

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

2 Carsten Fortmann-Grote, Julia [von Irmer](#), 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*

16

Deleted: a

Deleted: Balk

Formatted: German

Formatted: German

19 **Abstract**

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

Deleted: provide

Deleted:Page Break.....

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

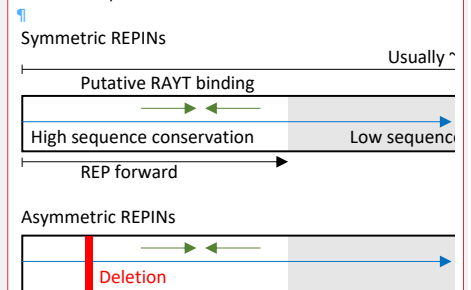
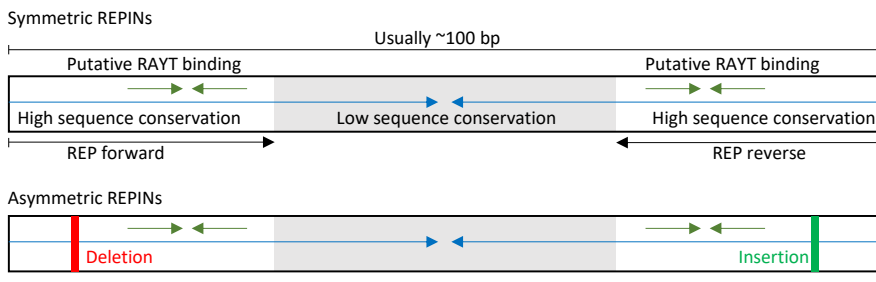


Figure 1. The structure of symmetric and asymmetric REPINs. ...

56 **Introduction**

57 Repetitive sequences in bacteria are rare compared to most eukaryotic genomes. In eukaryotic
58 genomes, repetitive sequences are the result of the activities of persistent parasitic transposable
59 elements. In bacteria, in contrast, parasitic transposable elements cannot persist for long periods
60 of time (Park *et al.* 2021; van Dijk *et al.* 2022). To persist in the gene pool, transposable elements
61 have to constantly infect novel hosts (Sawyer *et al.* 1987; Lawrence *et al.* 1992; Bichsel *et al.*
62 2010; Rankin *et al.* 2010; Wu *et al.* 2015; Park *et al.* 2021). Yet, there is at least one exception: a
63 class of transposable elements called REPINs.

64



65

Figure 1. The structure of symmetric and asymmetric REPINs. A typical REPIN consists of two highly conserved regions at the 5' and 3' end (white), separated by a spacer region of lower sequence conservation (grey). The entire REPIN is an imperfect 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 *Escherichia 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 for example *Pseudomonas fluorescens* SBW25 are symmetric while REPINs in *E. coli* are asymmetric. Asymmetric REPINs make analyses with RAREFAN more challenging. Figure adapted from (Bertels, Rainey 2023), which is licensed under CC BY 4.0.

66 REPINs are short (~100 bp) nested palindromic sequences (**Figure 1**) that consist of two inverted
67 REP (repetitive extragenic palindromic (Higgins *et al.* 1982)) sequences that can be present
68 hundreds of times per genome (Bertels, Rainey 2011a). Most REPINs are symmetric where the 5'
69 REP sequences are identical to the 3' REP sequences, with the occasional substitution (Bertels,

Deleted: of the REPIN

Deleted: a

Deleted: , which makes

Deleted: (Bertels, Rainey 2022).

Deleted: ¶

Deleted: is

72 Rainey 2011a; b). However, there are also asymmetric REPINs where the 5' REP sequence differs
73 from the 3' REP sequence by a point deletion or insertion (Bertels, Rainey 2011a, 2023), which
74 makes the analysis and detection of REPINs significantly more difficult (e.g., *Escherichia coli*
75 REPINs). Isolated REP sequences, sometimes referred to as REP singlets, can also be found in the
76 genome. These sequences are decaying remnants of REPINs that are no longer mobile (Bertels,
77 Rainey 2011a). REPINs are non-autonomous mobile genetic elements, which means they require
78 a RAYT (REP Associated tYrosine Transposase) transposase gene (also referred to as $tnpA_{REP}$) to
79 replicate inside the genome (Nunvar *et al.* 2010; Bertels, Rainey 2011a; Ton-Hoang *et al.* 2012).

Deleted: However, there are also asymmetric REPINs where the 5' REP sequence differs from the 3' REP sequence by a point deletion or insertion (Bertels, Rainey 2011a, 2022), which makes the analysis and detection significantly more difficult (e.g., *Escherichia coli* REPINs). Isolated REP sequences, REP singlets can also be found in the genome. These sequences are decaying remnants of REPINs that are not mobile anymore...

81 Within a genome, each REPIN population is usually only associated with a single RAYT gene.
82 Hence, RAYT genes occur only in single copies per genome and do not copy themselves, unlike
83 for example insertion sequences where often multiple identical sequences are present inside the
84 genome. Unlike insertion sequences RAYT genes are almost exclusively inherited vertically,
85 meaning they are host-beneficial transposases that have been coopted by the host (Bertels,
86 Gallie, *et al.* 2017; Bertels, Rainey 2023). The fact that REPINs and their corresponding RAYT genes
87 are confined to a single bacterial lineage makes them unique, in comparison to all other parasitic
88 mobile genetic elements in bacterial genomes (Bertels, Rainey 2023).

Deleted: to

Deleted: also only

Deleted: are

Deleted: (Bertels, Gallie, *et al.* 2017; Bertels, Rainey 2022).

Deleted: very special

Deleted: 2022

Field Code Changed

90 Of a total of five different RAYT families, there are only two RAYT families that are associated
91 with repetitive sequences such as REPIN or REP sequences: Group 2 and Group 3 RAYTs (Bertels,
92 Gallie, *et al.* 2017). Group 2 RAYTs are present in most Enterobacteria and usually occur only once
93 per genome associated with a single REPIN population. In contrast, Group 3 RAYTs are found in
94 most *Pseudomonas* species and are usually present in multiple divergent copies per genome,
95 each copy associated with a specific REPIN population (Bertels, Gallie, *et al.* 2017).

Deleted: REPINs

97 REPINs and their corresponding RAYT genes occur exclusively in bacterial genomes and are
98 absent in eukaryotic or archaeal genomes (Bertels, Gallie, *et al.* 2017; Bertels, Rainey 2023).
99 Within bacterial genomes REPINs and RAYTs have been evolving in single bacterial lineages for

Deleted: (Bertels, Gallie, *et al.* 2017; Bertels, Rainey 2022).

116 millions of years (Bertels, Gallie, *et al.* 2017). The long term persistence of REPINs in single
117 bacterial lineages can also be observed when analyzing REPIN populations (Bertels, Gokhale, *et*
118 *al.* 2017; Bertels, Rainey 2023).

Deleted: maybe even for a billion

Field Code Changed

Deleted: 2022

120 Parasitic insertion sequences usually occur in identical copies in bacterial genomes, reflecting the
121 fact that insertion sequences persist only briefly before they are eradicated from the genome or
122 kill their host (Park *et al.* 2021). REPINs in contrast are only conserved at the ends of the sequence
123 (presumably due to selection for function), the rest of the sequence is highly variable and only
124 the hairpin structure is conserved (Bertels, Rainey 2011a). The sequence variability of REPINs
125 within the same genome reflects their long-term persistence in single bacterial lineages (Bertels,
126 Rainey 2023). REPINs cannot simply reinfect another bacterial lineage since they rely for mobility
127 on their corresponding RAYT, which itself is immobile.

Deleted: in the genome

Deleted: (Park *et al.* 2021)

Deleted: The sequence variability of REPINs within the same genome reflects their long-term persistence in single bacterial lineages (Bertels, Rainey 2022).

129 RAYTs and REPINs are distinct from typical parasitic insertion sequences, yet we know very little
130 about their evolution or biology. Currently, it is completely unclear what kind of beneficial
131 function maintains REPINs and RAYTs as well as their association with each other. The reason for
132 our lack of knowledge is not because REPINs and RAYTs are rare. They are ubiquitously found in
133 many important and well-studied model bacteria such as Enterobacteria, Pseudomonads,
134 Neisseriads, and Xanthomonads. Microbial molecular biologists presumably encounter REPINs
135 quite frequently. However, connecting the presence or absence of REPINs/RAYTs with
136 phenotypes is difficult if we do not know when it is a REPIN that is present close to a gene of
137 interest or a different type or repeat sequence. Even if the scientist knows about the presence of
138 a REPIN, it is also important to know whether a corresponding RAYT is present (Bertels, Rainey
139 2023).

Deleted: for millions of years in the genome.

Deleted: ,

Deleted: probably

Deleted: , since the function of REPINs largely depends on the function of the presence of a corresponding RAYT gene

Field Code Changed

Deleted: 2022

141 The identification of REPIN populations and their corresponding RAYTs can be rather
142 cumbersome if done from scratch. This is particularly true if the microbial molecular biologist is
143 not aware of all the details of REPIN and RAYT biology. Identifying REPINs starts with an analysis

Deleted: Yet, the

Deleted: ins and outs

159 of short repetitive sequences in a genome. If there are excessively abundant short sequences
160 present in the genome, the distribution of these sequences is then analyzed, if these sequences
161 are exclusively identical tandem repeats without sequence variation, and present in only one or
162 two loci in the genome, then these sequences are probably part of a CRISPR array and not REPINs.
163 If the sequences are distributed across the genome, highly diverse and often present as inverted
164 repeats then it is likely that the repeats are indeed REPINs.
165

Deleted: the

Deleted: next

Deleted: they

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

Formatted: English (US)

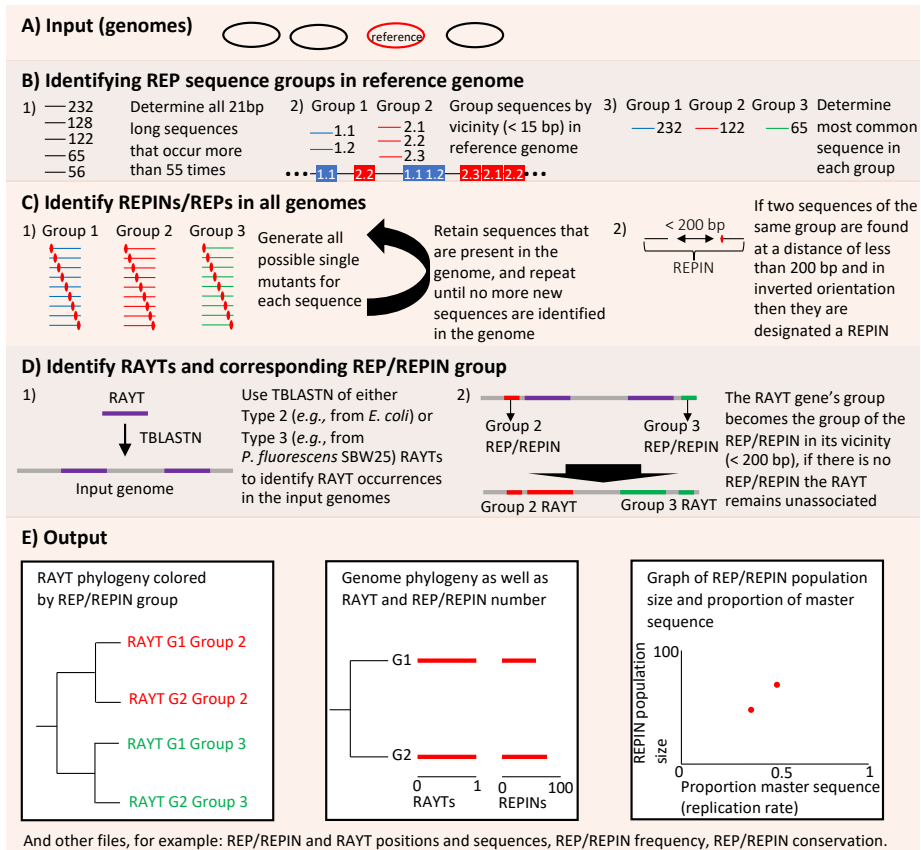


Figure 2. RAREFAN workflow. (a) By default RAREFAN requires the user to supply input sequences containing RAYTs and REPINS. These are fully sequenced and complete genomes. (b) RAREFAN then identifies seed sequence groups (potential REP sequences) in the reference genome by first isolating all 21 bp (adjustable 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 21 bp. Hence, we grouped all sequences together that occur within 15 bp (adjustable parameter) of each other anywhere in the genome. 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

Deleted: 20bp

Deleted: 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. ...

identify RAYT homologs (e-Value < 1e-30, adjustable parameter) of either *E. coli* (Group 2 RAYT) or from *P. 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) The first graph contains a RAYT phylogeny computed from a nucleotide alignment of all identified RAYT genes. 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. 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. RAREFAN also generates various files containing, for example, REP, REPIN, or RAYT sequences and their positions in the query genomes.

Deleted: *Pseudomonas*

Deleted: The alignment is calculated with MUSCLE (Edgar 2004) and a phylogeny with PHYML (Guindon *et al.* 2010).

Deleted: If no genome phylogeny is supplied RAREFAN calculates a whole genome phylogeny of the submitted genomes using andi (Haubold *et al.* 2015).

Deleted: 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).

Formatted: Font: +Body (Calibri), Bold

Deleted: ¶

Moved (insertion) [1]

176 Methods

177 Identification of REPs and REPINs

178 The algorithm to determine REP sequence groups has been described in previous papers and is
179 now slightly improved (Bertels, Rainey 2011a, 2023; Bertels, Gokhale, *et al.* 2017). The main
180 difference between the implementations is that RAREFAN automatically links REPs/REPIN
181 populations to RAYT genes, which was not possible previously.

182 The algorithm starts by extracting all N bp (21 bp by default) long seed sequences that occur
183 more than M times (55 by default) from the reference genome (Figure 2B). All sequences
184 occurring within the reference genome at least once within 15 bp of each other are then grouped
185 together into n REP sequence groups (numbered 0-[n-1]). The most common sequence in each
186 group, named REP seed sequence, is used for further analyses in each input genome.

Moved (insertion) [2]

187 In the next step all possible point mutants of the seed sequences are generated and searched for
188 in the genome (Figure 2C). If a sequence is found in the genome, then all possible point mutants
189 are generated for this sequence as well and searched against the genome and so on until no more
190 sequences can be identified. Once all sequences related to the seed sequence are found,
191 RAREFAN determines whether the sequences form REPINs. Two sequences form a REPIN when
192 they are located within 130 bp of each other in inverted orientation.

193 REP and REPIN sequences of the same type form REPIN populations. A REPIN population is
194 defined as the largest coherent sequence cluster. To identify the largest sequence cluster MCL is
195 applied to a network of REP/REPIN sequences where all sequences that differ by exactly one

196 nucleotide are connected using an inflation parameter of 1.2 (Van Dongen 2000). The clustering
197 results are stored in a file ending in `.mcl`. The sequences of the largest REPIN population
198 (excluding REP singlets) are stored in a file ending in `largestCluster.nodes`. The largest
199 REPIN populations are shown in the REPIN population plot and the master sequence correlation
200 plot (Figure 4).

Moved (insertion) [3]

201 Identification of RAYTs

Moved (insertion) [4]

202 RAYTs are identified using TBLASTN (Camacho *et al.* 2009) with either a protein sequence
203 provided by the user, a Group 2 RAYT from *E. coli* (yafM, Uniprot accession Q47152) or a Group
204 3 RAYT from *P. fluorescens* SBW25 (yafM, Uniprot accession C3JZ76). The presence of RAYTs in
205 the vicinity (default maximum of 200 bp) of a particular REPIN can be used to establish the
206 association between the RAYT gene and a REPIN group (Figure 2D). If a REP sequence or a REPIN
207 of a particular group occur within 200 bp (by default) of a RAYT gene then the RAYT gene is linked
208 to the REP/REPIN group. REP/REPIN associations are stored in the file
209 `repin_rayt_association.txt`. REPIN or REP sequences are almost always present in the
210 extragenic spaces of the RAYT and this linkage is consistent across the RAYT phylogeny (as shown
211 in Figure 5). However, what causes the REP/REPIN group to be linked with the RAYT gene is
212 unclear.

Moved (insertion) [5]

213 Implementation

214 RAREFAN is a modular webservice. It consists of a web frontend written in the `Python`
215 programming language (Van Rossum, Drake Jr 1995) using the `Flask` framework (Grinberg 2018),
216 a `Java` (Arnold *et al.* 2005) backend for genomic sequence analysis and an R (R Core Team 2016)
217 `Shiny` app (RStudio, Inc 2013) for data visualization. The software is developed and tested on the
218 Debian GNU/Linux operating system (Kleinmann *et al.* 2021). All components are released under
219 the MIT opensource license (Initiative 2021) and can be obtained from our public GitHub
220 repository at <https://github.com/mpievolbio-scicomp/rarefan>.

Deleted: python

Deleted: flask

Deleted: java

Deleted: shiny

221 The public RAREFAN instance at <http://rarefan.evolbio.mpg.de> runs on a virtual cloud server with
222 `four` single-threaded CPUs and 16GB of shared memory provided and maintained by the

Deleted: 4

228 Gesellschaft für Wissenschaftliche Datenverarbeitung Göttingen (GWDG) and running the Debian
229 GNU/Linux Operating System (Kleinmann *et al.* 2021).

230 The [Java](#) backend drives the sequence analysis. It makes system calls to TBLASTN (Altschul *et al.*
231 1990) to identify RAYT homologs and to MCL (Van Dongen 2000) for clustering REPIN sequences
232 in order to determine REPIN populations.

233 Jobs submitted through the web server are queued and executed as soon as the required
234 resources become available. Users are informed about the status of their jobs. After job
235 completion, the user can trigger the R [Shiny](#) app to visualize the results.

236 The [Java](#) backend can also be run locally *via* the command line interface (available for download
237 at <https://github.com/mpievbio-scicomp/rarefan/releases>).

238 *Usage of the webservice*

239 The front page of our webservice allows users to upload their bacterial genomes in FASTA (.fas)
240 format (**Figure 2A**). Optionally, users may also provide RAYT protein FASTA sequences (.faa) or
241 phylogenies in NEWICK (.nwk) format. After successful completion of the upload process, the
242 user fills out a web form to specify the parameters of the algorithm:

- 243 • Reference sequence: Which of the uploaded genome sequences will be designated as
244 reference genome (see below for explanations). Defaults to the first uploaded filename
245 in alphabetical order.
- 246 • Query RAYT: The RAYT gene that is used to identify homologous RAYTs in the query
247 genomes. If the user does not provide a protein sequence file then the user can choose
248 one of two RAYT sequences (one from Group 2 and one from Group 3 RAYTs (Bertels,
249 Gallie, *et al.* 2017)) as RAYT query.
- 250 • Tree file: A phylogenetic tree of the reference genomes that can be provided by the user,
251 otherwise the tree will be calculated using *andi* (Haubold *et al.* 2015).
- 252 • Minimum seed sequence frequency: Lower limit on seed sequence frequency in the
253 reference genome to be considered as a REP candidate. Default is 55.

Deleted: java

Deleted: shiny

Deleted: java

- 257 • Seed sequence length: The seed sequence length (in base pairs) is used to identify REPIN
258 candidates from the input genomes. Default is 21 bp.
- 259 • Distance group seeds: The maximum distance between a single occurrence of short
260 repetitive sequences to still be sorted into the same sequence group.
- 261 • REPIN-RAYT association distance: The maximum distance at which a REP sequence can be
262 located from a RAYT gene to be linked to that RAYT gene. Default is 200 bp.
- 263 • e-value cut-off: Alignment e-value cut-off for identifying RAYT homologs with TBLASTN.
264 Default is 1e-30.
- 265 • Analyse REPINs: Ticked REPINs will be analysed (two inverted REP sequences found at a
266 distance of less than 130 bp), if not ticked only short repetitive 21 bp long sequence will
267 be analysed.
- 268 • User email (optional): If provided, then the user will be notified by email upon run
269 completion.

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

277 Visualizations

278 For each REPIN-RAYT group summary plots are generated. These include plots showing the RAYT
279 phylogeny (calculated from a nucleotide alignment using MUSCLE (Edgar 2004) and PHYML
280 (Guindon *et al.* 2010) to generate a phylogeny), REPIN population sizes in relation to the genome
281 phylogeny (provided by the user or if not provided calculated by andi (Haubold *et al.* 2015)) as
282 well as the proportion of master sequences (most common REPIN in a REPIN population) in
283 relation to REPIN population size (Figure 2E).

Deleted: Association distance

Moved up [1]: <#>Identification of REPINs ¶

Moved up [2]: <#> The most common sequence in each group, named REP seed sequence, is used for further analyses in each input genome. ¶
In the next step all possible point mutants of the seed sequences are generated and searched for in the genome (Figure 2C).

Moved up [3]: <#> The clustering results are stored in a file ending in .mcl. The sequences of the largest REPIN population (excluding REP singlets) are

Moved up [4]: <#>The largest REPIN populations are shown in the REPIN population plot and the master sequence correlation plot (Figure 4). ¶

Moved up [5]: <#>Identification of RAYTs ¶
RAYTs are identified using TBLASTN (Camacho *et al.* 2009) with either a protein sequence provided by the user, a Group 2 RAYT from *E. coli* (yafM, Uniprot accession Q47152) or a Group 3 RAYT from *P. fluorescens* SBW25 (yafM, Uniprot accession C3JZZ6). The presence of RAYTs in the vicinity

Deleted: <#>The algorithm to determine REP sequence groups has been described in previous papers and is slightly improved (Bertels, Rainey 2011a, 2022; Bertels, Gokhale, *et al.* 2017). In our current implementation REPs/REPIN populations are now automatically linked to RAYT genes. ¶
First, all N bp (21 bp by default) long seed sequences that occur more than M times (55 by default) are extracted from the reference genome. N and M are the seed sequence length and minimum seed sequence frequency, respectively (Figure 2B). All sequences occurring within the reference genome at least once within 15 bp of each other are then grouped together into n REP sequence groups (numbered 0-(n-1)).

Deleted: <#>If a sequence is found in the genome, then all possible point mutations are generated for this sequence as well and so on until no more sequences can be identified. If two sequences are found within 130 bp of each other in inverted orientation, then these are designated REPINs. ¶
Among all identified REP and REPIN sequences REPIN populations can be isolated. REPIN populations are determined by applying MCL using an inflation parameter of 1.2 (Van Dongen 2000) to a network of REP/REPIN sequences where all sequences that differ by exactly one nucleotide are connected.

Deleted: <#>isolated in a file ending in largestCluster.nodes.

Deleted: <#> ¶

Deleted: <#>The presence of RAYTs in the vicinity (default 200 bp) of a particular REPIN can be used to establish the association between the RAYT gene and a REPIN group (Figure 2D). All positions of all REPINs and REP sequences of a REPIN group are checked whether they occur within 200 bp (by default) of a RAYT gene. If so then the RAYT gene is ... [1]

345 *Other outputs*

346 Identified REPINs, REP singlets as well as RAYTs are written to FASTA formatted sequence files
347 and to tab formatted annotation files that can be read with the Artemis genome browser
348 (Rutherford *et al.* 2000). The REPIN-RAYT associations as well as the number of RAYT copies per
349 genome are written to tabular data files. A detailed description of all output files is provided in
350 the manual (<http://rarefan.evolbio.mpg.de/manual>) and in the file “readme.md” in the output
351 directory.

352 **Sequence analysis and annotation**

353 For verification of RAREFAN results, REPIN-RAYT-systems were analysed in their corresponding
354 genomes using Geneious Prime version 2022.2.2 (Kearse *et al.* 2012). Nucleotide sequences and
355 positions of REP singlets, REPINs, and RAYTs were extracted from output files generated by
356 RAREFAN and mapped in the relevant *S. maltophilia* genome. Complete RAREFAN data used for
357 analysis can be accessed by using the run IDs listed in **Table 1**.

358
359 **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

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

361

362 **Results**

363 RAREFAN can identify REPINs and their corresponding RAYTs in a set of fully sequenced bacterial
364 genomes. The RAREFAN algorithm has been used in previous analyses to identify and
365 characterize REPINs and RAYTs in Pseudomonads (Bertels, Rainey 2011a, 2023), Neisseriads
366 (Bertels, Rainey 2023), and Enterobacteria (Bertels, Gallie, *et al.* 2017; Park *et al.* 2021). To
367 demonstrate RAREFAN’s capabilities, we are presenting an analysis of 49 strains belonging to the
368 opportunistic pathogen *S. maltophilia*.

Deleted: ¶

Formatted: Left

Formatted Table

Formatted: Left

Formatted: Left

Formatted: Left

Formatted: Left

Formatted: Left

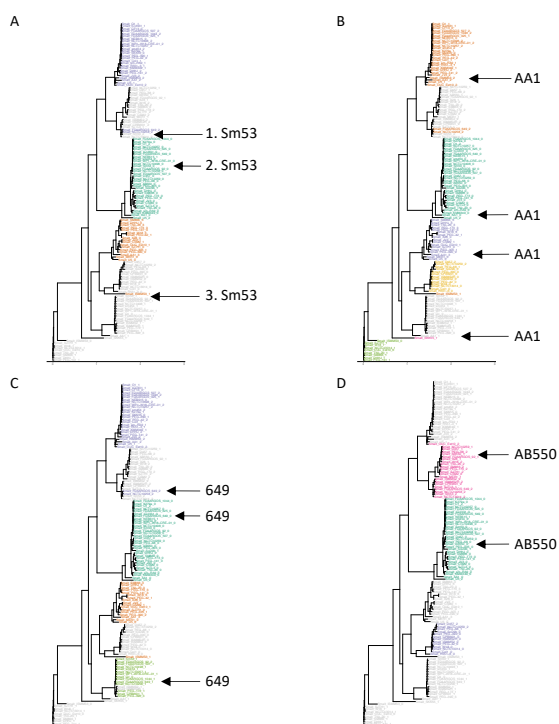
Deleted: (Bertels, Rainey 2011a, 2022)

Deleted: (Bertels, Rainey 2022), and Enterobacteria (Bertels, Gallie, *et al.* 2017; Park *et al.* 2021)

373 *S. maltophilia* strains contain Group 3 RAYTs, which are also commonly found in plant-associated
374 *Pseudomonas* species such as *P. fluorescens* or *P. syringae* (Bertels, Rainey 2011a, 2023). Similar
375 to Group 3 RAYTs in other species, *S. maltophilia* contains multiple REPIN-RAYT systems per
376 genome. Group 2 RAYTs, in contrast, tend to contain only one REPIN-RAYT system per genome
377 (Bertels, Rainey 2023).

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

378



379

Figure 3. Phylogenetic trees built from RAYT genes extracted from *S. maltophilia* genomes.

RAYT genes are coloured by RAREFAN according to their association with REPIN populations in the reference genome. If a REPIN population of a query genome is not present in the reference genome, then the REPIN population cannot be identified in the query genome and the corresponding 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 with 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.

Deleted: to

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

Deleted: ¶

388 REPIN-RAYT systems in *S. maltophilia* are surprisingly diverse compared to other species. For
389 example, *Pseudomonas chlororaphis* contains three separate REPIN populations that are present
390 in all *P. chlororaphis* strains, each associated with its cognate RAYT gene (Bertels, Rainey 2023).
391 *S. maltophilia*, in contrast, contains only one REPIN-RAYT system that is present across almost
392 the entire species (green clade in **Figure 3**), and at least eight REPIN-RAYT systems that are
393 present in subsets of strains (nine clades in **Figure 5**).

Deleted: (Bertels, Rainey 2022).

394 The patchy presence-absence pattern of REPIN-RAYT systems in *S. maltophilia*, makes the dataset
395 quite challenging to analyse. If a REPIN population is not present in the reference strain then
396 RAREFAN will not be able to detect it in any other strain. Yet, it is possible to detect RAYT genes
397 in all strains of a species independent of the reference strain selection. RAYT genes that are not
398 associated with a REPIN population are displayed in grey (**Figure 3A**). While these RAYT genes are
399 not associated with REPIN populations detected in the reference strain, they might still be
400 associated with a yet unidentified REPIN type present in the genome the unassociated RAYT gene
401 is located in.

Deleted: to

Deleted: to

402 In order to identify all REPIN populations across a species, multiple RAREFAN runs with different
403 reference strains **should be performed**. The RAREFAN web interface supports re-launching a
404 given job with modified parameters. To identify as many different REPIN-RAYT systems as
405 possible in each subsequent run the reference should be set to a genome that contains RAYTs
406 that were not associated with a REPIN population previously (*i.e.*, genomes containing grey RAYTs
407 in **Figure 3**). However, this strategy may also fail when the REPIN population size falls below the
408 RAREFAN seed sequence frequency threshold. In that case reducing the frequency threshold will
409 be more productive.

Deleted: we suggest to perform

410 For example, *S. maltophilia* Sm53 contains three RAYTs only one of which is associated with a
411 REPIN population (RAYT genes indicated in **Figure 3A**). However, the remaining two RAYTs are
412 indeed associated with a REPIN population, but these REPIN populations are too small to be

418 detected in *S. maltophilia* Sm53 (the seed sequence frequency threshold is set to 55 by default).
 419 In other *S. maltophilia* strains the REPIN populations are large enough to exceed the threshold.
 420 For example, if *S. maltophilia* AB550 is set as reference, RAYT number 1 from Sm53 (**Figure 3A**) is
 421 associated with the pink REPIN population (**Figure 3D**). If *S. maltophilia* 649 is set as reference
 422 RAYT number 3 from Sm53 (**Figure 3A**) is associated with the light green REPIN population (**Figure**
 423 **3C**). RAYTs from the bottom clade are only associated with REPIN populations when *S. maltophilia*
 424 AA1 is chosen as reference (**Figure 3B**). While lower thresholds can guarantee that all REPINs will
 425 be identified in the genome, the number of sequence groups that are not REPINs quickly
 426 explodes. This is especially true for genomes that contain large numbers of mobile genetic
 427 elements or CRISPRs (Bertels, Rainey 2023).

Deleted: Especially

Deleted: 2022

Field Code Changed

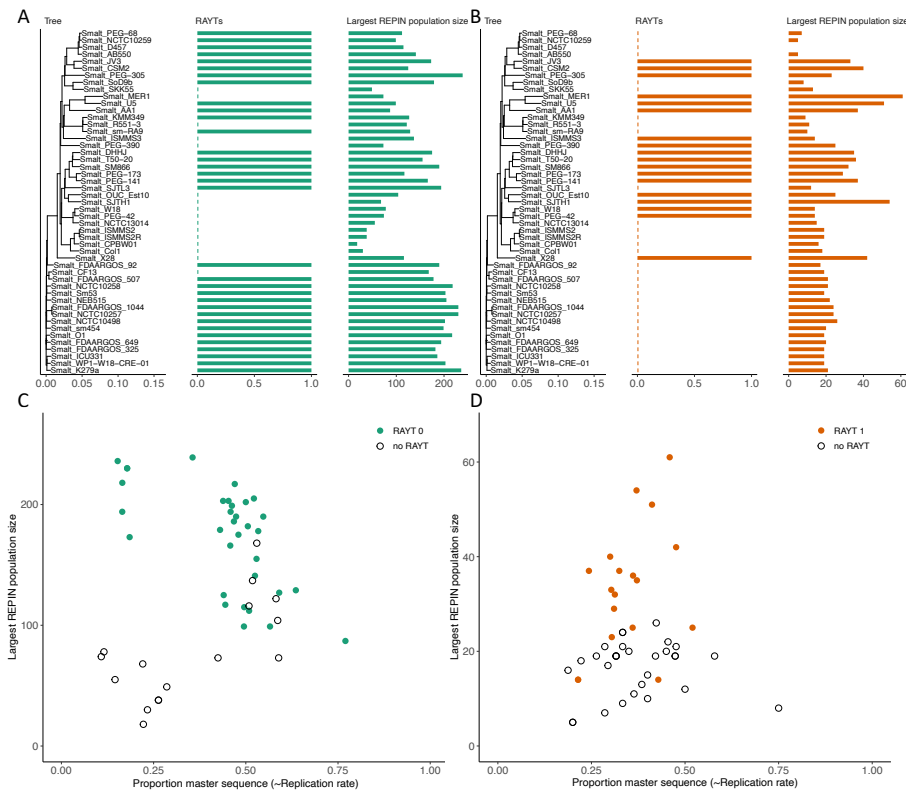


Figure 4. REPIN population sizes and conservation. The plots show two REPIN populations and their associated RAYTs that were identified in *S. maltophilia* using *S. maltophilia* Sm53 as reference. **(A)** The phylogenetic tree on the left side is a whole genome phylogeny generated by *andi* (Haubold *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 orange 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 orange. **(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. At 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.

Deleted: In an

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

432 The RAREFAN webserver visualizes REPIN population size and RAYT numbers in barplots. Barplots
433 are ordered by the phylogenetic relationship of the submitted bacterial strains using ggtree (Yu
434 *et al.* 2018). RAREFAN detects three populations when *S. maltophilia* Sm53 is selected as
435 reference strain (**Figure 3A**). The largest REPIN population (calculated by MCL from all REPINs of
436 that type) has a corresponding RAYT gene in almost all strains (first barplot in **Figure 4A**) and
437 most REPIN populations contain more than 100 REPINs (second barplot in **Figure 4A**). The second
438 largest REPIN population in Sm53 (orange population in **Figure 4B**) is significantly smaller and
439 contains no more than 61 REPINs in any strain and most strains do not contain a corresponding
440 RAYT for this population.

441 RAREFAN also provides information on REPIN replication rate (**Figure 4C and D**). REPIN replication
442 rate can be estimated by dividing the number of the most common REPIN sequence (master
443 sequence) by the REPIN population size if the population is in mutation selection balance (Bertels,
444 Gokhale, *et al.* 2017). If a REPIN replicates very fast most of the population will consist of identical
445 sequences because mutations do not have enough time to accumulate between replication
446 events. Hence, the proportion of master sequences will be high in populations that have a high
447 replication rate. Transposable elements such as insertion sequences consist almost exclusively of

448 identical master sequences because the time between replication events is not sufficient to
449 accumulate mutations and because quick extinction of the element usually prevents the
450 accumulation of mutations after replication (Park et al. 2021; Bertels, Rainey 2023). Sequence
451 diversity of REPIN populations in contrast is much higher suggesting that REPINs replicate slowly
452 and persist for long periods of time.

Deleted: (Park et al. 2021; Bertels, Rainey 2022). REPIN populations in contrast replicate slowly and persist for long periods of time, which means that a high proportion of master sequences suggests a high REPIN replication rate.

453 In *S. maltophilia* the proportion of master sequences in the population does not seem to correlate
454 well with REPIN population size, both in the green and the orange population (Figures 4C and D).
455 Similar observations have been made in *P. chlororaphis* (Bertels, Rainey 2023), and may suggest
456 that an increase in population size is not caused by an increase in replication rate. Population size
457 is likely to be more strongly affected by other factors such as the loss of the corresponding RAYT
458 gene, which leads to the decay of the REPIN population. One could even speculate that high
459 REPIN replication rates are more likely to lead to the eventual demise of the population due to
460 the negative fitness effect of high replication rates on the host (Park et al. 2021; Bertels, Rainey
461 2023).

Deleted: (Bertels, Rainey 2022)

462 The presence of RAYTs and the size of the corresponding REPIN population do correlate
463 surprisingly well (Figure 4A and B, p-value = 0.008 of a linear model of independent contrasts
464 (Felsenstein 1985) of green RAYT and REPIN number, p-value = 0.003 for orange REPIN
465 populations). Green RAYTs are absent from an entire *S. maltophilia* clade (middle of Figure 4A).
466 This clade has also lost a significant amount of green REPINs, and the proportion of the master
467 sequences in these populations is low (Figure 4C). Similarly, genomes without orange RAYTs have
468 smaller REPIN populations in the orange population than genomes with the corresponding RAYT
469 (Figure 4D). A similar observation has been made previously in *E. coli*, *P. chlororaphis*, *N.*
470 *meningitidis* and *N. gonorrhoeae* where the loss of the RAYT gene is followed by a decay of the
471 associated REPIN population (Park et al. 2021; Bertels, Rainey 2023).

Formatted: Font: Calibri

Formatted: Font: Calibri

Field Code Changed

Deleted: 2022

Formatted: Font: Calibri

Deleted: Value

Deleted: the

Deleted: Value

Deleted: 2022

Field Code Changed

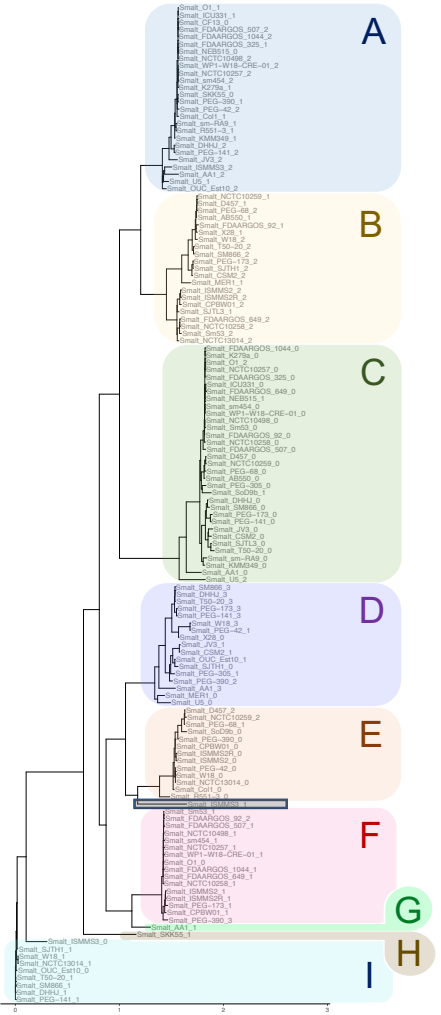


Figure 5. Phylogeny of RAYT genes and their associated REPINs. The tree shows RAYT genes from 49 *S. maltophilia* strains. Colours of clades A-I are assigned according to their association with a REPIN found within 200 bp of the RAYT gene (see Table 2). Except for a single RAYT gene ISMMS3_1 (grey box), which could not be linked to a REPIN population.

Deleted: 130

Table 2. REPIN palindromes associated with each RAYT clade.

RAYT population	REPIN palindromes
-----------------	-------------------

A	CCGACCAACGGTCGG
B	CCAACCAAGGTTGGC
C	CCGGCCAGCGGCCGG
D	TCCACGCATGGCGTGGA
E	CCGAGCCCATGCTCGG
F	TCGACTAACAGTCGA
G	TCGACCAACGGTCGA
H	GCCGGGCATGGCCCGGC
I	AGTCGAGCTTGCTCGACT

484 Each RAYT clade from **Figure 5** is associated with a unique imperfect palindrome that is present
 485 at the 5' and/or 3' end of the RAYT gene.

486

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

488 Unfortunately, RAREFAN is not always able to link the correct REPIN population with the correct
 489 RAYT gene. In some RAREFAN runs, associations between RAYTs and REPINs are not
 490 monophyletic, as for example the red RAYT clade in **Figure 3A**. However, the same clade of RAYTs
 491 is uniformly coloured in yellow in **Figure 3D**, suggesting that the entire RAYT clade is associated
 492 with the same REPIN group.

493 An analysis of all REPIN groups that were identified by RAREFAN across four different RAREFAN
 494 runs (**Table 1**, one additional analysis was performed with ISMMS3) showed that there are a total
 495 of nine different REPIN groups, each defined by an individual central palindrome (**Table 2**). Each
 496 REPIN group is associated with a monophyletic RAYT group (**Figure 5**). Only a single RAYT is not
 497 associated with a REPIN population (ISMMS3_1).

498 RAREFAN could not link a REPIN to the RAYT gene ISMMS3_1 (**Figure 5**, grey box). While there is
 499 a sequence that resembles the A palindrome as well as variants of the C palindrome flanking both
 500 sides of the RAYT gene (**Supplementary Figure 2**), none of the sequences formed REPIN
 501 populations large enough to be identified by RAREFAN. Presumably the RAYT ISMMS3_1, which

Deleted: RAYTs

503 is only present in a single *S. maltophilia* strain, is at the early stages of establishing a REPIN
504 population, and the REPIN population has not spread to a considerable size yet.

505 If the maximum REPIN-RAYT distance parameter is too small then RAREFAN will also fail to
506 correctly link REPINs and RAYTs. For example, when the maximum REPIN-RAYT distance
507 parameter is set to 130bp there are two cases where RAREFAN fails to link RAYT genes with
508 REPINs (ISSMS2_ and ISSMS2R_1, Supplementary Figure 1 D and E). When the parameter is set
509 to a distance of 200 bp (default RAREFAN setting), RAREFAN correctly links these REPINs to the
510 RAYT gene.

511 In three cases the wrong REPIN population was linked to a RAYT gene. In our dataset this can
512 happen when RAYTs are flanked by seed sequences from two different REPIN populations
513 **(Supplementary Figure 1 A-C)**. A single REP sequence from the “wrong” (non-monophyletic
514 RAYT) clade occurs together with multiple REP or REPIN sequences from the “right”
515 (monophyletic in a different RAREFAN run) clade. REPINs are linked to the “wrong” RAYT when
516 the correct REPIN population is absent in the chosen reference genome. This problem can be
517 alleviated by performing analyses with multiple reference genomes and comparing the results.

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

519 The seed distance parameter determines whether two highly abundant sequences are sorted
520 into the same or different REPIN groups (**Figure 2B**). If two REPINs from two different groups
521 occur next to each other, at a distance of less than the seed distance parameter, then the two
522 seeds are erroneously sorted into the same group. If two different REPIN groups are sorted into
523 the same group then one of the groups will be ignored by RAREFAN, because only the most
524 abundant seed in each group will be used to identify REPINs.

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

Deleted: There

Deleted: more

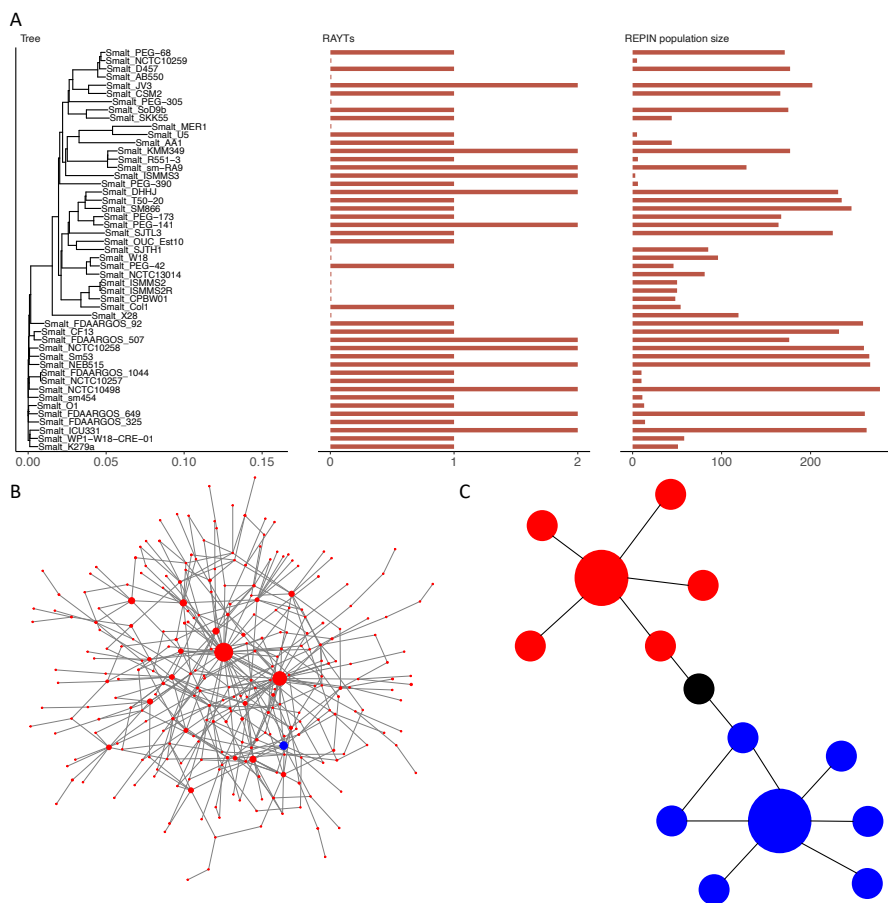
Deleted: failed

Deleted: any

Deleted: Detailed sequence analyses showed that the respective REPINs are located at a distance of more than 130 bp (an adjustable parameter in RAREFAN). These REPINs are ignored by RAREFAN by default. However, this parameter can be adjusted manually and when

Deleted: ,

540 A small seed distance parameter will separate seed sequences belonging to the same REPIN
541 group into different groups. Hence, RAREFAN will analyse the same REPIN group multiple times.
542 While this will lead to increased RAREFAN runtimes, these errors, are easy to spot, because (1)
543 the same RAYT gene will be associated to multiple REPIN groups, (2) the central palindrome
544 between the group is identical and (3) the master sequence between the groups will be very
545 similar.
546



547

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 clusters. 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.

- Deleted: Group
- Deleted: Group
- Deleted: Group
- Deleted: Group
- Deleted: Group
- Deleted: Group
- Deleted: cluster
- Deleted: Group
- Deleted: Group

548

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

550 Incorrect merging of REPIN groups can occur when two REPIN groups are closely related. We
 551 identified merged REPIN groups in *S. maltophilia* because RAREFAN linked some REPIN groups
 552 with two RAYT genes in the same genome (**Figure 6A**). While REPIN groups linked to two RAYTs
 553 have been observed before in *Neisseria meningitidis* (Bertels, Rainey 2023), it is particularly
 554 unusual in *S. maltophilia* due to some key differences between REPIN-RAYT in the two bacterial
 555 species. First, *N. meningitidis* contains Group 2 RAYTs and *S. maltophilia* only contains Group 3
 556 RAYTs (Bertels, Gallie, *et al.* 2017), two very divergent RAYT gene families. Second, RAYTs that
 557 are associated with the same REPIN group in *N. meningitidis* are almost identical, since they are
 558 copied by an insertion sequence *in trans* (Bertels, Rainey 2023), something that is not the case
 559 for *S. maltophilia*, where the two RAYTs are very distinct and quite distantly related from each
 560 other (green and red clade in **Figure 3A**, or clade A and C in **Figure 5**).

- Deleted: While REPIN groups linked to two RAYTs has been observed before in *Neisseria meningitidis* (Bertels, Rainey 2022)...
- Deleted: RAYTs in
- Deleted: belong to
- Deleted: RAYTs in
- Deleted: belong to
- Deleted: groups.
- Deleted: to
- Deleted: (Bertels, Rainey 2022), something that is not the case for *S. maltophilia*, where the two RAYTs are very distinct

561 A closer inspection of all sequences identified in REPIN group 2 shows that it also contains
 562 sequences belonging to REPIN group 0 (palindromes linked to clade A and C in **Table 2**). The
 563 relationship between the sequences shows that there is a chain of sequences that all differ by at
 564 most a single nucleotide between the most abundant sequence in group 2 to the most abundant
 565 sequence in group 0 (**Figure 6B and C**). Hence, the reason group 0 and group 2 are merged is that

577 they are too closely related to each other and hybrids of the two REPIN groups exist. Because
578 sequence groups are built by identifying all related sequences in the genome recursively, closely
579 related groups (the REPIN group 0 seed only differs in four nucleotides from the REPIN group 2
580 seed sequence) can be merged into a single REPIN group. REPIN population size and RAYT number
581 are the sum of REPIN group 0 and 2. There are various possibilities to resolve this issue: (1)
582 subtract sequences from group 0 (which does not contain group 2) from REPIN group 2; (2) use
583 a different sequence seed from the group 2 seed collection in the seed sequence file
584 (groupSeedSequences/Group_SmaI_t_Sm53_2.out); (3) sometimes it may be possible to
585 rerun RAREFAN with a different reference strain where the issue does not occur; or (4) increase
586 the length of the seed sequence.

587 *Performance*

588 We measured the elapsed time for a complete RAREFAN run for three different species and for
589 5, 10, 20, and 40 genomes with randomly selected reference strains and the two query RAYTs
590 (yafM_Ecoli and yafM_SBW25). For a given number N of submitted genomes of average
591 sequence length L (in megabases), a RAREFAN run completes in approximately $T = (8-10 \text{ seconds})$
592 $* N * L$ on our moderate server hardware (4 CPU cores, 16 GB shared RAM) (**Supplementary**
593 **Figure 3 and 4**).

594 **Discussion**

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

599 One limitation of RAREFAN is that REPINs can only be identified in genomes when they are
600 symmetric (**Figure 1**). Symmetric REPINs have seed sequences that can morph into each other by
601 a series of single substitutions (intermediate sequences need to be present in the genome). A
602 REPIN consists of a 5' and a 3' REP sequence. If one of these REP sequences contains an insertion
603 or deletion, which the other REP sequence does not contain then RAREFAN will not recognize the
604 second repeat of the seed sequence. In this case, RAREFAN will not be able to identify REPINs but

605 can still be used to analyze REP singlet populations. To date, the only asymmetric REPIN
606 populations known to us are found in *E. coli*. However, it is likely that asymmetric REPINs also
607 exist in other microbial species.

Deleted: known

608 RAREFAN sometimes cannot correctly divide REPINs into REPIN groups. Either because REPINs
609 from different groups occur in close proximity in the genome, an issue that can easily be solved
610 by adjusting a RAREFAN parameter, or because two REPIN groups are very closely related (**Figure**
611 **6**). Unfortunately, RAREFAN is not able to automatically detect and resolve the assignment of
612 closely related REPINs into groups yet. Hence it is advisable to manually check associations
613 between REPIN groups and RAYT genes by analyzing the composition of REPIN groups.

614 In the future we aim to make RAREFAN even more versatile and easier to use by, for example,
615 automatically integrating data from public databases such as GenBank, and creating a RAREFAN
616 Galaxy workflow (Afgan *et al.* 2018).

Deleted: Genbank

Deleted: integrating

Deleted: into workflows such as

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

620 **Acknowledgements**

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

623 **References**

624 Afgan E, Baker D, Batut B, van den Beek M, Bouvier D, Cech M, Chilton J, Clements D, Coraor N,
625 Gruning BA, Guerler A, Hillman-Jackson J, Hiltmann S, Jalili V, Rasche H, Soranzo N,
626 Goecks J, Taylor J, Nekrutenko A, Blankenberg D (2018) The Galaxy platform for
627 accessible, reproducible and collaborative biomedical analyses: 2018 update. *Nucleic*
628 *Acids Res.*, **46**, W537-W544-W537-W544. <https://doi.org/10.1093/nar/gky379>

633 Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ (1990) Basic local alignment search tool.
634 *Journal of Molecular Biology*, **215**, 403–410. <https://doi.org/10.1006/jmbi.1990.9999>
635 Arnold K, Gosling J, Holmes D (2005) *The Java programming language*. Addison Wesley
636 Professional.
637 Bertels F, Gallie J, Rainey PB (2017) Identification and Characterization of Domesticated Bacterial
638 Transposases. *Genome Biology and Evolution*, **9**, 2110–2121.
639 <https://doi.org/10.1093/gbe/evx146>
640 Bertels F, Gokhale CS, Traulsen A (2017) Discovering Complete Quasispecies in Bacterial
641 Genomes. *Genetics*, **206**, 2149–2157. <https://doi.org/10.1534/genetics.117.201160>
642 Bertels F, Rainey PB (2011a) Within-Genome Evolution of REPINs: a New Family of Miniature
643 Mobile DNA in Bacteria. *PLoS genetics*, **7**, e1002132.
644 <https://doi.org/10.1371/journal.pgen.1002132>
645 Bertels F, Rainey PB (2011b) Curiosities of REPINs and RAYTs. *Mobile Genetic Elements*, **1**, 262–
646 268. <https://doi.org/10.4161/mge.18610>
647 Bertels F, Rainey PB (2023) Ancient Darwinian replicators nested within eubacterial genomes.
648 [BioEssays](https://doi.org/10.1002/bies.202200085), **45**, 2200085. <https://doi.org/10.1002/bies.202200085>
649 Bichsel M, Barbour AD, Wagner A (2010) The early phase of a bacterial insertion sequence
650 infection. *Theoretical Population Biology*. <https://doi.org/10.1016/j.tpb.2010.08.003>
651 Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL (2009) BLAST+:
652 architecture and applications. *BMC Bioinformatics*, **10**, 421–9.
653 <https://doi.org/10.1186/1471-2105-10-421>

Deleted: 2022

Deleted: , 2021.07.10.451892.

Deleted: 1101/2021.07.10.451892

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

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

662 Felsenstein J (1985) Phylogenies and the comparative method. *American Naturalist*, 1–15.

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

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

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

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

674 Initiative TOS (2021) The MIT License.

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

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

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

681 Nunvar J, Huckova T, Licha I (2010) Identification and characterization of repetitive extragenic
682 palindromes (REP)-associated tyrosine transposases: implications for REP evolution and
683 dynamics in bacterial genomes. *BMC Genomics*, **11**, 44. [https://doi.org/10.1186/1471-](https://doi.org/10.1186/1471-2164-11-44)
684 [2164-11-44](https://doi.org/10.1186/1471-2164-11-44)

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

687 R Core Team (2016) R: A Language and Environment for Statistical Computing.

688 Rankin DJ, Bichsel M, Wagner A (2010) Mobile DNA can drive lineage extinction in prokaryotic
689 populations. *Journal of Evolutionary Biology*. [https://doi.org/10.1111/j.1420-](https://doi.org/10.1111/j.1420-9101.2010.02106.x)
690 [9101.2010.02106.x](https://doi.org/10.1111/j.1420-9101.2010.02106.x)

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

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

695 Sawyer SA, Dykhuizen DE, DuBose RF, Green L, Mutangadura-Mhlanga T, Wolczyk DF, Hartl DL
696 (1987) Distribution and Abundance of Insertion Sequences Among Natural Isolates of
697 *Escherichia coli*. *Genetics*, **115**, 51–63. <https://doi.org/10.1093/genetics/115.1.51>

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

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

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

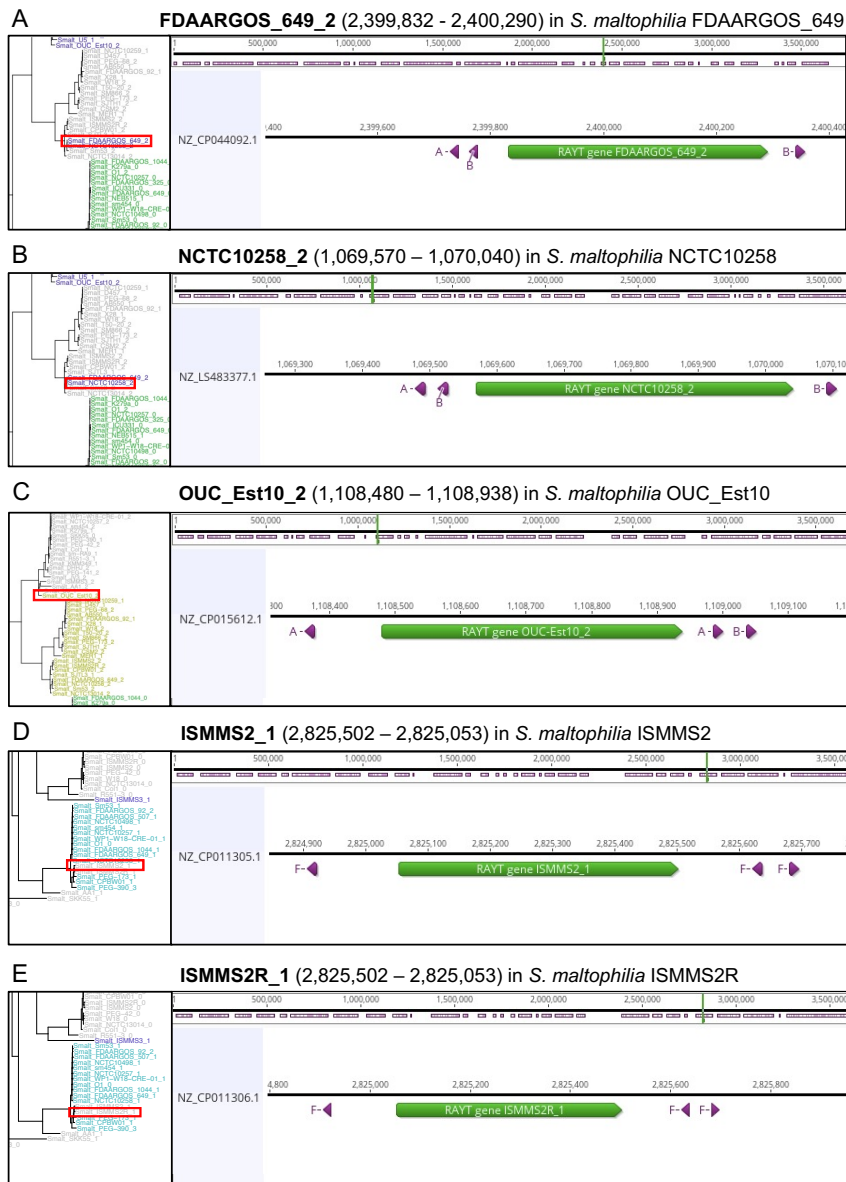
704 Wu Y, Aandahl RZ, Tanaka MM (2015) Dynamics of bacterial insertion sequences: can
705 transposition bursts help the elements persist? *BMC Evolutionary Biology*, **15**, 288–12.
706 <https://doi.org/10.1186/s12862-015-0560-5>

707 Yu G, Lam TT-Y, Zhu H, Guan Y (2018) Two Methods for Mapping and Visualizing Associated Data
708 on Phylogeny Using Ggtree. (FU Battistuzzi, Ed.). *Molecular biology and evolution*, **35**,
709 3041–3043. <https://doi.org/10.1093/molbev/msy194>

710

711

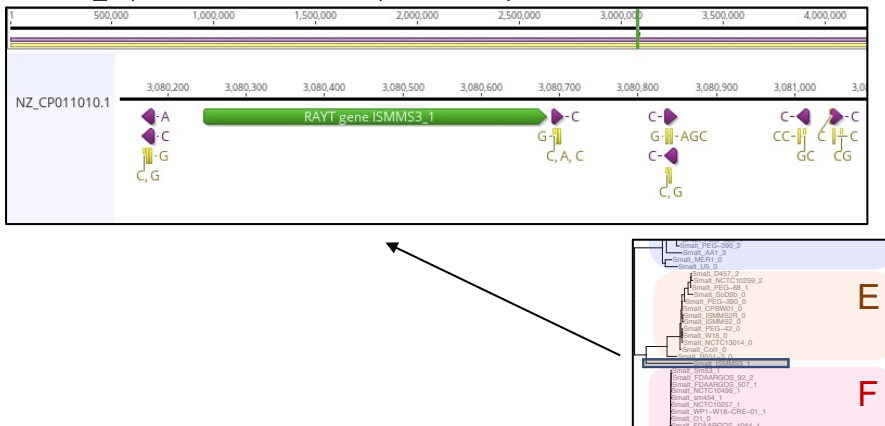
712 **Supplementary Figures**



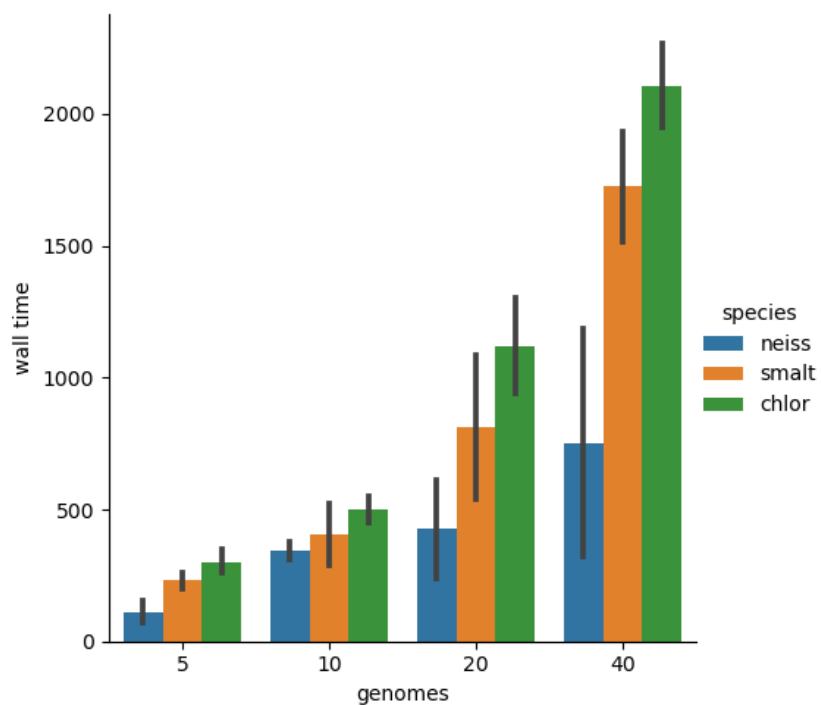
714 **Supplementary Figure 1. Sequence analysis shows REPIN groups are indeed associated with**
 715 **monophyletic RAYTs.** Non-monophyletic or missing associations to REPIN populations identified
 716 by RAREFAN were investigated in the corresponding genomes using Geneious (Kearse *et al.*
 717 2012). Red boxes mark the position of the atypical RAYT that is being analyzed in detail. Mapping
 718 of REPIN palindromes A-I (with zero mismatches) shows FDAARGOS_649_2 (A), NCTC10258_2
 719 (B), and OUC_Est_2 (C) are linked to the wrong REPIN group because REP singlets that are
 720 ordinarily linked to a RAYT sister clade are found in close proximity to the RAYT. These wrong
 721 associations between REPIN and RAYT usually occur when the correct REPIN population is absent
 722 from the reference genome. ISMMS2R_1 (D) and ISMMS2_1 (E) are not linked to REPIN
 723 populations by RAREFAN when the maximum REPIN-RAYT distance parameter is set to 130 bp.
 724 The RAYTs are linked to the correct REPIN populations when the REPIN-RAYT distance parameter
 725 is set to 200 bp (default). Nucleotide sequences and positions were extracted from output files
 726 generated by RAREFAN. Complete genome sequences are available in NCBI Nucleotide Database
 727 using Accessions: (A) NZ_CP044092.1, (B) NZ_LS483377.1, (C) NZ_CP015612.1, (D)
 728 NZ_CP011306.1, (E) NZ_CP011305.1.

- Deleted: were
- Deleted: population
- Deleted: because
- Deleted: corresponding seed sequences were located at a
- Deleted: of more than
- Deleted: from
- Deleted: RAYT gene.

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



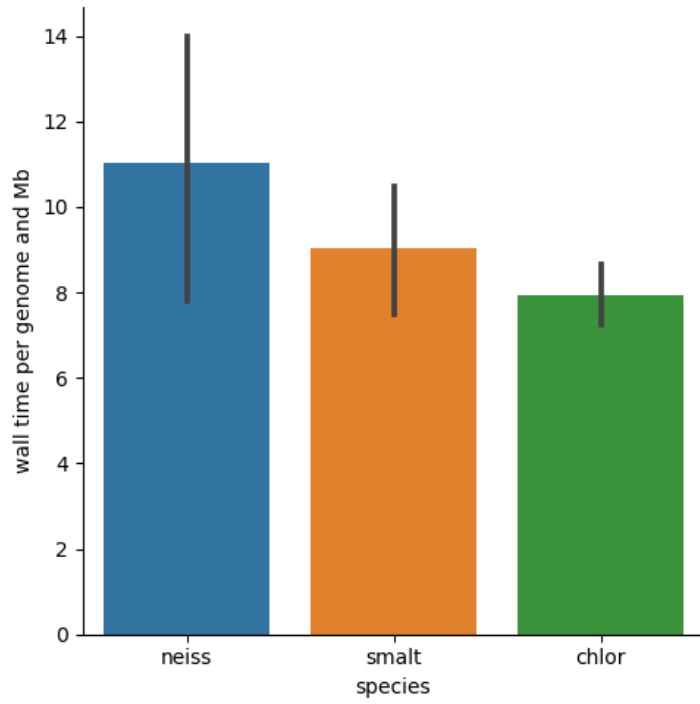
729
 730 **Supplementary Figure 2. RAYT gene ISMMS3_1 cannot be linked to a REPIN population.** The
 731 sequence of the RAYT gene ISMMS3_1 and its flanking sequences were analysed in Geneious
 732 (Kearse *et al.* 2012). The inset shows the location of ISMMS3_1 in the RAYT phylogeny (grey box).
 733 When mapping all of the identified palindromes to the RAYT region and allowing up to four
 734 mismatches (yellow annotations), various mutants of palindrome C were found in close proximity
 735 of the RAYT gene. However, we could not identify a corresponding REPIN population, which may
 736 indicate that the population has not yet expanded in the genome.



744 **Supplementary Figure 3.** Average time (in seconds) it takes RAREFAN to complete for different
 745 genome numbers from three bacterial species (*N. meningitidis*, *N. gonorrhoeae*, *S. maltophilia*,
 746 *Pseudomonas chlororaphis*). Black bars indicate the 95% CI across four runs, where two runs
 747 share the same query RAYT. For each run reference and query strains were randomly selected.
 748 All measurements were performed on 4CPU cores with 16 GB of shared memory.
 749

750

Formatted: Font: Italic



751

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

754

755

|

