occurring within the~~
203 ~~reference genome at least once within 15 bp of each other are then grouped together into n REP~~
204 ~~sequence groups (numbered 0-(n-1)). The most common sequence in each group, named REP~~
205 ~~seed sequence, is used for further analyses in each input genome.~~

206 ~~In the next step all possible point mutants of the seed sequences are generated and searched for~~
207 ~~in the genome (Figure 1C). The algorithm to determine REP sequence groups has been described~~
208 ~~in previous papers and is slightly improved (Bertels, Rainey 2011a, 2022; Bertels, Gokhale, et al.~~
209 ~~2017). In our current implementation REPs/REPIN populations are now automatically linked to~~
210 ~~RAYT genes.~~

211 ~~First, all N bp (21 bp by default) long seed sequences that occur more than M times (55 by default)~~
212 ~~are extracted from the reference genome. N and M are the seed sequence length and minimum~~
213 ~~seed sequence frequency, respectively (Figure 2B). All sequences occurring within the reference~~
214 ~~genome at least once within 15 bp of each other are then grouped together into n REP sequence~~
215 ~~groups (numbered 0-(n-1)). The most common sequence in each group, named REP seed~~
216 ~~sequence, is used for further analyses in each input genome.~~

217 ~~In the next step all possible point mutants of the seed sequences are generated and searched for~~
218 ~~in the genome (Figure 2C). If a sequence is found in the genome, then all possible point mutations~~
219 ~~are generated for this sequence as well and so on until no more sequences can be identified. If~~
220 ~~two sequences are found within 130 bp of each other in inverted orientation, then these are~~
221 ~~designated REPINs.~~

222 *Identification of RAYTs*

223 ~~RAYTs are identified using TBLASTN (Camacho et al. 2009) with either a protein sequence~~
224 ~~provided by the user or a Group 2 RAYT from *Escherichia coli* (yafM, Uniprot accession Q47152)~~
225 ~~or a Group 3 RAYT from *P. fluorescens* SBW25 (yafM, Uniprot accession C3JZZ6). The presence of~~
226 ~~RAYTs in the vicinity of a particular REPIN can be used to establish the association between the~~
227 ~~RAYT gene and a REPIN group (Figure 1D).~~

228 ~~Among all identified REP and REPIN sequences REPIN populations can be isolated. REPIN~~
229 ~~populations are determined by applying MCL using an inflation parameter of 1.2 (Van Dongen~~

230 2000) to a network of REP/REPIN sequences where all sequences that differ by exactly one
231 nucleotide are connected. The clustering results are stored in a file ending in .mc1. The
232 sequences of the largest REPIN population (excluding REP singlets) are isolated in a file ending in
233 largestCluster.nodes. The largest REPIN populations are shown in the REPIN population
234 plot and the master sequence correlation plot (Figure 4).

235

236 Identification of RAYTs

237 RAYTs are identified using TBLASTN (Camacho *et al.* 2009) with either a protein sequence
238 provided by the user, a Group 2 RAYT from *E. coli* (yafM, Uniprot accession Q47152) or a Group
239 3 RAYT from *P. fluorescens* SBW25 (yafM, Uniprot accession C3JZZ6). The presence of RAYTs in
240 the vicinity (default 200 bp) of a particular REPIN can be used to establish the association
241 between the RAYT gene and a REPIN group (Figure 2D). All positions of all REPINs and REP
242 sequences of a REPIN group are checked whether they occur within 200 bp (by default) of a RAYT
243 gene. If so then the RAYT gene is linked to the REPIN group in the file
244 repin rayt association.txt.

245 *Visualizations*

246 For each REPIN-RAYT group summary plots are generated. These include plots showing the RAYT
247 phylogeny, ~~REPIN population sizes in relation to the genome phylogeny (calculated from a~~
248 nucleotide alignment using MUSCLE (Edgar 2004) and PHYML (Guindon *et al.* 2010) to generate
249 a phylogeny), REPIN population sizes in relation to the genome phylogeny (provided by the user
250 or if not provided calculated by andi (Haubold *et al.* 2015)) as well as the proportion of master
251 sequences (most common REPIN in a REPIN population) in relation to REPIN population size
252 **(Figure 4E).**

253 *Other outputs*

254 Identified REPINs, REP singlets as well as RAYTs are written to FASTA formatted sequence files
255 and to tab formatted annotation files that can be read with the Artemis genome browser
256 (Rutherford *et al.* 2000). The REPIN-RAYT associations as well as the number of RAYT copies per
257 genome are written to tabular data files. A detailed description of all output files is provided in

258 the manual (<http://rarefan.evolbio.mpg.de/manual>;) and in the file “readme.md” in the output
259 [directory](#).

260

261 **Sequence analysis and annotation**

262 For verification of RAREFAN results, REPIN-RAYT-systems were ~~analyzed~~[analysed](#) in their
263 corresponding genomes using Geneious Prime (~~Java Version 11.0.12+7 (64 bit);~~
264 ~~Biomatters)-version 2022.2.2 (Kearse et al. 2012)~~). Nucleotide sequences and positions of REP
265 singlets, REPINs, and RAYTs were extracted from output files generated by RAREFAN and mapped
266 in the relevant *S. maltophilia* genome. ~~The association of a RAYT gene to a REPIN population has~~
267 ~~been validated when the corresponding seed sequence is flanking both ends of the RAYT gene~~
268 ~~within 130 bp~~. Complete RAREFAN data used for analysis can be accessed by using the run IDs
269 listed in **Table 1**.

270

271 **Table 1.- RAREFAN IDs linking to the raw data of the presented analyses.**

Run ID	Reference genome
1a8l7wu	<i>S. maltophilia</i> Sm53
mknhxp8	<i>S. maltophilia</i> AA1
pgfmaxx5	<i>S. maltophilia</i> FDAARGO_649
yy72i755	<i>S. maltophilia</i> AB550
78eu9zl0	<i>S. maltophilia</i> ISMMS3

272 Associated data can accessed by entering the run ID at <http://rarefan.evolbio.mpg.de/results>.

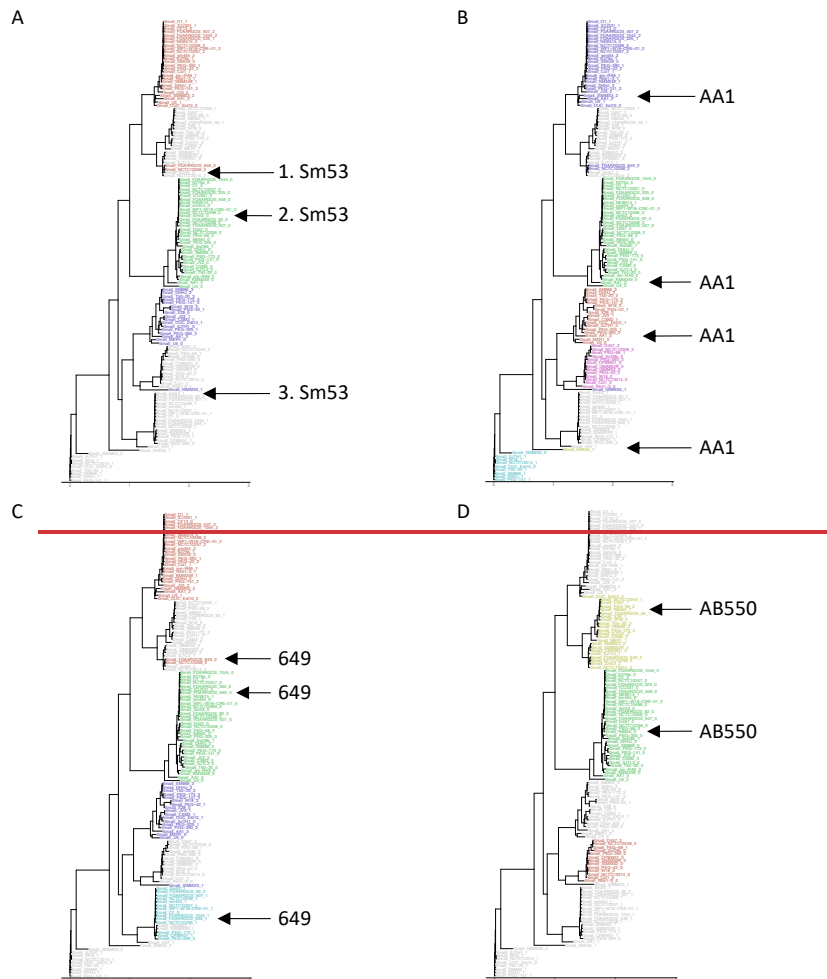
273

274 **Results**

275 RAREFAN can identify REPINs and their corresponding RAYTs in a set of ~~ideally~~ fully sequenced
276 ~~–~~ bacterial genomes. The RAREFAN algorithm has been used in previous analyses to identify and
277 characterize REPINs and RAYTs in ~~*Pseudomonas* (Bertels, Rainey 2011, 2022), *Neisseria* (Bertels,~~
278 ~~*Rainey 2022*), and *Enterobacteria* (Bertels, Gallie, et al. 2017; Park et al. 2021)~~[Pseudomonads](#)
279 [\(Bertels, Rainey 2011a, 2022\)](#), [Neisseriads \(Bertels, Rainey 2022\)](#), and [Enterobacteria \(Bertels,](#)
280 [Gallie, et al. 2017; Park et al. 2021\)](#). To demonstrate RAREFAN’s capabilities, we are presenting
281 an analysis of 49 strains belonging to the opportunistic pathogen *S. maltophilia*.

282 *S. maltophilia* strains contain Group 3 RAYTs, which are also commonly found in plant-associated
283 *Pseudomonas* species such as *P. fluorescens* or *P. syringae* (Bertels, Rainey 2011, 2022). Similar
284 to Group 3 RAYTs in other species, *S. maltophilia* contains multiple REPIN-RAYT systems per
285 genome. Group 2 RAYTs, in contrast, contain only ever one REPIN-RAYT system per genome
286 (Bertels, Rainey 2022).

287



288

289 *S. maltophilia* strains contain Group 3 RAYTs, which are also commonly found in plant-associated
290 *Pseudomonas* species such as *P. fluorescens* or *P. syringae* (Bertels, Rainey 2011a, 2022). Similar
291 to Group 3 RAYTs in other species, *S. maltophilia* contains multiple REPIN-RAYT systems per
292 genome. Group 2 RAYTs, in contrast, contain only ever one REPIN-RAYT system per genome
293 (Bertels, Rainey 2022).

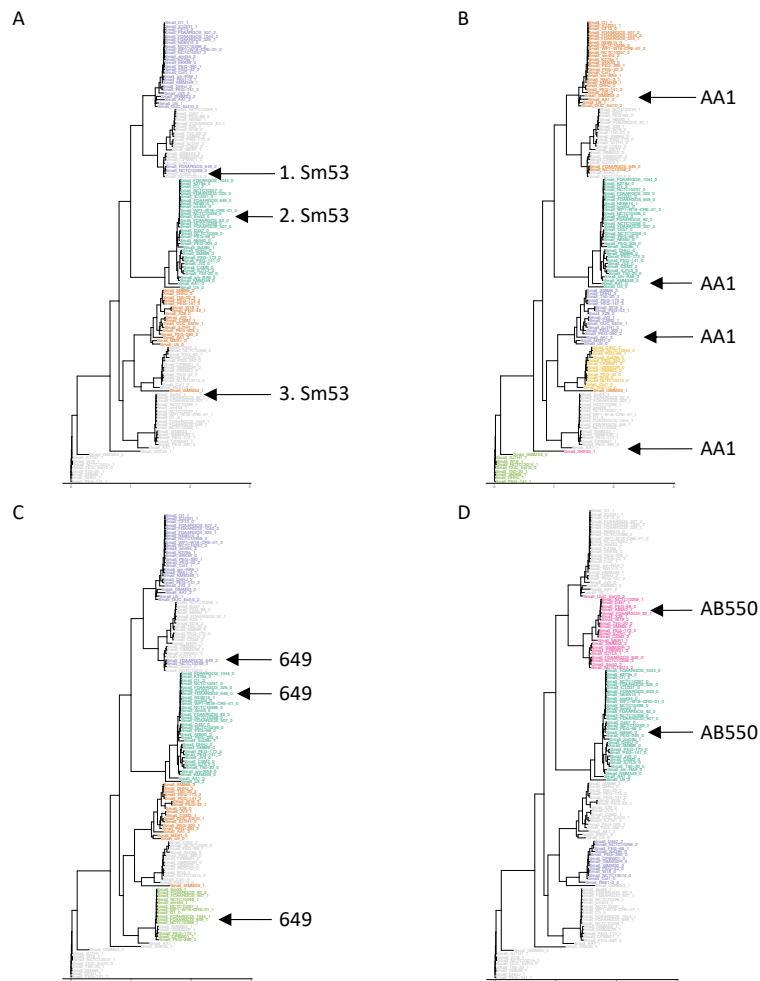


Figure 23. Phylogenetic trees built from RAYT genes extracted from *S. maltophilia* genomes. RAYT genes are coloured according to their association with REPIN populations in the reference genome. If ~~no~~ REPIN population of a query genome is not present in the reference genome ~~is~~ associated with a, then the REPIN population cannot be identified in the query genome and the corresponding RAYT gene the RAYT gene cannot be linked and is coloured in grey. The four panels **A-D** show phylogenies for four different reference strains. *S. maltophilia* strains Sm53, AA1, 649 and AB550 were used in panels **A** to **D**, respectively. Locations of a reference strain's RAYT genes in the tree are indicated by arrows. An association between almost all RAYTs and REPIN populations could be made by using four different reference genomes. Most of the RAYT genes are coloured (associated to a REPIN group) in at least one of the trees. The three numbered RAYT genes from the Sm53 RAREFAN run are referenced in the text.

297 *Nine different REPIN-RAYT systems in S. maltophilia*

298 REPIN-RAYT systems in *S. maltophilia* are surprisingly diverse compared to other species. For
299 example, *P. Pseudomonas chlororaphis* contains three separate REPIN populations that are
300 present in all *P. chlororaphis* strains, each associated with its cognate RAYT gene (Bertels, Rainey
301 2022). *S. maltophilia*, in contrast, contains only one REPIN-RAYT system
302 that is present across almost the entire species (green clade in **Figure 23**), and at least eight
303 REPIN-RAYT systems that are present in subsets of strains (nine clades in **Figure 45**).

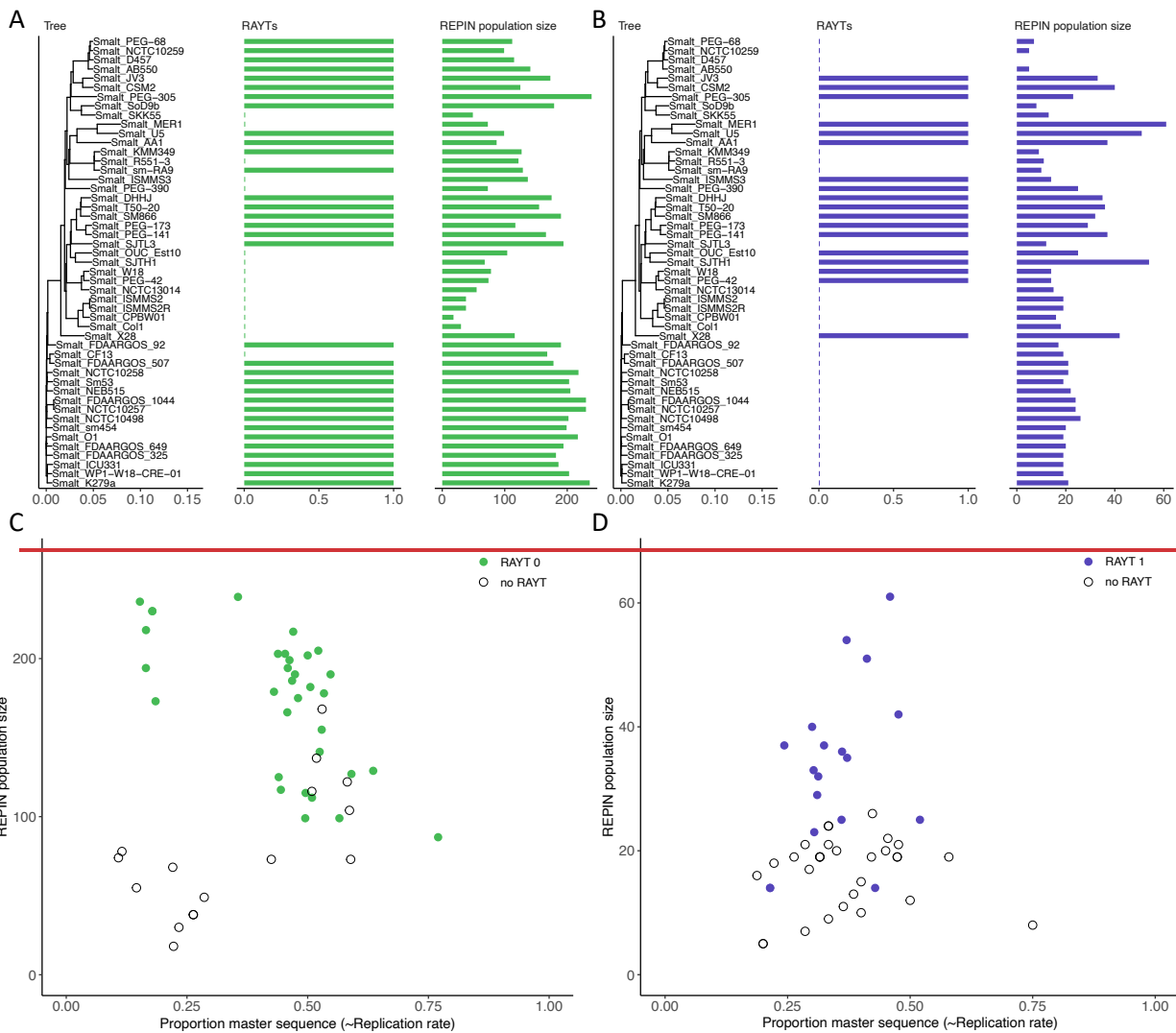
304 The patchy presence-absence pattern of REPIN-RAYT systems in *S. maltophilia*, makes the dataset
305 quite challenging to analyse. If a REPIN population is not present in the reference strain then
306 RAREFAN will not be able to detect it in any other strain. Yet, it is possible to detect RAYT genes
307 in all strains of a species independent of the reference strain selection. RAYT genes that are not
308 associated to a REPIN population are displayed in grey (**Figure 2A3A**). While these RAYT genes
309 are not associated to REPIN populations detected in the reference strain, they might still be
310 associated with a yet unidentified REPIN type— present in the genome the unassociated RAYT
311 gene is located in.

312 In order to identify all REPIN populations across a species, we suggest to perform multiple
313 RAREFAN runs with different reference strains. The RAREFAN web interface supports re-
314 launching a given job with modified parameters.

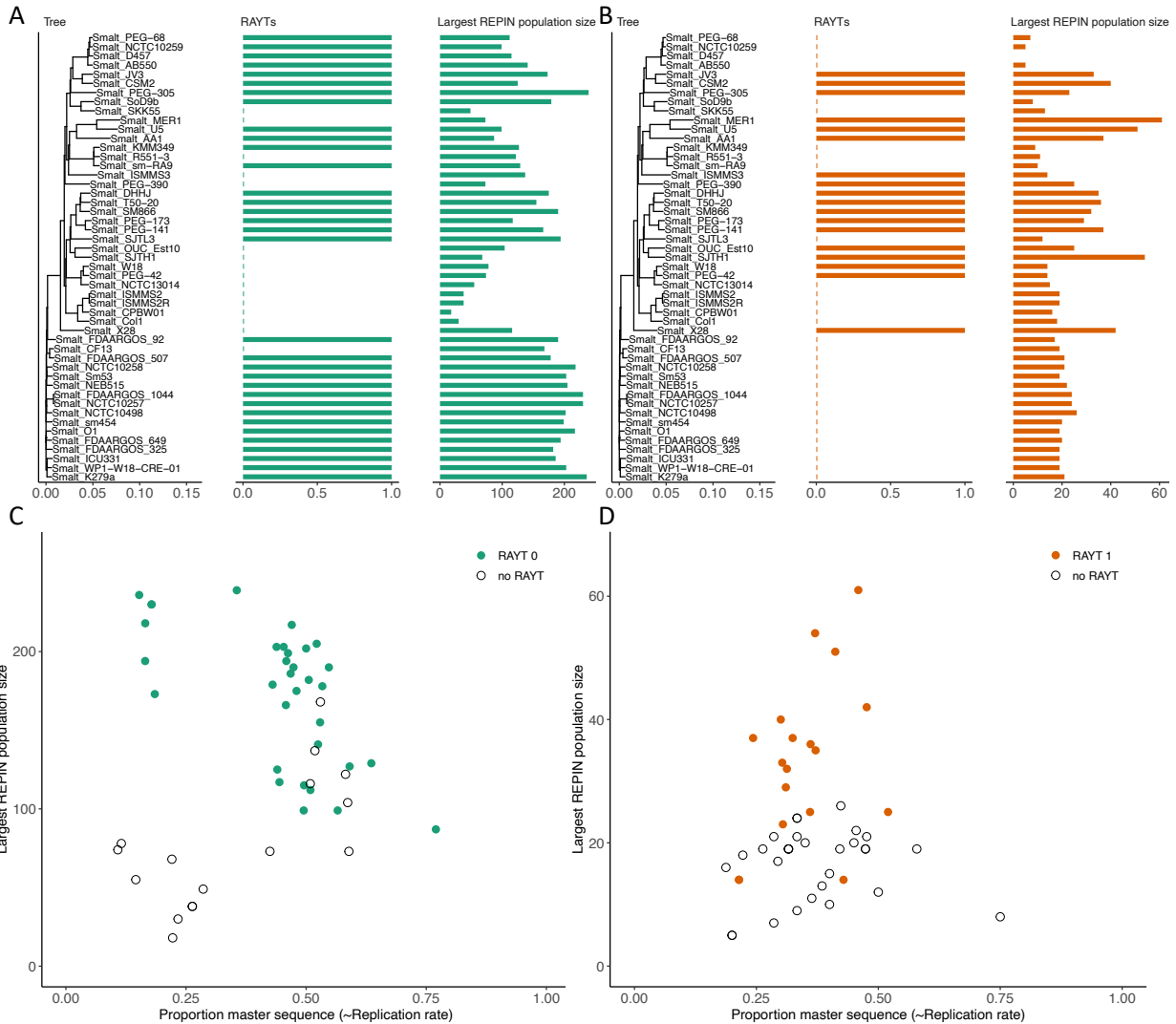
315 To identify as many different REPIN-RAYT systems as possible in each subsequent run the
316 reference should be set to a genome that contains RAYTs that were not associated with a REPIN
317 population previously (i.e., genomes containing grey RAYTs in **Figure 23**). However, this strategy
318 may also fail when the REPIN population size falls below the RAREFAN seed sequence frequency
319 threshold.

320 For example, *S. maltophilia* Sm53 contains three RAYTs only one of which is associated with a
321 REPIN population (RAYT genes indicated in **Figure 2A3A**). However, the remaining two RAYTs are
322 indeed associated with a REPIN population, but these REPIN populations are too small to be
323 detected in *S. maltophilia* Sm53 (the seed sequence frequency threshold is set to 55 by default).
324 In other *S. maltophilia* strains the REPIN populations are large enough to exceed the threshold.

325 For example, if *S. maltophilia* AB550 is set as reference, RAYT number 1 from Sm53 (**Figure 2A3A**)
 326 is associated with the **yellowpink** REPIN population (**Figure 2D3D**). If *S. maltophilia* 649 is set as
 327 reference RAYT number 3 from Sm53 (**Figure 2A3A**) is associated with the **turquoise/light green**
 328 REPIN population (**Figure 2C3C**). RAYTs from the bottom clade are only associated with REPIN
 329 populations when *S. maltophilia* AA1 is chosen as reference (**Figure 2B-3B**). While lower
 330 thresholds can guarantee that all REPINs will be identified in the genome, the number of
 331 sequence groups that are not REPINs quickly explodes. Especially for genomes that contain large
 332 numbers of mobile genetic elements or CRISPRs (Bertels, Rainey 2022).



333



334

Figure 34. REPIN population sizes and conservation. The plots show two REPIN populations and their associated RAYTs that were identified in *S. maltophilia* using *S. maltophilia* Sm53 as reference. **(A)** The phylogenetic tree on the left side is a whole genome phylogeny generated by andi (Haubold *et al.* 2015). Shown on the right are REPIN population sizes (which is the largest REPIN cluster calculated by MCL) and the number of associated RAYTs sorted by the genome phylogeny. The green REPIN populations and associated RAYTs are present in most strains in high abundance (maximum of 239 occurrences in *S. maltophilia* K279a, left panel). **(B)** The blueorange population in contrast is present in much lower numbers (maximum of 61 occurrences in *S. maltophilia* MER1, right panel). Note, REPIN populations are assigned consistent colours based on their abundance in the reference genome. For example, the most abundant REPIN population in the reference is always coloured in green, and the second most abundant population is always coloured in blueorange. **(C and D)** Proportion of master sequence in *S. maltophilia* REPIN populations. The master sequence in a REPIN population is the most common REPIN sequence. TheIn an equilibrium the higher the proportion of the master sequence in the population the higher the replication rate (Bertels, Gokhale, *et al.* 2017). The presence and absence of an

associated RAYT is also indicated by the colours of the dots. Empty circles indicate that the REPIN population is not associated with a RAYT gene in that genome.

335 *RAREFAN visualizes REPIN population size and potential replication rate*

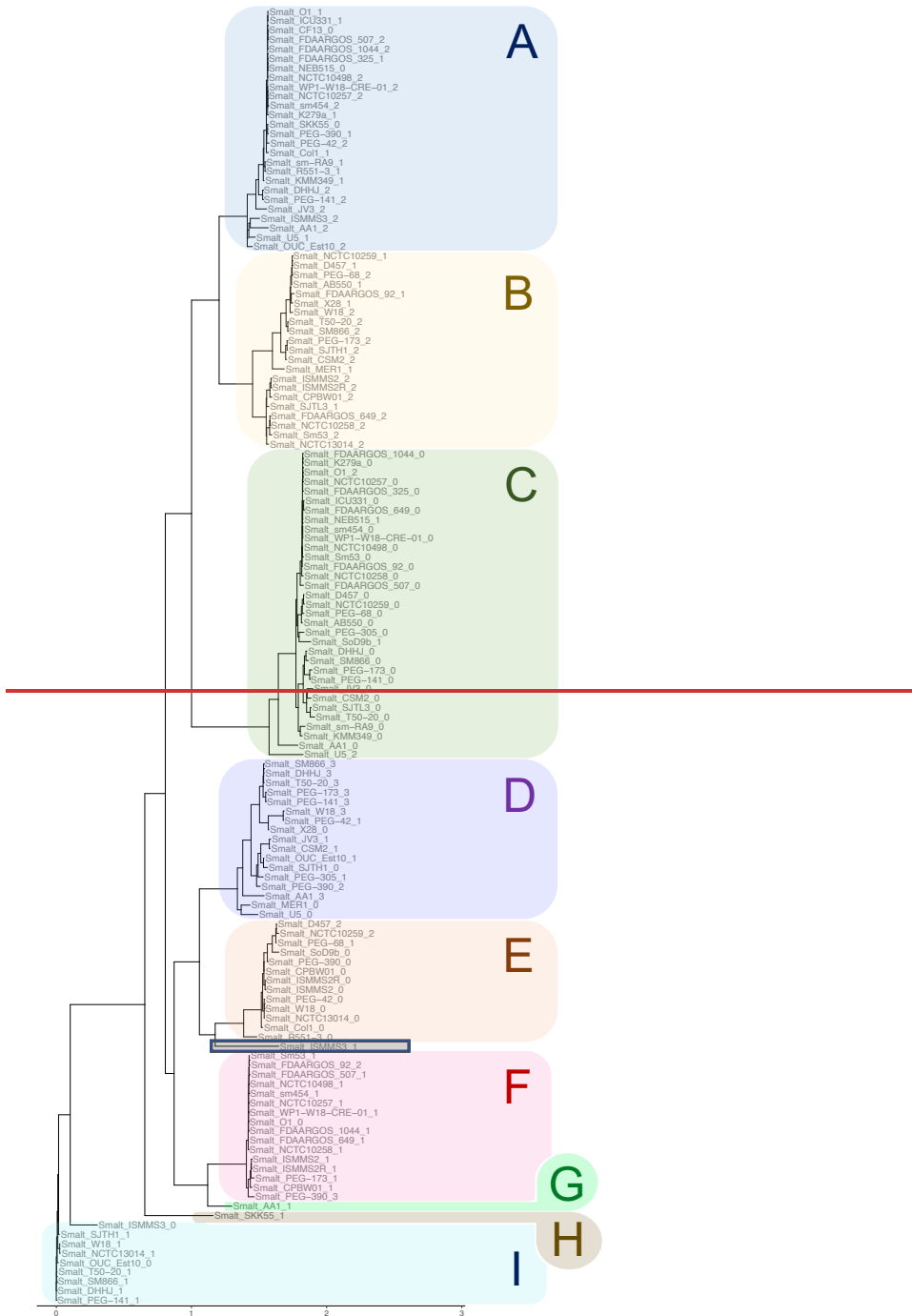
336 The RAREFAN webserver visualizes REPIN population size and RAYT numbers in barplots. ~~The~~
337 ~~barplot is~~Barplots are ordered by the phylogenetic relationship of the submitted bacterial strains
338 (Yu *et al.* 2018). RAREFAN detects three populations when *S. maltophilia* Sm53 is selected as
339 reference strain (**Figure 2A3A**). The largest REPIN population (calculated by MCL from all REPINs
340 of that type) has a corresponding RAYT gene in almost all strains (first barplot in **Figure 3A4A**)
341 and most REPIN populations contain more than 100 REPINs (second barplot in **Figure 3A4A**). The
342 second largest REPIN population in Sm53 (purpleorange population in **Figure 3B4B**) is
343 significantly smaller and contains no more than 61 REPINs in any strain and most strains do not
344 contain a corresponding RAYT for this population.

345 RAREFAN also provides information on REPIN replication rate (**Figure 3C4C** and **D**). REPIN
346 replication rate can be estimated by dividing the number of the most common REPIN sequence
347 (master sequence) by the REPIN population size if the population is in mutation selection balance
348 (Bertels, Gokhale, *et al.* 2017). If a REPIN replicates very fast most of the population will consist
349 of identical sequences because mutations do not have enough time to accumulate between
350 replication events. Hence, the proportion of master sequences will be high in populations that
351 have a high replication rate. Transposable elements such as insertion sequences consist almost
352 exclusively of identical master sequences because the time between replication events is not
353 sufficient to accumulate mutations and because quick extinction of the element usually prevents
354 the accumulation of mutations after replication (~~Park *et al.* 2021; Bertels, Rainey 2022~~)(Park *et*
355 *al.* 2021; Bertels, Rainey 2022). REPIN populations in contrast replicate slowly and persist for long
356 periods of time, which means that a high proportion of master sequences suggests a high REPIN
357 replication rate.

358 In *S. maltophilia* the proportion of master sequences in the population does not seem to correlate
359 well with REPIN population size, both in the green and the purpleorange population (**Figures**
360 **3C4C** and **D**). Similar observations have been made in *P. chlororaphis* (~~Bertels, Rainey~~
361 ~~2022~~)(Bertels, Rainey 2022), and may suggest that an increase in population size is not caused by

362 an increase in replication rate. Population size is likely to be more strongly affected by other
363 factors such as the loss of the corresponding RAYT gene, which leads to the decay of the REPIN
364 population. One could even speculate that high REPIN replication rates are more likely to lead to
365 the eventual demise of the population due to the negative fitness effect on the host (Bertels,
366 Rainey 2022)(Bertels, Rainey 2022).

367 ~~The presence of RAYTs and the size of the corresponding REPIN population do correlate~~
368 ~~surprisingly well. RAYTs are absent from an entire *S. maltophilia* clade (middle of Figure 3A). This~~
369 ~~clade has also lost a significant amount of green REPINs, and the proportion of the master~~
370 ~~sequences in these populations is low (Figure 3C). Similarly, genomes without RAYTs have smaller~~
371 ~~REPIN populations in the purple population than genomes with the corresponding RAYT (Figure~~
372 ~~3D). A similar observation has been made previously in *E. coli*, *P. chlororaphis* and *Neisseria*~~
373 ~~where the loss of the RAYT gene is followed by a decay of the REPIN population (Bertels, Rainey~~
374 ~~2022).~~



375

376

377

378

379

380

Figure 4 The presence of RAYTs and the size of the corresponding REPIN population do correlate surprisingly well (**Figure 4A and B**, p-Value = 0.008 of the linear model of independent contrasts (Felsenstein 1985) of green RAYT and REPIN number, p-Value = 0.003 for orange REPIN populations). Green RAYTs are absent from an entire *S. maltophilia* clade (middle of **Figure 4A**). This clade has also lost a significant amount of green REPINs, and the proportion of the master

381 sequences in these populations is low (**Figure 4C**). Similarly, genomes without orange RAYTs have
 382 smaller REPIN populations in the orange population than genomes with the corresponding RAYT
 383 (**Figure 4D**). A similar observation has been made previously in *E. coli*, *P. chlororaphis*, *N.*
 384 *meningitidis* and *N. gonorrhoeae* where the loss of the RAYT gene is followed by a decay of the
 385 associated REPIN population (Park *et al.* 2021; Bertels, Rainey 2022).

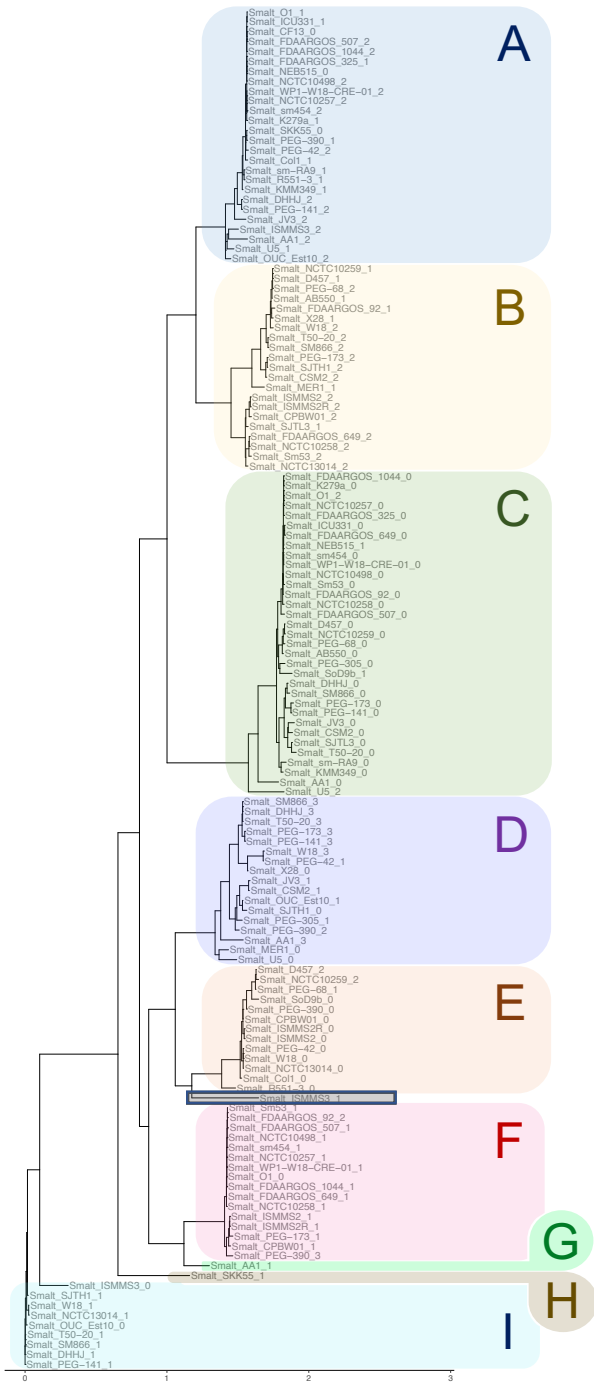


Figure 5. Phylogeny of RAYT genes and their associated REPINs. The tree shows RAYT genes from 49 *S. maltophilia* strains. Colours of clades A-I are assigned according to their association with a REPIN found within 130 bp ~~to of~~ the RAYT gene (see **Table 2**). Except for a single RAYT gene ISMMS3_1 (grey box), which could not be linked to a REPIN population.

387 **Table 2. REPIN palindromes associated with each RAYT clade.**

RAYT population	REPIN palindromes
A	CCGACCAACGGTCGG
B	CCAACCAAGGTTGGC
C	CCGGCC AG CGGCCGG
D	TCCACGC AT GGCGTGGA
E	CCGAGCC AT GCTCGG
F	TCGACT AA CAGTCGA
G	TCGACCAACGGTCGA
H	GCCGGGC AT GGCCCGGC
I	AGTCGAGCTTGCTCGACT

388 Each RAYT clade from **Figure 45** is associated with a unique imperfect palindrome that is present
 389 at the 5' and/or 3' end of the ~~REPIN (Figure 5)-RAYT gene.~~

390

391 *Linking REPIN populations with RAYT genes can be challenging*

392 Unfortunately, RAREFAN is not always able to link the correct REPIN population with the correct
 393 RAYT gene. ~~As shown in Figure 2 in~~ some RAREFAN runs ~~indicate that sometimes~~ associations
 394 between RAYTs and REPINs ~~seem to are~~ not ~~be~~ monophyletic, as for example ~~the red~~ RAYTs
 395 ~~highlighted in red~~ in **Figure 2A3A**. However, the same clade of RAYTs is uniformly coloured in
 396 yellow in **Figure 2D3D**, suggesting that the entire RAYT clade is associated with the same REPIN
 397 group. ~~To investigate the true relationships between REPINs and RAYTs we first mapped REPIN~~
 398 ~~groups to RAYT genes.~~

399 An analysis of all REPIN groups that were identified by RAREFAN across ~~fivefour~~ different
 400 RAREFAN runs (**Table 1, one additional analysis was performed with ISMMS3**) showed that there
 401 are a total of nine different REPIN groups, each defined by an individual central palindrome (**Table**

402 2). Each REPIN group is associated with a ~~mono-phyletic~~monophyletic RAYT group (**Figure 4**).
403 ~~There is only~~5). Only a single RAYT ~~that we could is not identify~~associated with a REPIN
404 ~~for~~population (ISMMS3_1). ~~Although there seems to be a one to one mapping between RAYT~~
405 ~~clades and REPIN groups, the question remains why RAREFAN sometimes links RAYT genes to the~~
406 ~~wrong REPIN group.~~

407 RAREFAN could not link a REPIN to the RAYT gene ISMMS3_1 (**Figure 5**, grey box). While there is
408 a sequence that resembles the A detailed analysispalindrome as well as variants of the extragenic
409 spaceC palindrome flanking both sides of the RAYT gene (**Supplementary Figure 2**), none of the
410 sequences formed REPIN populations large enough to be identified by RAREFAN. Presumably the
411 RAYT ISMMS3_1, which is only present in a single *S. maltophilia* strain, is at the early stages of
412 establishing a REPIN population~~“wrongly” associated,~~ and the REPIN population has not spread
413 to a considerable size yet.

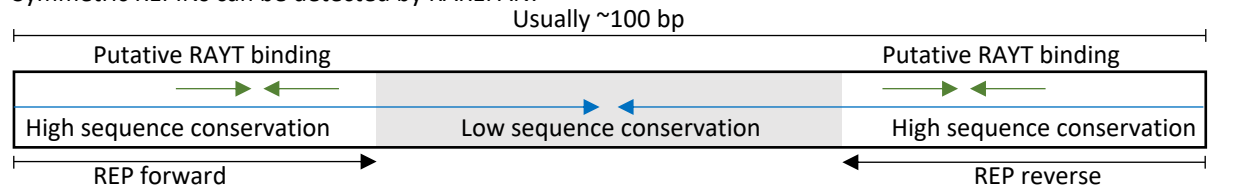
414 There are two more cases where RAREFAN failed to link RAYT genes with any REPINs (ISSMS2
415 and ISSMS2R_1, **Supplementary Figure 1 D and E**). Detailed sequence analyses showed that the
416 respective REPINs are located at a distance of more than 130 bp (an adjustable parameter in
417 RAREFAN). These REPINs are ignored by RAREFAN by default. However, this parameter can be
418 adjusted manually and when set to a distance of 200 bp, RAREFAN correctly links these
419 genes~~REPINs~~ to the RAYT gene.

420 In three cases the wrong REPIN population was linked to a RAYT gene. In our dataset this can
421 happen when RAYTs are flanked by seed sequences from two different REPIN populations
422 (**Supplementary Figure 1 A-C**). A single REP sequence from the “wrong” (non-monophyletic
423 RAYT) clade occurs together with multiple REP or REPIN sequences from the “right”
424 (monophyletic in a different RAREFAN run) clade. REPINs are linked to the “wrong” RAYT when
425 the correct REPIN population is absent in the chosen reference genome. ~~Moreover, the “wrongly”~~
426 ~~associated REP singlets always show up as belonging to the REPIN population of a RAYT sister~~
427 ~~group and show high sequence similarity~~This problem can be alleviated by performing analyses
428 with multiple reference genomes and comparing the results.

429 Additionally to linking the wrong REPINs and RAYTs, RAREFAN sometimes failed to link RAYT
 430 genes with any REPINs. Detailed sequence analyses showed that in two out of a total of three
 431 such cases the REPIN was located at a distance of more than 130bp (a parameter set in RAREFAN
 432 that could be modified) (**Supplementary Figure 1 D-E**). REPINs that are located at a distance of
 433 more than 130bp are ignored by RAREFAN. In a third case there was no REPIN that could be linked
 434 to the RAYT gene ISMMS3_1 (**Figure 4**, grey box). While there is a sequence that resembles the
 435 A-palindrome as well as variants of the C-palindrome flanking both sides of the RAYT gene,
 436 (**Supplementary Figure 2**), none of the sequences formed REPIN populations large enough to be
 437 identified by RAREFAN. Presumably the RAYT ISMMS3_1, which is only present in a single *S.*
 438 *maltophilia* strain, is at the early stages of establishing a REPIN population. Based on our findings,
 439 RAREFAN users should always critically analyse RAREFAN results, particularly when the results
 440 require unusual evolutionary explanations.

441

Symmetric REPINs can be detected by RAREFAN:



Asymmetric REPINs cannot be detected by RAREFAN:



442

Figure 5. The structure of symmetric and asymmetric REPINs. A typical REPIN consists of two highly conserved regions at the 5' and 3' end of the REPIN (white), separated by a spacer region of lower sequence conservation (grey). The entire REPIN is a palindrome (blue arrows), which means it can form hairpin structures in single stranded DNA or RNA. Each 5' and 3' region contains a nested imperfect palindrome, which is referred to as REP (repetitive extragenic palindromic) sequence and has first been described in *E. coli* (Higgins *et al.* 1982). REPINs can be either symmetric or asymmetric. Asymmetric REPINs have a deletion and a corresponding insertion in the highly conserved 5' or 3' end, which leads to “bubbles” in the hairpin structure. REPINs in *E. coli* are asymmetric, which makes analyses with RAREFAN more challenging. Figure adapted from (Bertels, Rainey 2022).

443

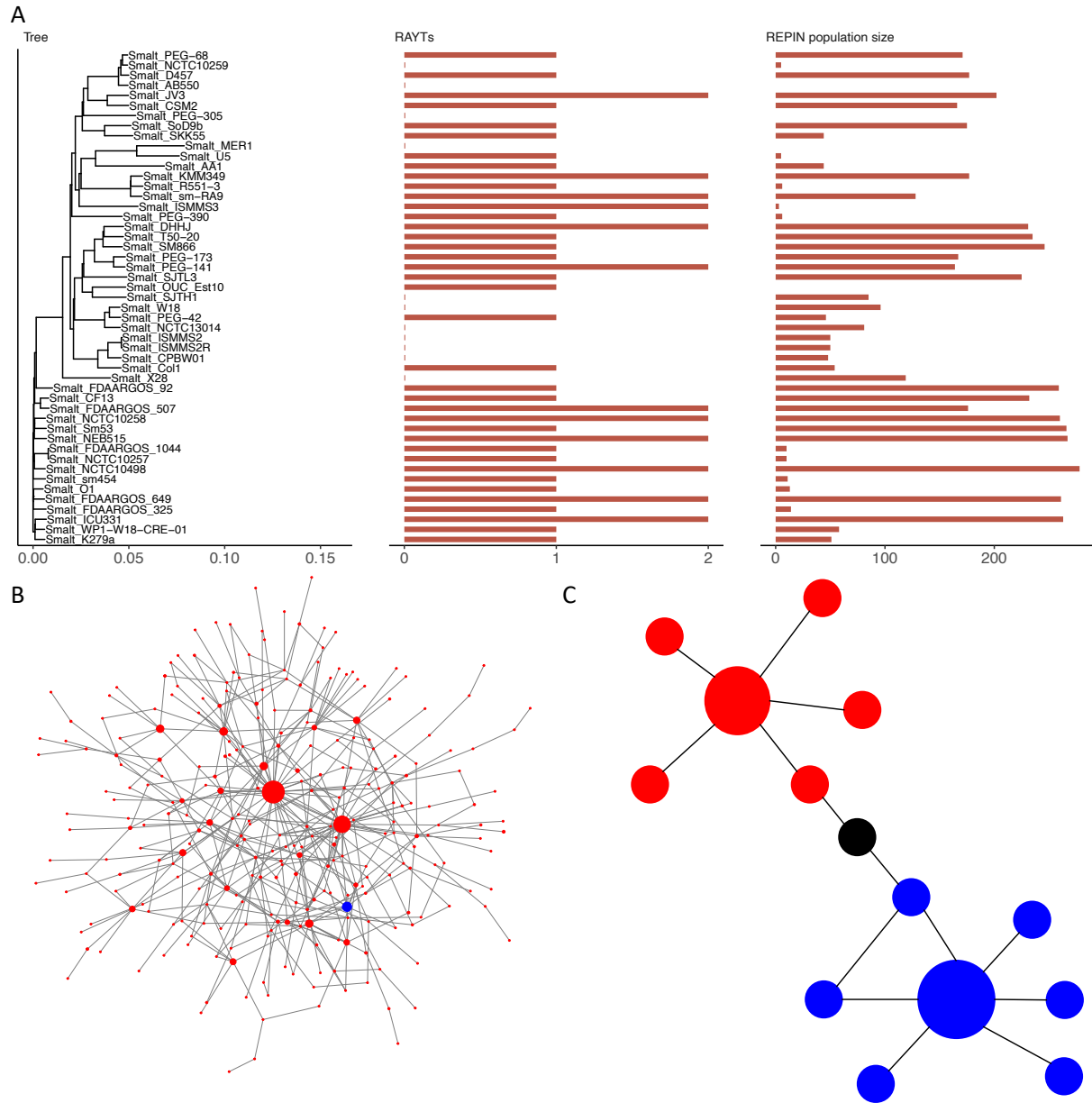
444 *REPIN groups may be lost when the seed distance is too large*

445 The seed distance parameter determines whether two highly abundant sequences are sorted
446 into the same or different REPIN groups (**Figure 1B2B**). If two REPINs from two different groups
447 occur next to each other, at a distance of less than the seed distance parameter, then the two
448 seeds are erroneously sorted into the same group. If two different REPIN groups are sorted into
449 the same group then one of the groups will be ignored by RAREFAN, because only the most
450 abundant seed in each group will be used to identify REPINs.

451 A manual analysis (*e.g.*, multiple sequence alignment) of sequences in the groupSeedSequences
452 folder of the RAREFAN output can identify erroneously merged REPIN groups. In *S. maltophilia*,
453 groups are separated well when the distance parameter is set to 15 bp and Sm53 is used as a
454 reference. When the parameter is set to 30 bp instead, one of the REPIN groups will be missed
455 by RAREFAN.

456 A small seed distance parameter will separate seed sequences belonging to the same REPIN
457 group into different groups. Hence, RAREFAN will analyse the same REPIN group multiple times.
458 While this will lead to increased RAREFAN runtimes, these errors, are easy to spot, because (1)
459 the same RAYT gene will be associated to multiple REPIN groups, (2) the central palindrome
460 between the group is identical and (3) the master sequence between the groups will be very
461 similar.

462



463

Figure 6. Closely related REPIN populations may be merged by RAREFAN. (A) REPIN group 2 identified in a *S. maltophilia* Sm53 RAREFAN run. The RAREFAN result suggests that REPIN group 2 is sometimes associated with two RAYTs. **(B)** A closer inspection of the data shows that Group 2 is a combination of two different REPIN groups, the “real” Group 2 and Group 0. The network shown, visualizes all REP sequences identified as Group 2. Nodes in the network represent 21 bp long REP sequences. Two nodes are connected if the sequences they represent differ by exactly one nucleotide. The node size indicates the abundance of the sequence in the genome. The blue node represents the most common Group 2 sequence, occurring 65 times in the genome. The largest red node occurs 407 times in the genome and resembles a Group 0 REP sequence. **(C)** Illustration of how small changes to a single sequence can connect two sequence cluster. The most common 21 bp long sequence in Group 0 differs in only four positions from the

most common 21 bp long sequence in Group 2. There is a set of sequences that connects these two groups that only differ in exactly one position each (nodes connected by an edge), which passes through the black node. If there is such an unbroken path between REP sequences, then REPIN groups will be merged.

464

465 *Closely related REPIN groups may be merged into a single group by RAREFAN*

466 Incorrect merging of REPIN groups can occur when two REPIN groups are closely related. We
467 identified merged REPIN groups in *S. maltophilia* because RAREFAN linked some REPIN groups
468 with two RAYT genes in the same genome (**Figure 6A**). ~~While REPIN groups linked to two RAYTs~~
469 ~~has been observed before in *Neisseria meningitidis* (Bertels, Rainey 2022), it is particularly~~
470 ~~unusual in *S.*~~ While REPIN groups linked to two RAYTs has been observed before in *Neisseria*
471 *meningitidis* (Bertels, Rainey 2022), it is particularly unusual in *S. maltophilia* due to some key
472 differences between REPIN-RAYT in the two bacterial species. First, RAYTs in *N. meningitidis*
473 belong to Group 2 and RAYTs in *S. maltophilia* belong to Group 3 (Bertels, Gallie, *et al.* 2017), two
474 very divergent RAYT groups. Second, RAYTs that are associated to the same REPIN group in *N.*
475 ~~*meningitidis* are almost identical, since they are copied by an insertion sequence in trans (Bertels,~~
476 ~~*Rainey 2022)*~~ *meningitidis* are almost identical, since they are copied by an insertion sequence in
477 *trans* (Bertels, Rainey 2022), something that is not the case for *S. maltophilia*, where the two
478 RAYTs are very distinct from each other (green and red clade in **Figure 2A3A**, or clade A and C in
479 **Figure 45**).

480 A closer inspection of all sequences identified in REPIN group 2 shows that it also contains
481 sequences belonging to REPIN group 0 (palindromes linked to clade A and C in **Table 2**). The
482 relationship between the sequences shows that there is a chain of sequences that all differ by at
483 most a single nucleotide between the most abundant sequence in group 2 to the most abundant
484 sequence in group 0 (**Figure 6B and C**). Hence, the reason group 0 and group 2 are merged is that
485 they are too closely related to each other and hybrids of the two REPIN groups exist. Because
486 sequence groups are built by identifying all related sequences in the genome recursively, closely
487 related groups (the REPIN group 0 seed only differs in four nucleotides from the REPIN group 2
488 seed sequence) can be merged into a single REPIN group. REPIN population size and RAYT number
489 are the sum of REPIN group 0 and 2. There are various possibilities to resolve this issue: (1)

490 subtract sequences from group 0 (which does not contain group 2) from REPIN group 2; (2) use
491 a different sequence seed from the group 2 seed collection in the seed sequence file
492 (groupSeedSequences/Group_SmaI_t_Sm53_2.out); ~~or~~ (3) sometimes it may be possible
493 to rerun RAREFAN with a different reference strain where the issue does not occur; or (4) increase
494 the length of the seed sequence.

495

496 Performance

497 We measured the elapsed time for a complete RAREFAN run for three different species and for
498 5, 10, 20, and 40 genomes with randomly selected reference strains and the two query RAYTs
499 (yafM Ecoli and yafM SBW25). For a given number N of submitted genomes of average
500 sequence length L (in megabases), a RAREFAN run completes in approximately $T = (8-10 \text{ seconds})$
501 $* N * L$ on our moderate server hardware (4 CPU cores, 16 GB shared RAM) (**Supplementary**
502 **Figure 3 and 4**).

503 **Discussion**

504 RAREFAN allows users to quickly detect REPIN populations and RAYT transposases inside
505 bacterial genomes. It also links the RAYT transposase genes to the REPIN population it duplicates.
506 These data help the user to study REPIN-RAYT dynamics in their strains of interest without a
507 dedicated bioinformatician, and hence will render REPIN-RAYT systems widely accessible.

508 One limitation of RAREFAN is that REPINs can only be identified in genomes when they are
509 symmetric (**Figure 51**). Symmetric REPINs have seed sequences that can morph into each other
510 by a series of single substitutions (intermediate sequences need to be present in the genome). A
511 REPIN consists of a 5' and a 3' REP sequence. If one of these REP sequences contains an insertion
512 or deletion, which the other REP sequence does not contain then RAREFAN will not recognize the
513 second repeat of the seed sequence. In this case, RAREFAN will not be able to identify REPINs but
514 can still be used to analyze REP singlet populations. To date, the only known asymmetric REPIN
515 ~~population~~ populations are found in *E. coli* ~~REPINs~~. However, it is likely that asymmetric REPINs
516 also exist in other microbial species.

517 RAREFAN sometimes cannot correctly divide REPINs into REPIN groups. Either because REPINs
518 from different groups occur in close proximity in the genome, an issue that can easily be solved
519 by adjusting a RAREFAN parameter, or because two REPIN groups are very closely related (**Figure**
520 **6**). Unfortunately, RAREFAN is not able to automatically detect and resolve the assignment of
521 closely related REPINs into groups yet. Hence it is advisable to manually check associations
522 between REPIN groups and RAYT genes by analyzing the composition of REPIN groups.

523 In the future we aim to make RAREFAN even more versatile and easier to use by, for example,
524 automatically integrating data from public databases such as Genbank, and integrating RAREFAN
525 into workflows such as Galaxy (Afgan *et al.* 2018).

526 RAREFAN makes the study of REPIN-RAYT systems more accessible to any biologist or
527 bioinformatician interested in studying intragenomic sequence populations. Our tool will help
528 understand the purpose and evolution of REPIN-RAYT systems in bacterial genomes.

529 **Acknowledgements**

530 We would like to thank Prajwal Bharadwaj for assisting us with the sequence analysis and Jenna
531 Gallie for valuable feedback on the manuscript.

532 **References**

- 533 Afgan E, Baker D, Batut B, van den Beek M, Bouvier D, Cech M, Chilton J, Clements D, Coraor N,
534 Grüning BA, Guerler A, Hillman-Jackson J, Hiltemann S, Jalili V, Rasche H, Soranzo N,
535 Goecks J, Taylor J, Nekrutenko A, Blankenberg D (2018) The Galaxy platform for
536 accessible, reproducible and collaborative biomedical analyses: 2018 update. *Nucleic*
537 *Acids Res.*, **46**, W537-W544-W537–W544. <https://doi.org/10.1093/nar/gky379>
- 538 Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ (1990) Basic local alignment search tool.
539 *Journal of Molecular Biology*, **215**, 403–410. <https://doi.org/10.1006/jmbi.1990.9999>
- 540 Arnold K, Gosling J, Holmes D (2005) *The Java programming language*. Addison Wesley
541 Professional.

542 Bertels F, Gallie J, Rainey PB (2017) Identification and Characterization of Domesticated Bacterial
543 Transposases. *Genome Biology and Evolution*, **9**, 2110–2121.
544 <https://doi.org/10.1093/gbe/evx146>

545 Bertels F, Gokhale CS, Traulsen A (2017) Discovering Complete Quasispecies in Bacterial
546 Genomes. *Genetics*, **206**, 2149–2157. <https://doi.org/10.1534/genetics.117.201160>

547 Bertels F, Rainey PB (~~2011~~2011a) Within-Genome Evolution of REPINs: a New Family of Miniature
548 Mobile DNA in Bacteria. *PLoS genetics*, **7**, e1002132.
549 <https://doi.org/10.1371/journal.pgen.1002132>

550 [Bertels F, Rainey PB \(2011b\) Curiosities of REPINs and RAYTs. *Mobile Genetic Elements*, **1**, 262–](#)
551 [268. <https://doi.org/10.4161/mge.18610>](#)

552 Bertels F, Rainey PB (2022) Ancient Darwinian replicators nested within eubacterial genomes. ,
553 2021.07.10.451892. <https://doi.org/10.1101/2021.07.10.451892>

554 Bichsel M, Barbour AD, Wagner A (2010) The early phase of a bacterial insertion sequence
555 infection. *Theoretical Population Biology*. <https://doi.org/10.1016/j.tpb.2010.08.003>

556 Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL (2009) BLAST+:
557 architecture and applications. *BMC Bioinformatics*, **10**, 421–9.
558 <https://doi.org/10.1186/1471-2105-10-421>

559 van Dijk B, Bertels F, Stolk L, Takeuchi N, Rainey PB (2022) Transposable elements promote the
560 evolution of genome streamlining. *Philosophical Transactions of the Royal Society B:*
561 *Biological Sciences*, **377**, 20200477. <https://doi.org/10.1098/rstb.2020.0477>

562 [Edgar RC \(2004\) MUSCLE: multiple sequence alignment with high accuracy and high throughput.](#)
563 [Nucleic Acids Research](#), **32**, 1792–1797. <https://doi.org/10.1093/nar/gkh340>

564 [Felsenstein J \(1985\) Phylogenies and the comparative method. *American Naturalist*, 1–15.](#)

565 Grinberg M (2018) *Flask web development: developing web applications with python*. O'Reilly
566 Media, Inc.

567 [Guindon S, Dufayard J-F, Lefort V, Anisimova M, Hordijk W, Gascuel O \(2010\) New algorithms and
568 \[methods to estimate maximum-likelihood phylogenies: assessing the performance of
569 \\[PhyML 3.0. *Systematic Biology*, 59, 307–321. <https://doi.org/10.1093/sysbio/syq010>\\]\\(#\\)\]\(#\)](#)

570 Haubold B, Klötzl F, Pfaffelhuber P (2015) andi: fast and accurate estimation of evolutionary
571 distances between closely related genomes. *Bioinformatics*, **31**, 1169–1175.
572 <https://doi.org/10.1093/bioinformatics/btu815>

573 Higgins CF, Ames GF, Barnes WM, Clement JM, Hofnung M (1982) A novel intercistronic
574 regulatory element of prokaryotic operons. *Nature*, **298**, 760–762.
575 <https://doi.org/10.1038/298760a0>

576 Initiative TOS (2021) The MIT License.

577 Kearse M, Moir R, Wilson A, Stones-Havas S (2012) Geneious Basic: an integrated and extendable
578 desktop software platform for the organization and analysis of sequence data.

579 Kleinmann SG, Rudolph S, Vila S, Rodin J, Peña JF-S (2021) *The Debian GNU/Linux Operating
580 System Manual*.

581 Lawrence JG, Ochman H, Hartl DL (1992) The evolution of insertion sequences within enteric
582 bacteria. *Genetics*, **131**, 9–20. <https://doi.org/10.1093/genetics/131.1.9>

583 Nunvar J, Huckova T, Licha I (2010) Identification and characterization of repetitive extragenic
584 palindromes (REP)-associated tyrosine transposases: implications for REP evolution and

585 dynamics in bacterial genomes. *BMC Genomics*, **11**, 44. <https://doi.org/10.1186/1471->
586 2164-11-44

587 Park HJ, Gokhale CS, Bertels F (2021) How sequence populations persist inside bacterial genomes.
588 *Genetics*, **217**. <https://doi.org/10.1093/genetics/iyab027>

589 R Core Team (2016) R: A Language and Environment for Statistical Computing. ~~R Foundation for~~
590 ~~Statistical Computing, Vienna, Austria.~~

591 Rankin DJ, Bichsel M, Wagner A (2010) Mobile DNA can drive lineage extinction in prokaryotic
592 populations. *Journal of Evolutionary Biology*. <https://doi.org/10.1111/j.1420->
593 9101.2010.02106.x

594 RStudio, Inc (2013) Easy web applications in R.

595 Rutherford K, Parkhill J, Crook J, Horsnell T, Rice P, Rajandream MA, Barrell B (2000) Artemis:
596 sequence visualization and annotation. *Bioinformatics*, **16**, 944–945.
597 <https://doi.org/10.1093/bioinformatics/16.10.944>

598 Sawyer SA, Dykhuizen DE, DuBose RF, Green L, Mutangadura-Mhlanga T, Wolczyk DF, Hartl DL
599 (1987) Distribution and Abundance of Insertion Sequences Among Natural Isolates of
600 *Escherichia coli*. *Genetics*, **115**, 51–63. <https://doi.org/10.1093/genetics/115.1.51>

601 Ton-Hoang B, Siguier P, Quentin Y, Onillon S, Marty B, Fichant G, Chandler M (2012) Structuring
602 the bacterial genome: Y1-transposases associated with REP-BIME sequences. *Nucleic*
603 *Acids Research*, **40**, 3596–3609. <https://doi.org/10.1093/nar/gkr1198>

604 Van Dongen S (2000) A cluster algorithm for graphs. *Report-Information systems*, 1–40.

605 Van Rossum G, Drake Jr FL (1995) *Python reference manual*. Centrum voor Wiskunde en
606 Informatica Amsterdam.

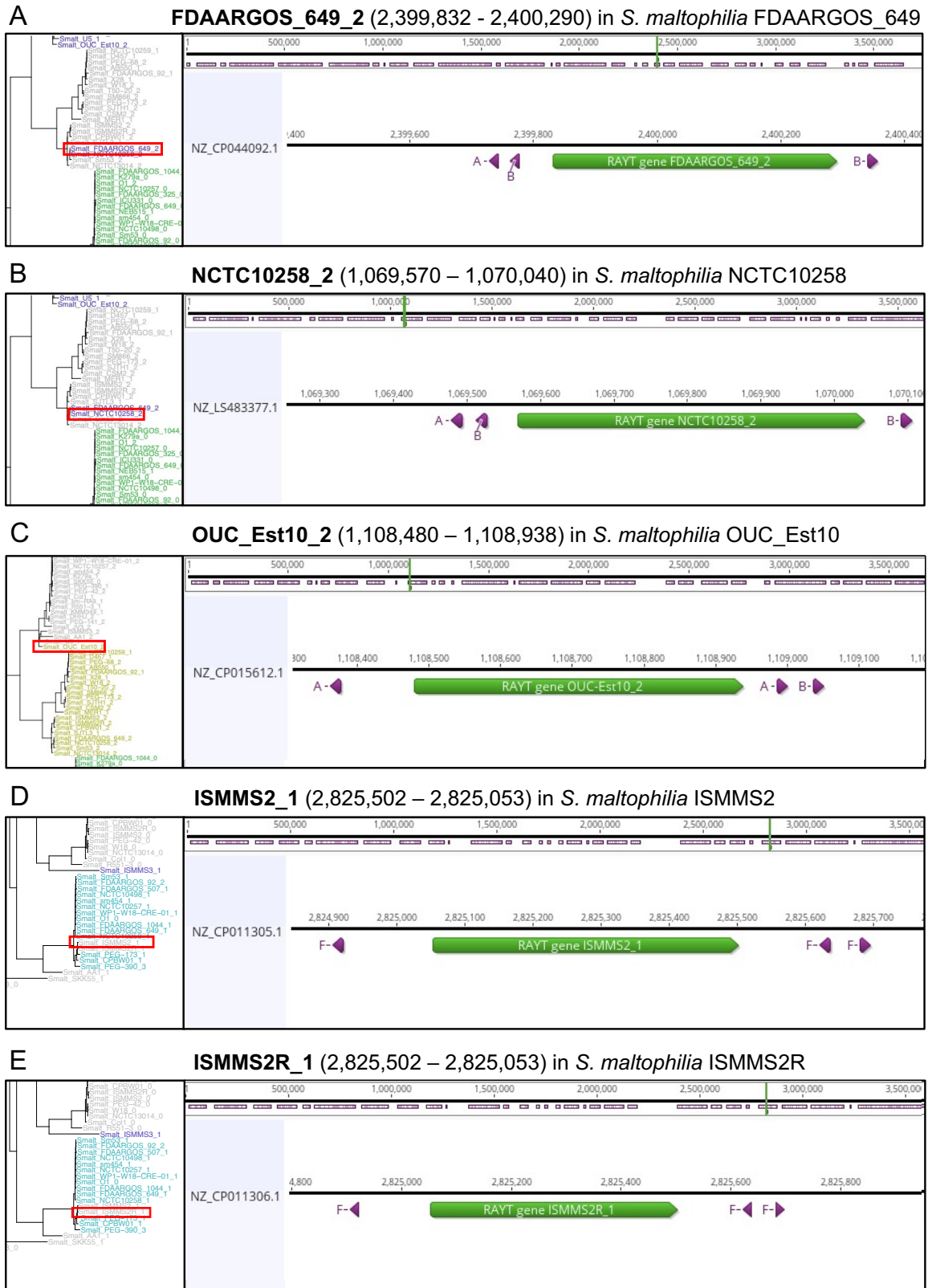
607 Wu Y, Aandahl RZ, Tanaka MM (2015) Dynamics of bacterial insertion sequences: can
608 transposition bursts help the elements persist? *BMC Evolutionary Biology*, **15**, 288–12.
609 <https://doi.org/10.1186/s12862-015-0560-5>

610 Yu G, Lam TT-Y, Zhu H, Guan Y (2018) Two Methods for Mapping and Visualizing Associated Data
611 on Phylogeny Using Ggtree. (FU Battistuzzi, Ed,). *Molecular biology and evolution*, **35**,
612 3041–3043. <https://doi.org/10.1093/molbev/msy194>

613

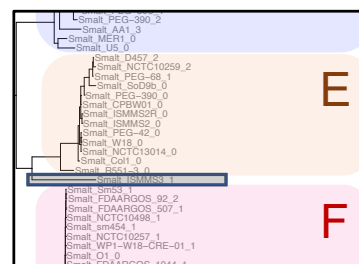
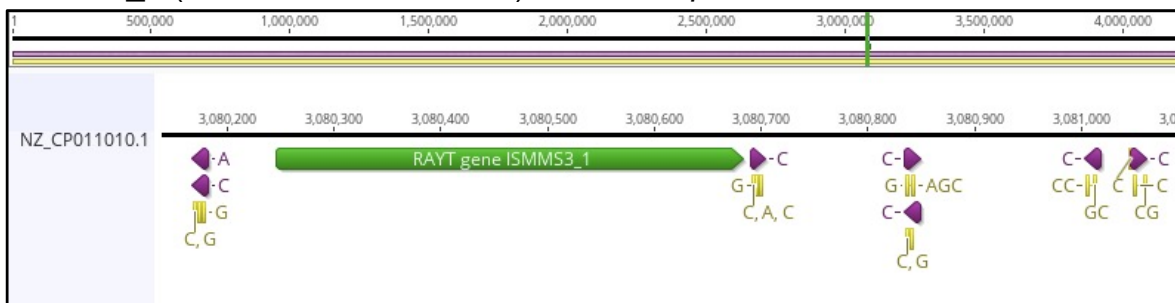
614

615 **Supplementary Figures**



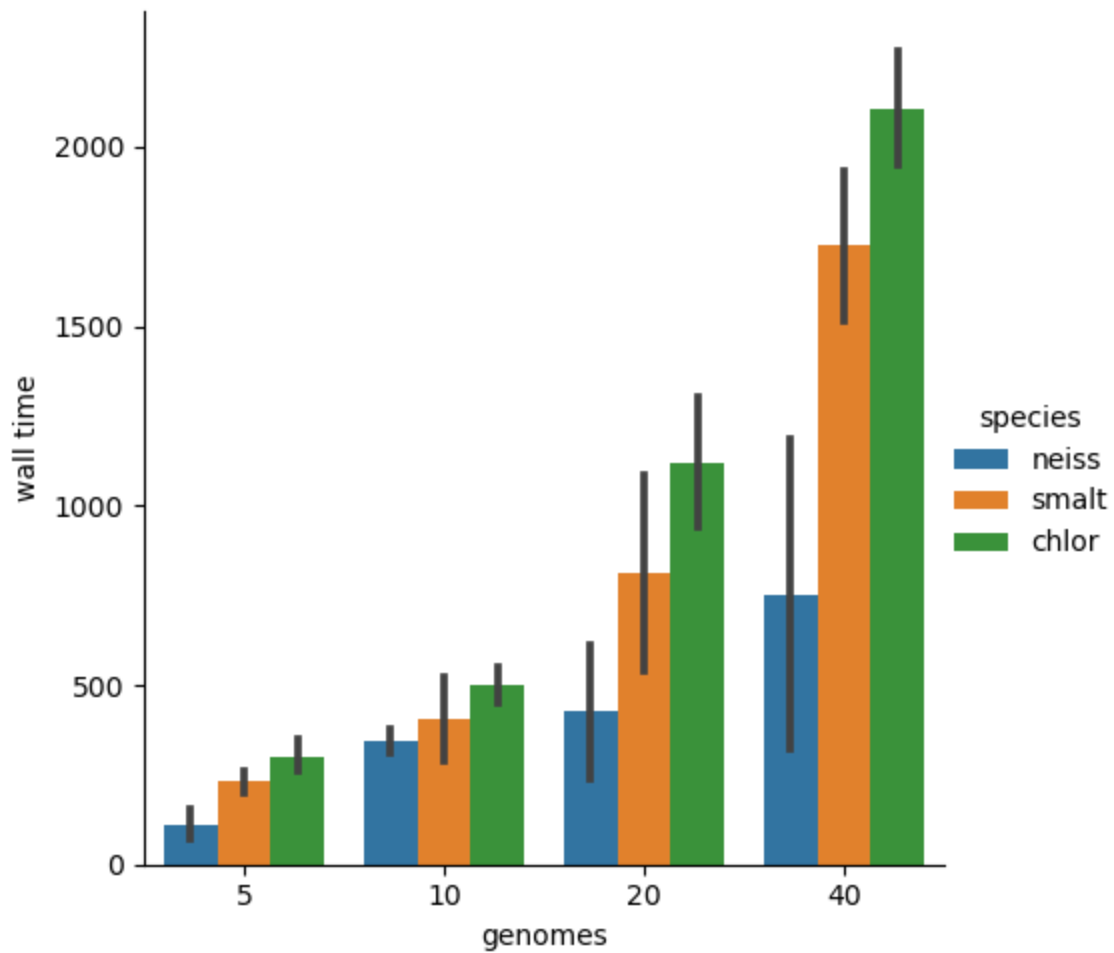
617 **Supplementary Figure 1:** Sequence analysis shows REPIN groups are indeed associated with
 618 **monophyletic RAYTs**. Non-monophyletic or missing associations to REPIN populations identified
 619 by RAREFAN were investigated in the corresponding genomes using Geneious (Kearse *et al.*
 620 2012). Red boxes mark the position of the atypical RAYT that is being analyzed in detail. Mapping
 621 of REPIN palindromes A-I (with zero mismatches) shows FDAARGOS_649_2 (A), NCTC10258_2
 622 (B), and OUC_Est_2 (C) are linked to the wrong REPIN group because REP singlets that are
 623 ordinarily linked to a RAYT sister clade are found in close proximity to the RAYT. These wrong
 624 associations between REPIN and RAYT usually occur when the correct REPIN population is absent
 625 from the reference genome. ISMMS2R_1 (D) and ISMMS2_1 (E) were not linked to REPIN
 626 population by RAREFAN because the corresponding seed sequences were located at a distance
 627 of more than 130 bp from the RAYT gene. Nucleotide sequences and positions were extracted
 628 from output files generated by RAREFAN. Complete genome sequences are available in NCBI
 629 Nucleotide Database using Accessions: (A) NZ_CP044092.1, (B) NZ_LS483377.1, (C)
 630 NZ_CP015612.1, (D) NZ_CP011306.1, (E) NZ_CP011305.1.

ISMMS3_1 (3,080,683 – 3,080,246) in *S. maltophilia* ISMMS3



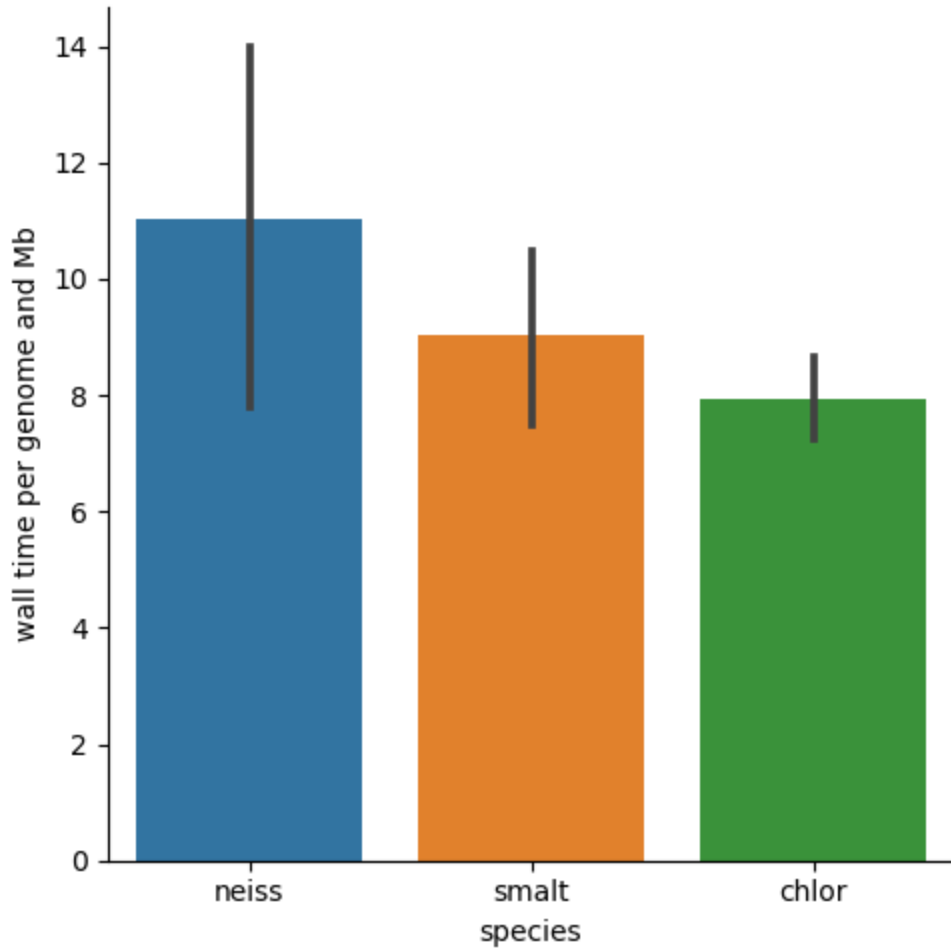
631

632 **Supplementary Figure 2:** RAYT gene ISMMS3_1 cannot be linked to a REPIN population. The
 633 sequence of the RAYT gene ISMMS3_1 and its flanking sequences ~~was analyzed~~ were analysed in
 634 Geneious (Kearse *et al.* 2012). The inset shows the location of ISMMS3_1 in the RAYT phylogeny
 635 (grey box). When mapping all of the identified palindromes to the RAYT region and allowing up
 636 to four mismatches (yellow annotations), various mutants of palindrome C were found in close
 637 proximity of the RAYT gene. However, we could not identify a corresponding REPIN population,
 638 which may indicate that the population has not yet expanded in the genome.



639 Supplementary Figure 3. Average time (in seconds) it takes RAREFAN to complete for different
 640 genome numbers from three bacterial species (*N. meningitidis*, *N. gonorrhoeae*, *S. maltophilia*,
 641 *Pseudomonas chlororaphis*). Black bars indicate the 95% CI across four runs, where two runs
 642 share the same query RAYT. For each run reference and query strains were randomly selected.
 643 All measurements were performed on 4CPU cores with 16 GB of shared memory.
 644

645



646

647 Supplementary Figure 4. Approximate elapsed run time per megabase sequence length
648 calculated from the same runtime data generated in Supplementary Figure 3.

649

650