

1 **Palidis: fast discovery of novel insertion sequences**

2

3 **Authors**

4 Victoria R. Carr^{1,2,*}, Solon P. Pissis^{3,4}, Peter Mullany⁵, Saeed Shoaie^{2,6}, David Gomez-
5 Cabrero^{2,7,8}, David L. Moyes²

6

7 ¹Parasites and Microbes, Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton,
8 UK

9

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

12

13 ³Centrum Wiskunde en Informatica, Amsterdam, NL

14

15 ⁴Vrije Universiteit, Amsterdam, NL

16

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

19

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

22

23 ⁷Bioscience 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 ⁸Translational Bioinformatics Unit, Navarrabiomed, Complejo Hospitalario de Navarra
28 (CHN), Universidad Pública de Navarra (UPNA), IdiSNA, Pamplona, Spain

29

30 *Corresponding author: vc11@sanger.ac.uk

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 bioinformatics

37

38 pipeline called Palidis that recognises insertion sequences in metagenomic sequence data
39 rapidly by identifying inverted terminal repeat regions from mixed microbial community

39

40 genomes. Applying Palidis to 264 human metagenomes identifies 879 unique insertion
41 sequences, with 519 being novel and not previously characterised. Querying this catalogue

41

42 against a large database of isolate genomes reveals evidence of horizontal gene transfer
43 events across bacterial classes. We will continue to apply this tool more widely, building the

43

44 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

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

59 **Data Summary**

60 1. Palidis is available here: github.com/blue-moon22/palidis

61 2. The Insertion Sequence Catalogue is available to download here:
62 <https://github.com/blue-moon22/ISC>

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

65 4. The analysis for this paper is available here: [github.com/blue-](https://github.com/blue-moon22/palidis_paper_analysis)
66 [moon22/palidis_paper_analysis](https://github.com/blue-moon22/palidis_paper_analysis)

67 5. The output of Palidis that was run on these reads is available in Supplementary Data 2

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 generates a larger, more comprehensive catalogue of
77 insertion sequences based on a fast algorithm harnessing genomic diversity in mixed

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

84

85 **Introduction**

86

87 Swapping genetic information between members of a microbial community, a mechanism
88 referred to as horizontal gene transfer (HGT), is a key process in the microbiome. It allows
89 for the spread of new genes and functionality throughout the community. The result of HGT
90 can be acquisition of a new gene, duplication of an existing gene or even interruption of a
91 current gene. The mechanisms that support HGT have been well described and involve the
92 transfer of mobile genetic elements (MGEs). MGEs are best defined as broadly as possible,
93 as any genetic element that can mediate its own transfer from one part of a genome to another
94 or between different genomes. The most complex elements are conjugative plasmids and
95 Integrative Conjugative Elements (ICEs) which can mediate their transfer between bacterial
96 cells¹. The simplest and most abundant MGEs are the insertion sequences which only contain
97 enough genetic information for their own transposition. MGEs are best thought of as a
98 continuum ranging from the relatively simple insertion sequences right up to conjugative
99 elements and everything in between². MGEs are crucially important in bacterial evolution as
100 a result of the extensive diversity they generate, an aspect of this is their central role in the
101 spread of antimicrobial resistance genes (ARGs) between microbial genomes³.

102

103 Insertion sequences are short transposable elements between 700-2,500 bp in length
104 containing genes that code for the proteins involved in their own transposition they are found
105 in both chromosomes, ICEs and plasmids⁴. Most insertion sequences contain one or
106 sometimes two genes encoding transposases, the most ubiquitous genes in prokaryotic and
107 eukaryotic genomes⁵. Insertion sequences and transposons (transposons are defined as genetic
108 elements that can transpose from one part of the genome to another but carry sequences other
109 than those involved in transposition, unlike insertion sequences which just encode the genetic
110 information for their own translocation) can be broadly classified by the amino acids in their
111 transposase, commonly DDE (aspartic acid, aspartic acid and glutamic acid), DEDD or HUH

112 (two histidine residues separated by any large hydrophobic amino acid) motifs, and their
113 mechanism of transposition (either conservative or replicative)³. Common DDE insertion
114 sequences contain two inverted terminal repeats (ITRs) at each end of a 10-50 bp size DNA
115 sequence that are reverse complement sequences of each other. Some insertion sequences are
116 flanked by unique shorter direct repeat sequences, also known as target site duplications
117 (TSDs), which are formed by the duplication of the insertion sequence target site upon
118 insertion⁴. Unit transposons are a similar type of transposable element to insertion sequences
119 containing a pair of ITRs but can also carry ARGs as well as transposases³. For simplicity,
120 the abbreviation “IS” will be used hereafter to mean insertion sequence or unit transposon.
121 ARGs can also be carried by composite transposons that are bounded by two copies of two
122 different insertion sequences which can move together in a single unit⁶. A composite
123 transposon can contain one or more passenger genes, such as ARGs, flanked by two insertion
124 sequences and with two TSDs at both ends³.

125

126 Microbial genomes can be annotated for ISs by querying reference databases of known
127 transposable elements, such as ISfinder⁷, but these databases are small and do not represent
128 many transposable elements in nature. As transposable elements are the most ubiquitous and
129 abundant MGE, it is a continual effort to catalogue them all using common methods. Novel
130 ISs containing ITRs can be detected using computational tools, such as EMBOSS⁸, that
131 search for palindromic sequences representing ITRs⁹. However, transposable elements in
132 isolated genomes that are assembled from short reads can be misassembled or incomplete,
133 since assembly algorithms struggle to resolve repeated elements¹⁰. Additionally, ITR pairs
134 are not typically exact reverse complements, and algorithms that only detect perfect
135 palindromes may fail to identify many insertion sequences. Alternatively, novel ISs can be
136 identified by manually searching for ITRs or flanking regions of interest (such as ARGs)
137 using a genome browser, but this can be a difficult and tedious process. Alternatively, Hidden
138 Markov Models (HMMs) have been used to identify transposases within these elements,
139 include those without ITRs⁹. However, the presence of a transposase is not sufficient
140 evidence for a transposition event to have occurred.

141

142 In this paper, we present a tool called Palidis (Palindromic Detection of Insertion Sequences)
143 that finds ISs using an efficient maximal exact matching algorithm to identify ITRs across
144 different genomic loci in reads sequenced from mixed microbial communities. These ISs can
145 then be pooled and clustered to create a non-redundant catalogue of ISs. PaliDIS can also

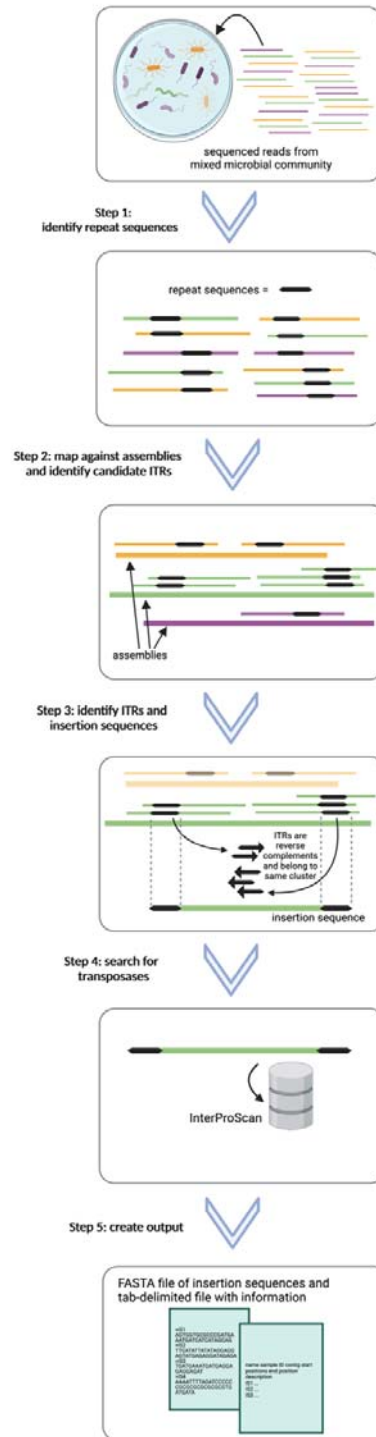
146 predict the origins of these ISs by querying them against ISfinder or a COmpact Bit-sliced
147 Signature (COBS) index¹¹ of 661,405 microbial genomes¹². Here, we present the theory and
148 implementation of this tool on 264 short read metagenomes to generate 879 unique ISs
149 included in the first release of the Insertion Sequence Catalogue (ISC). Beyond this paper,
150 Palidis will continue identifying ISs to expand ISC.

151

152 **Theory and Implementation**

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

158



159

160 **Figure 1.** Steps summarising Palidis. Step 1: Reads from mixed microbial communities are pre-processed and
161 run through pal-MEM to identify reads containing repeat sequences. Step 2: Reads containing repeat sequences
162 are mapped against the assemblies using Bowtie2 to find their positions and proximity filters applied to identify
163 candidate ITRs. Step 3: Candidate ITRs are clustered using CD-HIT-EST. ISs are identified by ITRs that are of
164 the same cluster and are reverse complements of each other. Step 4: Search of transposases using InterProScan.
165 Step 5: Final outputs of a FASTA file with insertion sequences and tab-delimited file with information are
166 created.

167

168 **Step 1: Reads from mixed microbial communities are pre-processed and run through**
169 **pal-MEM to identify reads containing repeat sequences**

170 Firstly, the FASTQ files are converted to FASTA files with headers prepended with their
171 sequence order (e.g. Seq1, Seq2 etc.). A software tool, called pal-MEM
172 (<https://github.com/blue-moon22/pal-mem>), was developed and applied an efficient maximal
173 exact matching algorithm¹³ to identify repeat sequences that may represent ITRs. A maximal
174 exact match (MEM) between two strings is an exact match (i.e. an exact local alignment),
175 which cannot be extended on either side without introducing a mismatch (or a gap).

176

177 **Preparing the reference and query data structures**

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

185

186 **(eq. 1)** $b_r \leq j((L - k) + 1) \leq e_r - k + 1$

187

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

194

195 **Applying the algorithm to find repeat sequences**

196 Each k -mer from the query read is looked up against the reference hash table to retrieve a
197 matching k -mer. The first k -mer window starts from the beginning of the query and continues
198 to shift every two bits, but skips the positions within the random sequences. These matching
199 k -mers are then extended in both directions to make larger sequence matches until
200 mismatches disrupt the extension, making a MEM. The algorithm performs this process using
201 an interval halving approach. The sequence is extended to the left end position of the shortest

202 of the two sequences. If there is no match, the extension is halved until a match is made. The
203 extension is elongated by one nucleotide at a time until no more exact matches can be made.
204 This is repeated on the right side. A repeat sequence is found once a MEM has a length
205 greater than or equal to the minimum ITR length and less than or equal to the maximum ITR
206 length as defined by the user. If a repeat sequence is found, pal-MEM will move on to the
207 next read in the query, given it is expected that a read from short-read sequencing would
208 contain only one ITR.

209

210 **Dealing with technical repeats from amplified read libraries**

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

216

217 **Step 2: Reads containing repeat sequences are mapped against the assemblies using** 218 **Bowtie2 to find their positions and proximity filters applied to identify candidate ITRs**

219 Reads containing repeat sequences identified in Step 1 are mapped using Bowtie2¹⁴ against
220 their associated assemblies. A Python script uses the output of Bowtie2 to identify mapped
221 reads with candidate ITRs where the positions of the repeats are located between the
222 minimum and maximum IS length as defined by the user.

223

224 **Step 3: Candidate ITRs are clustered using CD-HIT-EST. ISs are identified by ITRs** 225 **that are of the same cluster and are reverse complements of each other**

226 The candidate ITRs are clustered using CD-HIT-EST¹⁵ where nucleotide sequences that meet
227 a 1) sequence identity threshold c , 2) a global G 1 or local alignment G 0, 3) alignment
228 coverage for the longer sequence aL , 4) alignment coverage for the shorter sequence aS and
229 5) minimal alignment coverage control for the both sequences A (that can be specified by the
230 user). The ISs are generated in a FASTA format with an accompanying tab-delimited file
231 containing the sample ID, assembly name, start and end positions of the ITRs and their
232 cluster using a Python script. The ISs must contain ITRs that 1) belong to the same cluster, 2)
233 are within the minimum and maximum specified ITR length, 3) are within the minimum and
234 maximum IS length, and 4) are reverse complements of each other where the two sequences
235 aligned using BLASTn¹⁶ (with parameters `-task blastn -word_size 4`) have

236 “*Strand=Plus/Minus*” and “*Identities*” greater than or equal to the specified minimum ITR
237 length.

238

239 **Step 4: Search of transposases using InterProScan**

240 Candidate ISs are queried for transposases using the InterProScan¹⁷, a tool that combines
241 multiple search tools to predict protein family membership. Putative ISs must have
242 Transposase, Integrase-like and/or Ribonuclease H within at least one protein family
243 description.

244

245 **Step 5: Final output**

246 A Python script generates a FASTA file of ISs and a tab-delimited file of information
247 including: 1) their name (containing information on the length of the IS, the InterPro or
248 PANTHER accession(s) identified in Step 4 and their position(s)), 2) sample ID, 3) contig, 4)
249 start and end positions of the two ITRs on the corresponding contig, and 5) description of the
250 protein family represented by the accession(s).

251

252 **Using Palidis to create Insertion Sequence Catalogue v1.0.0**

253 A catalogue of insertion sequences was generated using Palidis applied to 264 human oral
254 and gut metagenomic reads from the Human Microbiome Project (Supplementary Data 1)¹⁸.
255 The reads were quality controlled, filtered and assembled as previously described¹⁹. A total of
256 2,517 ISs were identified from 1,837 contigs across 218 (out of 264) samples with Palidis
257 v3.1.0 using default parameters (`--min_itr_length 25 --max_itr_length 50 --kmer_length 15 --`
258 `min_is_len 500 --max_is_len 3000 --cd_hit_G 0 --cd_hit_c 0.9 --cd_hit_G 0 --cd_hit_aL 0.0 --`
259 `cd_hit_aS 0.9`) (Supplementary Data 2).

260

261 The ISs were then clustered using CD-HIT-EST v4.8.1 (with a sequence identity threshold `-c`
262 `0.95` and default parameters) to create the Insertion Sequence Catalogue (ISC). ISC contains a
263 FASTA file of 879 unique ISs between 524 and 2999 bp in length (Fig. 2a) and containing 87
264 unique transposases (Fig. 2b) (<https://github.com/blue-moon22/ISC>).

265

266 In order to identify ISs from ISC that have previously been discovered, ISC was queried
267 against ISfinder using their online BLAST search tool (<https://www-is.biotoul.fr/search.php>)
268 (with *e-value* 0.01 and using default parameters on 7th October 2022). 360 (41.0 %) ISs have
269 hits in ISfinder, while the remaining 519 are novel. 60 have a strict homology with ISs from

270 ISfinder ($e\text{-value} < 1e\text{-}50$), while the other 300 have a loose homology ($0.01 > e\text{-value} \geq 1e\text{-}$
271 50).

272

273 The origins of the ISs were determined by querying the ISC against a COBS index of
274 661,405 bacterial genomes (referred to as the 661k database) from European Nucleotide
275 Archive (ENA) ([http://ftp.ebi.ac.uk/pub/databases/ENA2018-bacteria-](http://ftp.ebi.ac.uk/pub/databases/ENA2018-bacteria-661k/661k.cobs_compact)
276 [661k/661k.cobs_compact](http://ftp.ebi.ac.uk/pub/databases/ENA2018-bacteria-661k/661k.cobs_compact))¹² using cobs query v0.1.2 (with $-t\ 0.9$ and using default
277 parameters).

278

279 155 ISs were located in 791 samples (NCBI BioSample IDs) from the 661k database.
280 Metadata was available for 785 for these samples from the 661k database study¹². 684
281 samples came from isolates that had an associated taxonomic identity, whereas the other 85
282 samples were sequenced from microbial communities or were unclassified (known as
283 “bacterium”). These 684 samples originate from 63 known genera (Fig. 3a) with 138 known
284 species (not labelled *sp.*) (Fig. 3b).

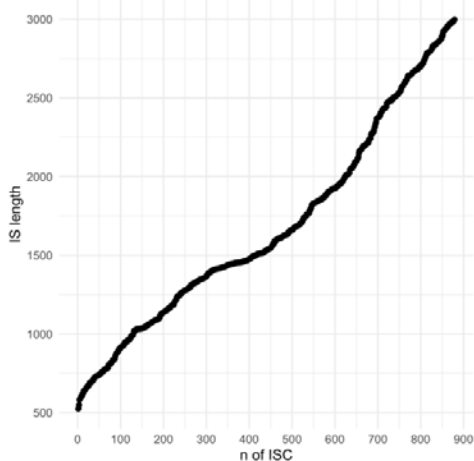
285

286 Many ISs were shared across bacterial genomes of different species, genera and classes. 52
287 and 70 ISs originate from more than one known genus (Fig. 4a) and known species (not
288 labelled *sp.*), respectively. The IS that is shared across most genera (21 genera and 46
289 species) and was also found in ISfinder as IS1249, is IS_length_1391-IPR001207_495_804
290 (Fig 4b, blue square), containing a Transposase, mutator type. The IS that is shared across the
291 most genera but was not found in ISfinder, is IS_length_2555-IPR036397_1291_1768-
292 PTHR35004_877_2065-IPR012337_1315_1720, containing a RV3428C-related Transposase
293 (Fig 3b, red square). Many ISs found in multiple *Bacteroides*, *Corynebacterium*,
294 *Curibacterium* and *Prevotella* species are solely represented by those not found in ISfinder,
295 suggesting that in its current release (7th October 2022), ISfinder database is
296 underrepresenting ISs in these genera (Fig. 4c).

297

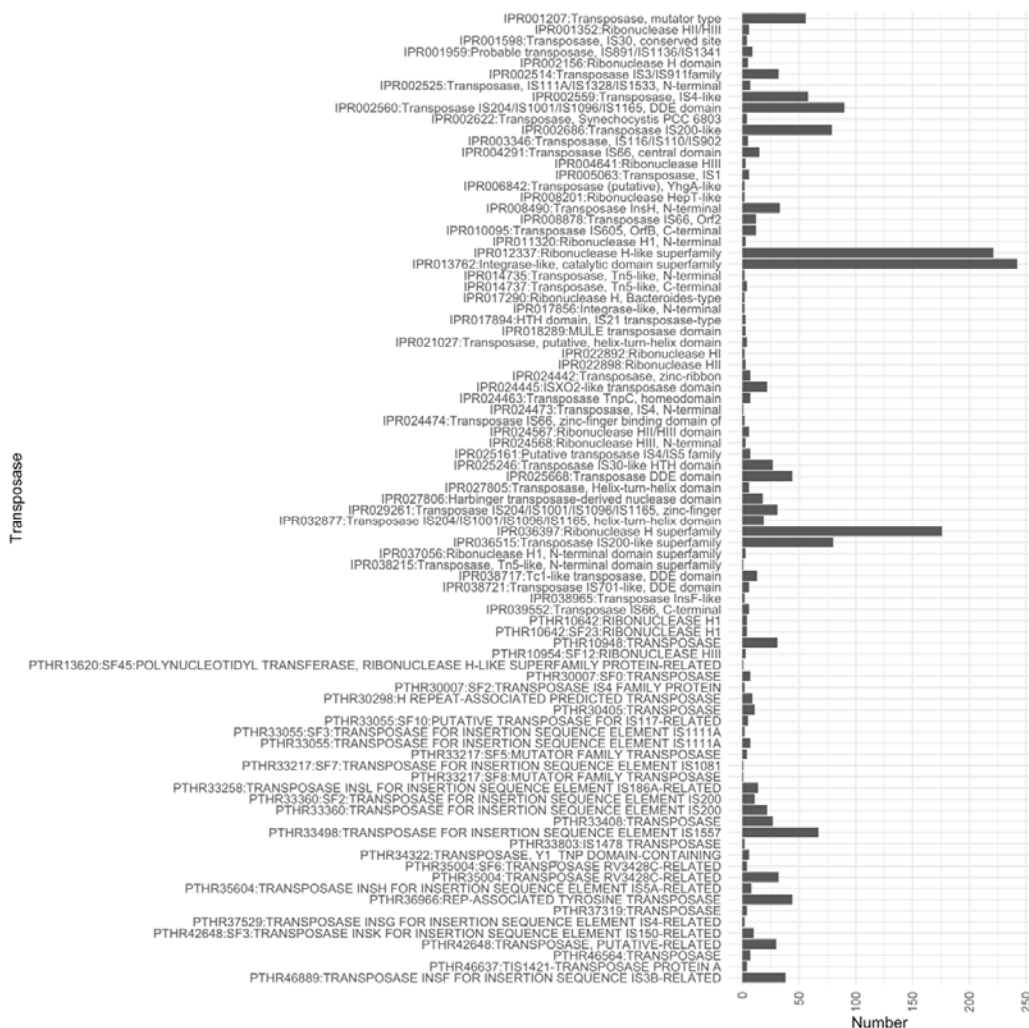
298

299 **a**



300

301 **b**



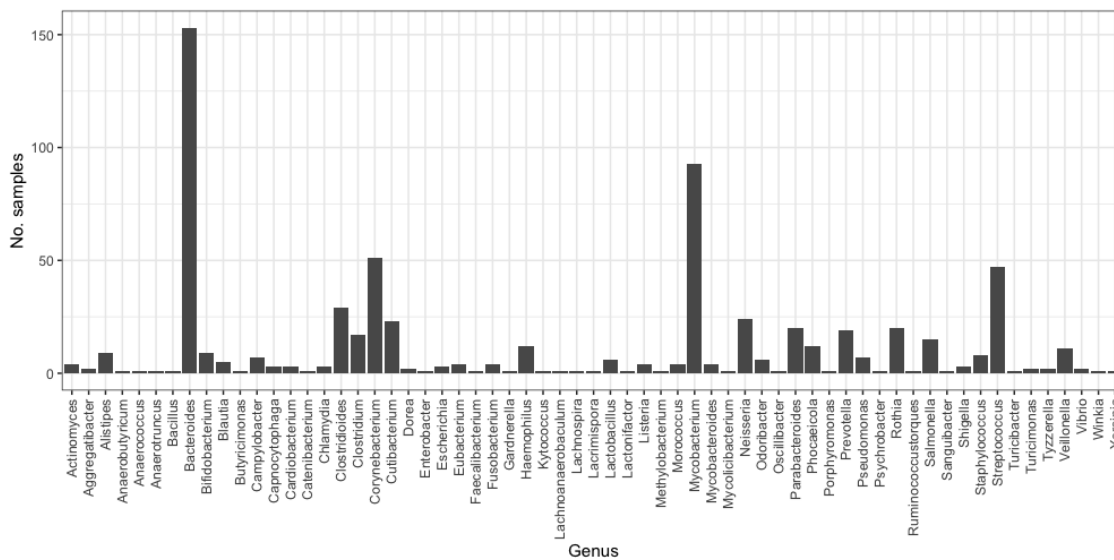
302

303

304

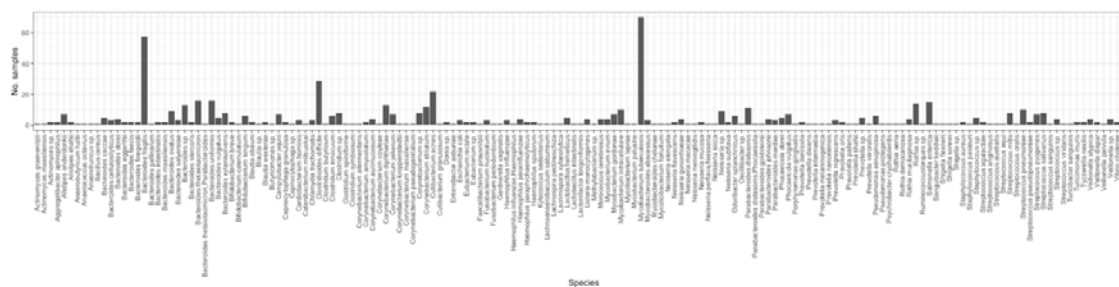
Figure 2. a) Length of insertion sequences and **b)** Number of transposases within the ISC

305 **a**



306

307 **b**

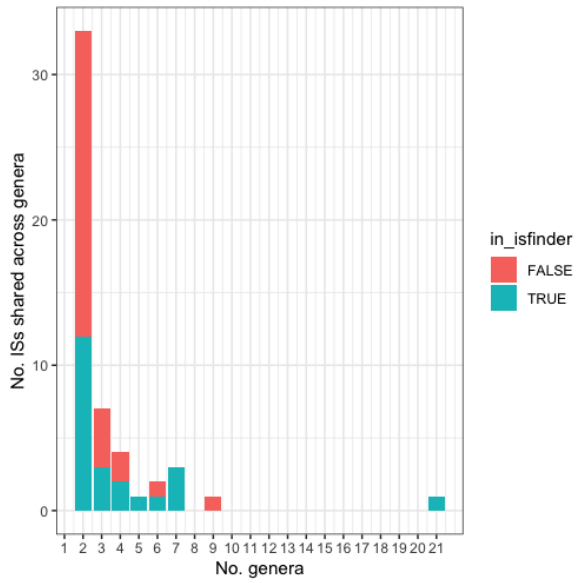


308

309 **Figure 3.** Number of samples from the 661k database that contain an IS found in across **a)** genera; **b)** species in
 310 the 661k database

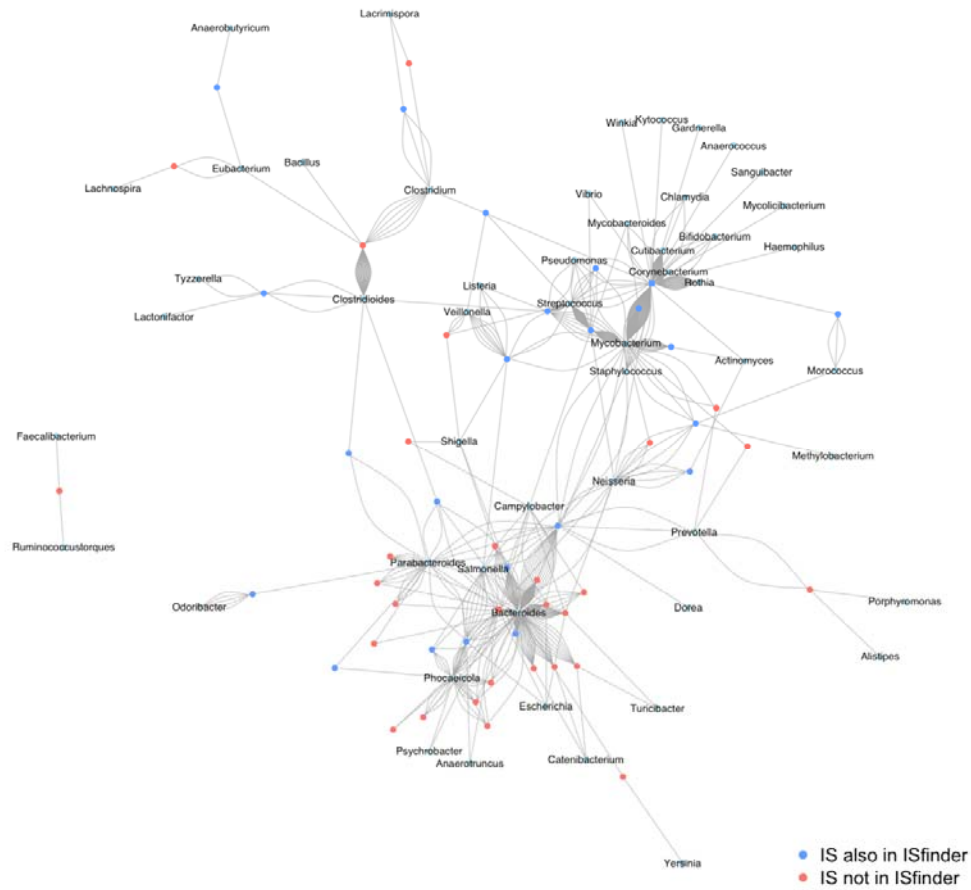
311

312 **a**



313

314 **b**



315

316

324 **Discussion**

325 Identification of transposable elements, including insertion sequences, in metagenomic
326 datasets is critical in our ability to accurately define the profile of mobile genetic elements.
327 In turn, accurate and complete characterisation of mobile genetic elements (i.e. the
328 mobilome) of a community is central to understanding the spread and epidemiology of
329 different genes in microbial communities, such as virulence genes and antimicrobial
330 resistance genes. Here, we describe a tool and subsequent catalogue that enables this to
331 proceed. Palidis is a tool that discovers novel ISs from mixed microbial communities by
332 applying a fast maximal exact matching algorithm to identify ITRs. As a result, we have
333 released the first version of ISC, a catalogue containing 879 ISs. Already, this is a valuable
334 resource for researchers to search for ISs in isolated genomes. However, since Palidis was
335 only applied to metagenomes sequenced from the healthy human oral cavity and stool
336 samples, it is recommended ISC v1.0.0 is used as a reference for annotating isolates sourced
337 from human oral and stool samples.

338

339 The main limitation of the current ISC is that it only contains common DDE types of ISs with
340 ITRs, although these mobile genetic elements make up a large proportion of ISs³. Palidis is
341 currently only equipped with discovering ISs with ITRs. We are planning to include other
342 databases into the catalogue, such as ISfinder, and we invite the research community to
343 contribute and submit ISs to the catalogue. Another limitation is that the catalogue currently
344 contains ISs with ITRs that are 25 or greater nucleotides in length as generated by Palidis,
345 although ITRs can be as short as 10 nucleotides in length. It is possible to run Palidis with a
346 lower minimum ITR length threshold and smaller *k*-mer length, but at these smaller sizes, it
347 becomes more computationally intensive, especially with more complex mixed microbial
348 genomes. However, we will run Palidis with a lower minimum ITR length threshold on less
349 complex genomes to discover ISs with smaller ITRs.

350

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

357

358 In light of current times, disruptive sequencing technologies are advancing rapidly by
359 becoming more accurate and generating longer reads. Very soon, Palidis could be applied to
360 longer reads of isolates to identify novel ISs, as well as generating reference databases (like
361 ISC) from mixed microbial genomes. However, the ISC could become a valuable resource
362 for querying microbial genomes for ARGs that have been acquired through transposition. We
363 will continue to enrich the ISC towards a comprehensive catalogue by applying Palidis with
364 different parameters to more mixed microbial genomes from a diverse range of sources, and
365 encouraging submission of ISs from the scientific community.

366

367 **Acknowledgements**

368 The project was supported by the Centre for Host-Microbiome Interactions, King's College
369 London, funded by the Biotechnology and Biological Sciences Research Council (BBSRC)
370 grant BB/M009513/1 awarded to D.L.M. S.S. was supported by Engineering and Physical
371 Sciences Research Council (EPSRC), EP/S001301/1, Biotechnology Biological Sciences
372 Research Council (BBSRC) BB/S016899/1 and Science for Life Laboratory (SciLifeLab).
373 S.P.P. is supported in part by the PANGAIA and ALPACA projects that have received
374 funding from the European Union's Horizon 2020 research and innovation programme under
375 the Marie Skłodowska-Curie grant agreements No. 872539 and 956229, respectively.

376

377 **Conflicts of interest**

378 The authors declare no conflicting interests.

379

380 **References**

- 381 1. Roberts, A. P. & Mullany, P. Tn916-like genetic elements: a diverse group of modular
382 mobile elements conferring antibiotic resistance. *FEMS Microbiol. Rev.* **35**, 856–871
383 (2011).
- 384 2. Frost, L. S., Leplae, R., Summers, A. O. & Toussaint, A. Mobile genetic elements: the
385 agents of open source evolution. *Nat. Rev. Microbiol.* **3**, 722–732 (2005).
- 386 3. Partridge, S. R., Kwong, S. M., Firth, N. & Jensen, S. O. Mobile Genetic Elements
387 Associated with Antimicrobial Resistance. *Clin. Microbiol. Rev.* **31**, (2018).

- 388 4. Mahillon, J. & Chandler, M. Insertion Sequences. *Microbiol Mol Biol Rev* **62**, 725–774
389 (1998).
- 390 5. Aziz, R. K., Breitbart, M. & Edwards, R. A. Transposases are the most abundant, most
391 ubiquitous genes in nature. *Nucleic Acids Res.* **38**, 4207–4217 (2010).
- 392 6. Tansirichaiya, S., Mullany, P. & Roberts, A. P. PCR-based detection of composite
393 transposons and translocatable units from oral metagenomic DNA. *FEMS Microbiol.*
394 *Lett.* **363**, (2016).
- 395 7. Siguier, P., Perochon, J., Lestrade, L., Mahillon, J. & Chandler, M. ISfinder: the
396 reference centre for bacterial insertion sequences. *Nucleic Acids Res.* **34**, D32–36 (2006).
- 397 8. Rice, P., Longden, I. & Bleasby, A. EMBOSS: the European Molecular Biology Open
398 Software Suite. *Trends Genet. TIG* **16**, 276–277 (2000).
- 399 9. Kamoun, C., Payen, T., Hua-Van, A. & Filée, J. Improving prokaryotic transposable
400 elements identification using a combination of de novo and profile HMM methods. *BMC*
401 *Genomics* **14**, 700 (2013).
- 402 10. Treangen, T. J. & Salzberg, S. L. Repetitive DNA and next-generation sequencing:
403 computational challenges and solutions. *Nat. Rev. Genet.* **13**, 36–46 (2012).
- 404 11. Bingmann, T., Bradley, P., Gauger, F. & Iqbal, Z. COBS: A Compact Bit-Sliced
405 Signature Index. in *String Processing and Information Retrieval* (eds. Brisaboa, N. R. &
406 Puglisi, S. J.) 285–303 (Springer International Publishing, 2019). doi:10.1007/978-3-030-
407 32686-9_21.
- 408 12. Blackwell, G. A. *et al.* Exploring bacterial diversity via a curated and searchable snapshot
409 of archived DNA sequences. *PLOS Biol.* **19**, e3001421 (2021).
- 410 13. Khiste, N. & Ilie, L. E-MEM: efficient computation of maximal exact matches for very
411 large genomes. *Bioinformatics* **31**, 509–514 (2015).

- 412 14. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods*
413 **9**, 357–359 (2012).
- 414 15. Fu, L., Niu, B., Zhu, Z., Wu, S. & Li, W. CD-HIT: accelerated for clustering the next-
415 generation sequencing data. *Bioinformatics* **28**, 3150–3152 (2012).
- 416 16. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local
417 alignment search tool. *J. Mol. Biol.* **215**, 403–410 (1990).
- 418 17. Jones, P. *et al.* InterProScan 5: genome-scale protein function classification.
419 *Bioinformatics* **30**, 1236–1240 (2014).
- 420 18. Turnbaugh, P. J. *et al.* The Human Microbiome Project. *Nature* **449**, 804–810 (2007).
- 421 19. Carr, V. R. *et al.* Abundance and diversity of resistomes differ between healthy human
422 oral cavities and gut. *Nat. Commun.* **11**, 693 (2020).
- 423