

1 **PaliDIS: A tool for fast discovery of novel insertion sequences**

2

3 **Authors**

4 Victoria R. Carr<sup>1,2,\*</sup>, Solon P. Pissis<sup>3,4</sup>, Peter Mullany<sup>5</sup>, Saeed Shoaie<sup>2,6</sup>, David Gomez-  
5 Cabrero<sup>2,7,8</sup>, David L. Moyes<sup>2</sup>

6

7 <sup>1</sup>Parasites and Microbes, Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton,  
8 UK

9

10 <sup>2</sup>Centre for Host-Microbiome Interactions, Faculty of Dentistry, Oral & Craniofacial  
11 Sciences, King's College London, SE1 9RT, UK

12

13 <sup>3</sup>Centrum Wiskunde en Informatica, Amsterdam, NL

14

15 <sup>4</sup>Vrije Universiteit, Amsterdam, NL

16

17 <sup>5</sup>Department of Microbial Diseases, Eastman Dental Institute, University College London,  
18 256 Gray's Inn Road, London, WC1X 8LD, UK

19

20 <sup>6</sup>Science for Life Laboratory, KTH – Royal Institute of Technology, Stockholm, SE-171 21,  
21 Sweden

22

23 <sup>7</sup>Bioscience Program, Bioengineering Program, Biological and Environmental Science and  
24 Engineering Division, King Abdullah University of Science and Technology (KAUST),  
25 Thuwal 23955-6900, Saudi Arabia

26

27 <sup>8</sup>Translational Bioinformatics Unit, Navarrabiomed, Complejo Hospitalario de Navarra  
28 (CHN), Universidad Pública de Navarra (UPNA), IdISNA, Pamplona, Spain

29

30 \*Corresponding author

31

32

33

33 **Abstract**

34 The diversity of microbial insertion sequences, crucial mobile genetic elements in generating  
35 diversity in microbial genomes, needs to be better represented in current microbial databases.

36

36 Identification of these sequences in microbiome communities presents some significant  
37 problems that have led to their underrepresentation. Here, we present a software tool called

37

38 PaliDIS that recognises insertion sequences in metagenomic sequence data rapidly by  
39 identifying inverted terminal repeat regions from mixed microbial community genomes.

39

40 Applying this software to 266 human metagenomes identifies 11,681 unique insertion  
41 sequences. Querying this catalogue against a large database of isolate genomes reveals

41

42 evidence of horizontal gene transfer events of clinically relevant antimicrobial resistance  
43 genes between classes of bacteria. We will continue to apply this tool more widely, building

43

44 the Insertion Sequence Catalogue, a valuable resource for researchers wishing to query their  
45 microbial genomes for insertion sequences.

46

#### 47 **Keywords**

48 Insertion sequences, transposon, metagenome, horizontal gene transfer, mobile genetic  
49 element, antimicrobial resistance, software

50

#### 51 **Abbreviations**

52 ARG – antimicrobial resistance gene

53 bp – base pairs

54 ISC –Insertion Sequence Catalogue

55 IS – insertion sequence/unit transposon

56 ITR – inverted terminal repeat

57 MEM – maximal exact match

58 PaliDIS - Palindromic Detection of Insertion Sequences

59

#### 60 **Data Summary**

61 1. The PaliDIS software is available here: [github.com/blue-moon22/palidis](https://github.com/blue-moon22/palidis)

62 2. The Insertion Sequence Catalogue is available to download here:

63 <https://github.com/blue-moon22/ISC>

64 3. The raw reads from the Human Microbiome Project can be retrieved using the  
65 download links provided in Supplementary Data 1

66 4. The 21 contig files can be retrieved using the download links provided in  
67 Supplementary Data 3

68

#### 69 **Impact Statement**

70 Insertion sequences are a class of transposable element that play an important role in the  
71 dissemination of antimicrobial resistance genes. However, it is challenging to completely  
72 characterise the transmission dynamics of insertion sequences and their precise contribution  
73 to the spread of antimicrobial resistance. The main reasons for this are that it is impossible to  
74 identify all insertion sequences based on limited reference databases and that *de novo*  
75 computational methods are ill-equipped to make fast or accurate predictions based on  
76 incomplete genomic assemblies. PaliDIS is a new software tool that is generating a larger,  
77 more comprehensive catalogue of insertion sequences based on a fast algorithm harnessing

78 genomic diversity in mixed microbial communities. This catalogue will enable genomic  
79 epidemiologists and researchers to annotate genomes for insertion sequences more  
80 extensively and advance knowledge of how insertion sequences contribute to bacterial  
81 evolution in general and antimicrobial resistance spread across microbial lineages in  
82 particular. This will be useful for genomic surveillance, and for development of microbiome  
83 engineering strategies targeting inactivation or removal of important transposable elements  
84 carrying antimicrobial resistance genes.

85

## 86 **Introduction**

87

88 Swapping genetic information between members of a microbial community, a mechanism  
89 referred to as horizontal gene transfer (HGT), is a key process in the microbiome. It allows  
90 for the spread of new genes and functionality throughout the community. The result of HGT  
91 can be acquisition of a new gene, duplication of an existing gene or even interruption of a  
92 current genes. The mechanisms that support HGT have been well described and involve the  
93 transfer of mobile genetic elements (MGEs). MGEs are best defined as broadly as possible,  
94 as any genetic element that can mediate its own transfer from one part of a genome to another  
95 or between different genomes. The most complex elements are conjugative plasmids and  
96 Integrative Conjugative Elements (ICEs) which can mediate their transfer between bacterial  
97 cells<sup>1</sup>. The simplest and most abundant MGEs are the insertion sequences which only contain  
98 enough genetic information for their own transposition. MGEs are best thought of as a  
99 continuum ranging from the relatively simple insertion sequences right up to conjugative  
100 elements and everything in between<sup>2</sup>. MGEs are crucially important in bacterial evolution as  
101 a result of the extensive diversity they generate, an aspect of this is their central role in the  
102 spread of antimicrobial resistance genes (ARGs) between microbial genomes.

103

104 Insertion sequences are short transposable elements between 700-2,500 bp in length  
105 containing genes that code for the proteins involved in their own transposition they are found  
106 in both chromosomes, ICEs and plasmids<sup>3</sup>. Most insertion sequences contain one or  
107 sometimes two genes encoding transposases, the most ubiquitous genes in prokaryotic and  
108 eukaryotic genomes<sup>4</sup>. Insertion sequences and transposons (transposons are defined as genetic  
109 elements that can transpose from one part of the genome to another but carry sequences other  
110 than those involved in transposition, unlike insertion sequences which just encode the genetic  
111 information for their own translocation) can be broadly classified by the amino acids in their

112 transposase, commonly DDE (aspartic acid, aspartic acid and glutamic acid), DEDD or HUH  
113 (two histidine residues separated by any large hydrophobic amino acid) motifs, and their  
114 mechanism of transposition (either conservative or replicative)<sup>5</sup>. Common DDE insertion  
115 sequences contain two inverted terminal repeats (ITRs) at each end of a 10-50 bp size DNA  
116 sequence that are reverse complement sequences of each other. Some insertion sequences are  
117 flanked by unique shorter direct repeat sequences, also known as target site duplications  
118 (TSDs), which are formed by the duplication of the insertion sequence target site upon  
119 insertion<sup>3</sup>. Unit transposons are a similar type of transposable element to insertion sequences  
120 containing a pair of ITRs but can also carry ARGs as well as transposases. For simplicity, the  
121 abbreviation “IS” will be used hereafter to mean insertion sequence or unit transposon. ARGs  
122 can also be carried by composite transposons that are bounded by two copies of two different  
123 insertion sequences which can move together in a single unit. A composite transposon can  
124 contain one or more passenger genes, such as ARGs, flanked by two insertion sequences and  
125 with two TSDs at both ends.

126

127 Microbial genomes can be annotated for ISs by querying reference databases of known  
128 transposable elements, such as ISfinder<sup>6</sup>, but these databases are small and do not represent  
129 many transposable elements in nature. As transposable elements are the most ubiquitous and  
130 abundant MGE, it is a continual effort to catalogue them all using common methods. Novel  
131 ISs containing ITRs can be detected using computational tools, such as EMBOSS<sup>7</sup>, that  
132 search for palindromic sequences representing ITRs<sup>8</sup>. However, transposable elements in  
133 isolated genomes that are assembled from short reads can be misassembled or incomplete,  
134 since assembly algorithms struggle to resolve repeated elements<sup>9</sup>. Additionally, ITR pairs are  
135 not typically exact reverse complements, and algorithms that only detect perfect palindromes  
136 may fail to identify many insertion sequences. Alternatively, novel ISs can be identified by  
137 manually searching for ITRs or flanking regions of interest (such as ARGs) using a genome  
138 browser, but this can be a difficult and tedious process. Alternatively, Hidden Markov  
139 Models (HMMs) have been used to identify transposases within these elements, include those  
140 without ITRs<sup>8</sup>. However, the presence of a transposase is not sufficient evidence for a  
141 transposition event to have occurred.

142

143 In this paper, we present a tool called PaliDIS (Palindromic Detection of Insertion  
144 Sequences) that finds ISs using an efficient maximal exact matching algorithm to identify  
145 ITRs across different genomic loci in reads sequenced from mixed microbial communities.

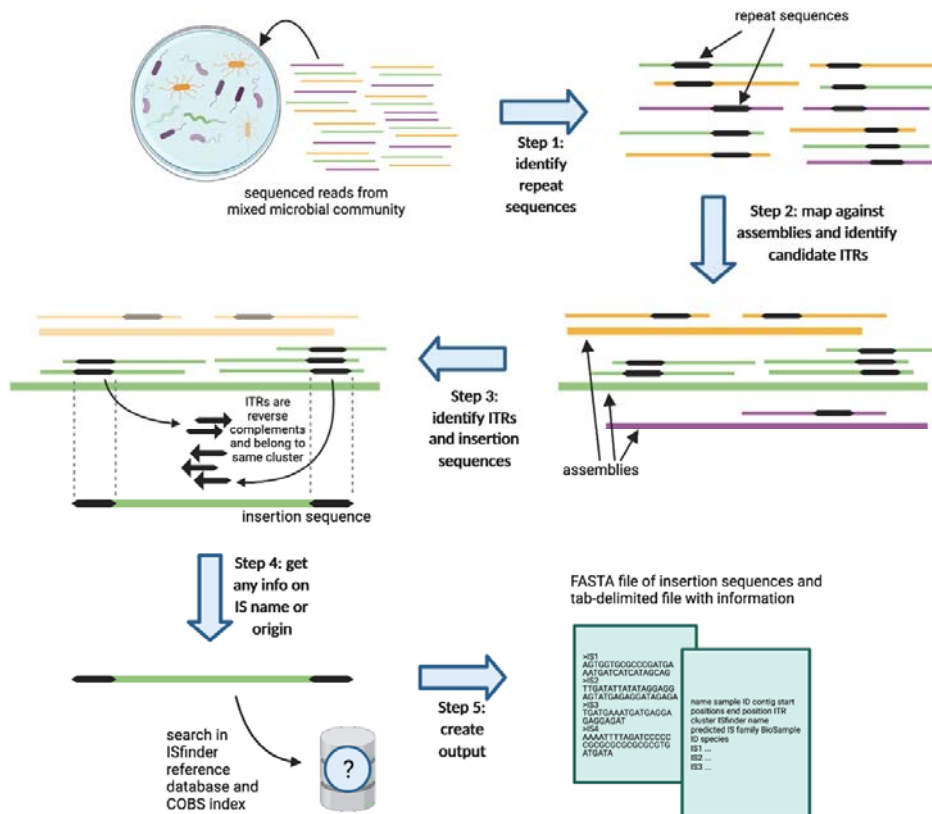
146 These ISs can then be pooled and clustered to create a non-redundant catalogue of ISs.  
147 PaliDIS can also predict the origins of these ISs by querying the them against ISfinder or a  
148 COmpact Bit-sliced Signature (COBS) index<sup>10</sup> of 661,405 microbial genomes<sup>11</sup>. Here, we  
149 present the theory and implementation of this tool on 266 short read metagenomes to  
150 generate 11,681 unique ISs included in the first release of the Insertion Sequence Catalogue  
151 (ISC). Beyond this paper, PaliDIS will continue generating ISs to expand ISC.

152

### 153 Theory and Implementation

154 PaliDIS is implemented as a Nextflow pipeline with all dependency software packaged in one  
155 container image. The input file of PaliDIS is a tab-delimited manifest text file that contains  
156 information on the read file IDs, the file paths to the read fastq.gz files, sample ID and file  
157 paths to the assemblies. The output files are a FASTA file of ISs and accompanying tab-  
158 delimited file of information. The following steps are also illustrated in Figure 1.

159



160

161 **Figure 1.** Steps summarising the PaliDIS software. Step 1: Reads from mixed microbial communities are pre-  
162 processed and run through an algorithm to identify reads containing repeat sequences. Step 2: Reads containing  
163 repeat sequences are mapped against the assemblies to find their positions and proximity filters applied to  
164 identify candidate ITRs. Step 3: Candidate ITRs are clustered. ISs are identified by ITRs that are of the same

165 cluster and are reverse complements of each other. Step 4 ISs are queried against existing databases to identify  
166 known ISs and predict their origin. Step 5: Final outputs of a FASTA file with insertion sequences and tab-  
167 delimited file with information are created.  
168

### 169 **Step 1: Reads from mixed microbial communities are pre-processed and run through** 170 **an algorithm to identify reads containing repeat sequences**

171 Firstly, the FASTQ files are converted to FASTA files with headers prepended with their  
172 sequence order (e.g. Seq1, Seq2 etc.). A software tool, called pal-MEM  
173 (<https://github.com/blue-moon22/pal-mem>), was developed and applied an efficient maximal  
174 exact matching algorithm<sup>12</sup> to identify repeat sequences that may represent ITRs. A maximal  
175 exact match (MEM) between two strings is an exact match (i.e. an exact local alignment),  
176 which cannot be extended on either side without introducing a mismatch (or a gap).  
177

### 178 **Preparing the reference and query data structures**

179 pal-MEM creates a reference hash table from the sequences for some integer  $k > 0$  defined by  
180 the user, in which  $k$ -mers are the keys and the corresponding occurring positions are their  
181 values. The nucleotides of  $k$ -mers are encoded as unique combinations of two bits (0 and 1),  
182 (where A is 00, C is 01, G is 10 and T is 11), reducing memory requirements. In addition, it is  
183 not required for all  $k$ -mers to be stored, reducing the demand on memory further. A  $k$ -mer is  
184 stored only when it has a position that is a multiple of  $(L - k) + 1$  (where  $k$  is the length of the  
185  $k$ -mer and  $L$  is the minimum ITR length), i.e.

186

$$187 \text{ (eq. 1) } b_r \leq j((L - k) + 1) \leq e_r - k + 1$$

188

189 where  $b_r$  and  $e_r$  are the start and end positions of a maximal exact match (MEM) and  $j \geq 1$ .  
190 The sequences are then also used to create a query data structure of unsigned 64-bit integers  
191 representing blocks of 32 nucleotides where each nucleotide is represented by two bits (A is  
192 00, C is 01, G is 10 and T is 11). Random 20-bit sequences are stored between the array of  
193 reads define their boundaries. The start and end positions for each read and random sequence  
194 are stored in another data structure.

195

### 196 **Applying the algorithm to find repeat sequences**

197 Each  $k$ -mer from the query read is looked up against the reference hash table to retrieve a  
198 matching  $k$ -mer. The first  $k$ -mer window starts from the beginning of the query and continues  
199 to shift every two bits, but skips the positions within the random sequences. These matching

200  $k$ -mers are then extended in both directions to make larger sequence matches until  
201 mismatches disrupt the extension, making a MEM. The algorithm performs this process using  
202 an interval halving approach. The sequence is extended to the left end position of the shortest  
203 of the two sequences. If there is no match, the extension is halved until a match is made. The  
204 extension is elongated by one nucleotide at a time until no more exact matches can be made.  
205 This is repeated on the right side. A repeat sequence is found once a MEM has a length  
206 greater than or equal to the minimum ITR length and less than or equal to the maximum ITR  
207 length as defined by the user. If a repeat sequence is found, pal-MEM will move on to the  
208 next read in the query, given it is expected that a read from short-read sequencing would  
209 contain only one ITR.

210

### 211 **Dealing with technical repeats from amplified read libraries**

212 Read libraries are dominated by technical as well as biological repeated sequences that are  
213 the result of sequencing amplified regions. To reduce the frequency of technical repeats being  
214 identified as biological repeats, MEMs are also excluded if their start or end positions are  
215 within a buffer length of 20 nucleotides (40 bits) from either end of the read. This model  
216 represents an alignment of the prefix or suffix of a read typical of a technical repeat.

217

### 218 **Step 2: Reads containing repeat sequences are mapped against the assemblies to find** 219 **their positions and proximity filters applied to identify candidate ITRs**

220

221 Reads containing repeat sequences identified in Step 1 are mapped using Bowtie2<sup>13</sup> against  
222 their associated assemblies. A Python script uses the output of Bowtie2 to identify mapped  
223 reads with candidate ITRs where the positions of the repeats are located between the  
224 minimum and maximum IS length as defined by the user.

225

### 226 **Step 3: Candidate ITRs are clustered and ISs are identified by ITRs that are of the** 227 **same cluster and are reverse complements of each other**

228 The candidate ITRs are clustered using CD-HIT-EST<sup>14</sup> where nucleotide sequences that meet  
229 a 1) sequence identity threshold  $c$ , 2) a global  $G$  1 or local alignment  $G$  0, 3) alignment  
230 coverage for the longer sequence  $aL$ , 4) alignment coverage for the shorter sequence  $aS$  and  
231 5) minimal alignment coverage control for the both sequences  $A$  (that can be specified by the  
232 user). The ISs are generated in a FASTA format with an accompanying tab-delimited file  
233 containing the sample ID, assembly name, start and end positions of the ITRs and their

234 cluster using a Python script. The ISs must contain ITRs that 1) belong to the same cluster, 2)  
235 are within the minimum and maximum specified ITR length, 3) are within the minimum and  
236 maximum IS length, and 4) are reverse complements of each other where the two sequences  
237 aligned using BLASTn<sup>15</sup> (with parameters *-task blastn -word\_size 4*) have  
238 “*Strand=Plus/Minus*” and “*Identities*” greater than or equal to the specified minimum ITR  
239 length.

240

#### 241 **Step 4: ISs are queried against existing databases to identify known ISs and predict** 242 **their origin**

243 ISs are queried against a non-redundant database of ISs from ISfinder<sup>6</sup> in 2020 using  
244 BLASTn (as documented here: [https://github.com/blue-](https://github.com/blue-moon22/PaliDIS/tree/master/db/ISfinder-sequences)  
245 [moon22/PaliDIS/tree/master/db/ISfinder-sequences](https://github.com/blue-moon22/PaliDIS/tree/master/db/ISfinder-sequences)). An IS that is a match with an ISfinder  
246 sequence is assigned as being a complete homolog if the alignment has an identity and a  
247 coverage of at least 99 %. Otherwise, the IS is assigned a predicted IS family. The origin of  
248 these ISs can be found by searching using cobs query<sup>10</sup> against a COBS index of microbial  
249 genomes with NCBI BioSample IDs. The taxonomy of those genomes containing these ISs  
250 are found by querying the BioSample IDs using a metadata retrieval tool, ffq<sup>16</sup>  
251 (<https://github.com/pachterlab/ffq>).

252

#### 253 **Step 5: Final output**

254

255 A Python script generates a FASTA file of ISs and a tab-delimited file of information  
256 including their header name, sample ID, contig, start and end positions of their ITRs on their  
257 contigs, ITR cluster, ISfinder name, predicted IS family from ISfinder, BioSample ID and  
258 origin species.

259

260

#### 261 **First release of the Insertion Sequence Catalogue using PaliDIS**

262

263 A catalogue of insertion sequences was generated using PaliDIS applied to 266 human oral  
264 and gut metagenomic reads from the Human Microbiome Project (Supplementary Data 1)<sup>17</sup>.  
265 The reads were quality controlled, filtered and assembled as previously described<sup>18</sup>. A total of  
266 25,650 ISs were identified from all these samples with PaliDIS v2.9.1 using default  
267 parameters (*--min\_itr\_length 25 --max\_itr\_length 50 --kmer\_length 15 --min\_is\_len 500 --*



268 `max_is_len 3000 --cd_hit_G 0 --cd_hit_c 0.9 --cd_hit_G 0 --cd_hit_aL 0.0 --cd_hit_aS 0.9 --`  
269 `cobs_threshold 1 --e_value 1e-50`) and a COBS index (specified by `--cobs_index`) of 661,405  
270 bacterial genomes created from European Nucleotide Archive (ENA) in 2018  
271 ([http://ftp.ebi.ac.uk/pub/databases/ENA2018-bacteria-661k/661k.cobs\\_compact](http://ftp.ebi.ac.uk/pub/databases/ENA2018-bacteria-661k/661k.cobs_compact))<sup>11</sup>.

272

273 The ISs were then clustered using CD-HIT-EST v4.8.1 (with a sequence identity threshold `-c`  
274 `0.99` and default parameters) to create the Insertion Sequence Catalogue (ISC). ISC contains a  
275 FASTA file of 11,681 unique ISs (<https://github.com/blue-moon22/ISC>) that were found in  
276 10,810 contigs across 253 (out of 266) samples (Supplementary Data 2).

277

278 In PaliDIS, the ISs were queried against existing databases, ISfinder and the COBS index, to  
279 identify known ISs and predict their origin. Only 8 ISs were found in ISfinder (ISBvu3,  
280 ISBf3, ISBf8, ISLh1, ISVesp1, ISBvu4, IS1249 and ISBuba1). Another 164 ISs were  
281 predicted to belong to 17 families in ISfinder (IS1182, ISAs1, IS110, IS5, IS630, ISLre2,  
282 IS256, IS200/IS605, IS30, IS3, ISL3, IS1595, IS982, IS1380, IS4, IS66 and IS481). 722 ISs  
283 were located in 16,803 unique sequenced sources (NCBI BioSample IDs). 505 of these ISs  
284 were found in 11,516 microbial isolate genomes with known taxonomy consisting of 61  
285 genera (Figure 2a) and 120 known species. 58 and 70 ISs originate from more than one genus  
286 (Figure 2b) and known species (not labelled *sp.*), respectively.

287

288 The IS shared across most genera is IS\_cluster\_192991\_length\_544 that was found within 18  
289 genera and 27 species in 290 unique biological sources (NCBI BioSample IDs) and was not  
290 identified in ISfinder. 21 assemblies out of 290 BioSamples were publicly available and  
291 downloaded (Supplementary Data 3). Despite being found in all 21 samples' reads,  
292 IS\_cluster\_192991\_length\_544 was only found in 17 assemblies (using blastn v2.13.0 with e-  
293 value cut-off 1e-10) (Supplementary Data 4). These 17 assemblies were then annotated with  
294 prokka v1.14.5 to find functional genes and genomic features<sup>19</sup>. The clinically relevant  
295 tetracycline-resistant gene *tet(O)* was found upstream of IS\_cluster\_192991\_length\_544 in 15  
296 assemblies across different classes: *Clostridium perfringens*, *Enterococcus gallinarum*,  
297 *Streptococcus agalactiae*, *Streptococcus azizii* and *Streptococcus suis* (Supplementary Figure  
298 1a, b, d-g, i-q). This IS may therefore have a role in the HGT of *tet(O)*, probably by  
299 mediating the transposition of *tet(O)* to conjugative elements. Use of PaliDIS and the  
300 associated ISC can thus reveal new and important ISs that function in the spread of AMR  
301 across species and lineages.



315 resistance genes. Here, we describe a tool and subsequent catalogue that enables this to  
316 proceed. PaliDIS is a tool that discovers novel ISs from mixed microbial communities by  
317 applying a fast maximal exact matching algorithm to identify ITRs. As a result, we have  
318 released the first version of ISC, a catalogue containing 11,681 ISs. Already, this is a  
319 valuable resource for researchers to search for ISs in isolated genomes. However, since  
320 PaliDIS was only applied to metagenomes sequenced from the healthy human oral cavity and  
321 stool samples, it is recommended ISC is used as a reference for annotating isolates sourced  
322 from human oral and stool samples.

323

324 The main limitation of the current ISC is that it only contains common DDE types of ISs with  
325 ITRs, although these mobile genetic elements make up a large proportion of ISs. PaliDIS is  
326 currently only equipped with discovering ISs with ITRs. We are planning to include other  
327 databases into the catalogue, such as ISfinder, and we invite the research community to  
328 contribute and submit ISs to the catalogue. Another limitation is that the catalogue currently  
329 contains ISs with ITRs that are 25 or greater nucleotides in length as generated by PaliDIS,  
330 although ITRs can be as short as 10 nucleotides in length. It is possible to run PaliDIS with a  
331 lower minimum ITR length threshold and smaller  $k$ -mer length, but at these smaller sizes, it  
332 becomes more computationally intensive, especially with more complex mixed microbial  
333 genomes. However, we will run PaliDIS with a lower minimum ITR length threshold on less  
334 complex genomes to discover ISs with smaller ITRs.

335

336 It is also important to note that all ISs in the catalogue contain a region that is flanked by  
337 ITRs within a 500 to 3000 bp proximity. Given the recursive mechanism of insertion events  
338 (i.e. ISs inserting within ISs), it is possible for a region to also contain another IS. Therefore,  
339 it is also possible for regions that have been lengthened by other insertion events to extend  
340 outside this proximity range and be missed by PaliDIS. Increasing the maximum IS length  
341 will account for this, and may be done for future iterations of the ISC.

342

343 In light of creating this tool and catalogue, we cannot turn a blind eye to the fact that  
344 disruptive sequencing technologies are advancing rapidly by becoming more accurate and  
345 generating longer reads. Very soon, it will be easy to apply tools for *de novo* discovery of ISs  
346 in genomic assemblies with resolved repeat regions, rather than relying on reference  
347 databases. However, reference catalogues, like ISC, generated from older data could be  
348 applied to monitor microbes that may acquire ISs that carry antimicrobial resistance genes,

349 which will be invaluable information for appropriate actions for tackling AMR. For instance,  
350 determining whether an IS carrying an ARG has already been in circulation that has to be  
351 controlled or is emerging that can be prevented from spreading early. Furthermore, having a  
352 catalogue of ISs will enable simple searches of genomic datasets, as well as comparisons with  
353 ISs from different species using less computationally intensive methods that are available to  
354 all in the community. We will continue to enrich the ISC towards a comprehensive catalogue  
355 by applying PaliDIS with different parameters to more mixed microbial genomes from a  
356 diverse range of sources, and encouraging submission of ISs from the scientific community.

357

### 358 **Acknowledgements**

359 The project was supported by the Centre for Host-Microbiome Interactions, King's College  
360 London, funded by the Biotechnology and Biological Sciences Research Council (BBSRC)  
361 grant BB/M009513/1 awarded to D.L.M. S.S. was supported by Engineering and Physical  
362 Sciences Research Council (EPSRC), EP/S001301/1, Biotechnology Biological Sciences Research  
363 Council (BBSRC) BB/S016899/1 and Science for Life Laboratory (SciLifeLab).

364

### 365 **Conflicts of interest**

366 The authors declare no conflicting interests.

367

### 368 **References**

- 369 1. Roberts, A. P. & Mullany, P. Tn916-like genetic elements: a diverse group of modular  
370 mobile elements conferring antibiotic resistance. *FEMS Microbiol. Rev.* **35**, 856–871  
371 (2011).
- 372 2. Frost, L. S., Leplae, R., Summers, A. O. & Toussaint, A. Mobile genetic elements: the  
373 agents of open source evolution. *Nat. Rev. Microbiol.* **3**, 722–732 (2005).
- 374 3. Mahillon, J. & Chandler, M. Insertion Sequences. *Microbiol Mol Biol Rev* **62**, 725–774  
375 (1998).
- 376 4. Aziz, R. K., Breitbart, M. & Edwards, R. A. Transposases are the most abundant, most  
377 ubiquitous genes in nature. *Nucleic Acids Res.* **38**, 4207–4217 (2010).

- 378 5. Partridge, S. R., Kwong, S. M., Firth, N. & Jensen, S. O. Mobile Genetic Elements  
379 Associated with Antimicrobial Resistance. *Clin. Microbiol. Rev.* **31**, (2018).
- 380 6. Siguier, P., Perochon, J., Lestrade, L., Mahillon, J. & Chandler, M. ISfinder: the  
381 reference centre for bacterial insertion sequences. *Nucleic Acids Res.* **34**, D32-36 (2006).
- 382 7. Rice, P., Longden, I. & Bleasby, A. EMBOSS: the European Molecular Biology Open  
383 Software Suite. *Trends Genet. TIG* **16**, 276–277 (2000).
- 384 8. Kamoun, C., Payen, T., Hua-Van, A. & Filée, J. Improving prokaryotic transposable  
385 elements identification using a combination of de novo and profile HMM methods. *BMC*  
386 *Genomics* **14**, 700 (2013).
- 387 9. Treangen, T. J. & Salzberg, S. L. Repetitive DNA and next-generation sequencing:  
388 computational challenges and solutions. *Nat. Rev. Genet.* **13**, 36–46 (2012).
- 389 10. Bingmann, T., Bradley, P., Gauger, F. & Iqbal, Z. COBS: A Compact Bit-Sliced  
390 Signature Index. in *String Processing and Information Retrieval* (eds. Brisaboa, N. R. &  
391 Puglisi, S. J.) 285–303 (Springer International Publishing, 2019). doi:10.1007/978-3-030-  
392 32686-9\_21.
- 393 11. Blackwell, G. A. *et al.* Exploring bacterial diversity via a curated and searchable snapshot  
394 of archived DNA sequences. *PLOS Biol.* **19**, e3001421 (2021).
- 395 12. Khiste, N. & Ilie, L. E-MEM: efficient computation of maximal exact matches for very  
396 large genomes. *Bioinformatics* **31**, 509–514 (2015).
- 397 13. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods*  
398 **9**, 357–359 (2012).
- 399 14. Fu, L., Niu, B., Zhu, Z., Wu, S. & Li, W. CD-HIT: accelerated for clustering the next-  
400 generation sequencing data. *Bioinformatics* **28**, 3150–3152 (2012).
- 401 15. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local  
402 alignment search tool. *J. Mol. Biol.* **215**, 403–410 (1990).

- 403 16. Gálvez-Merchán, Á., Min, K. H. (Joseph), Pachter, L. & Boeshaghi, A. S. Metadata  
404 retrieval from sequence databases with ffq. 2022.05.18.492548 (2022)  
405 doi:10.1101/2022.05.18.492548.
- 406 17. Turnbaugh, P. J. *et al.* The Human Microbiome Project. *Nature* **449**, 804–810 (2007).
- 407 18. Carr, V. R. *et al.* Abundance and diversity of resistomes differ between healthy human  
408 oral cavities and gut. *Nat. Commun.* **11**, 693 (2020).
- 409 19. Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinforma. Oxf. Engl.* **30**,  
410 2068–2069 (2014).
- 411