

1 **ISMMapper: Identifying insertion sequences in bacterial**
2 **genomes from short read sequence data**

3
4 **Jane Hawkey^{1,2§}, Mohammad Hamidian³, Ryan R. Wick¹, David J. Edwards¹,**
5 **Helen Billman-Jacobe², Ruth M. Hall³, Kathryn E. Holt¹**

6
7 ¹Department of Biochemistry and Molecular Biology, Bio21 Molecular Science and
8 Biotechnology Institute, University of Melbourne, Parkville, Victoria, Australia 3010

9 ²Faculty of Veterinary and Agricultural Science, University of Melbourne, Parkville,
10 Victoria, Australia 3010

11 ³School of Molecular Bioscience, The University of Sydney, Sydney, Australia 2006

12

13 [§]Corresponding author

14 Email addresses:

15 Jane Hawkey: hawkey.jane@gmail.com

16 Helen Billman-Jacobe: hbj@unimelb.edu.au

17 Mohammad Hamidian: mohammad.hamidian@sydney.edu.au

18 Ryan R. Wick: rwick@student.unimelb.edu.au

19 David J. Edwards: d.edwards2@student.unimelb.edu.au

20 Ruth M. Hall: ruth.hall@sydney.edu.au

21 Kathryn E. Holt: kholt@unimelb.edu.au

22 **Abstract**

23 **Background**

24 Insertion sequences (IS) are small transposable elements, commonly found in
25 bacterial genomes. Identifying the location of IS in bacterial genomes can be useful
26 for a variety of purposes including epidemiological tracking and predicting antibiotic
27 resistance. However IS are commonly present in multiple copies in a single genome,
28 which complicates genome assembly and the identification of IS insertion sites. Here
29 we present ISMapper, a mapping-based tool for identification of the site and
30 orientation of IS insertions in bacterial genomes, direct from paired-end short read
31 data.

32 **Results**

33 ISMapper was validated using three types of short read data: (i) simulated reads from
34 a variety of species, (ii) Illumina reads from 5 isolates for which finished genome
35 sequences were available for comparison, and (iii) Illumina reads from 7
36 *Acinetobacter baumannii* isolates for which predicted IS locations were tested using
37 PCR. A total of 20 genomes, including 13 species and 32 distinct IS, were used for
38 validation. ISMapper correctly identified 96% of known IS insertions in the analysis
39 of simulated reads, and 98% in real Illumina reads. Subsampling of real Illumina
40 reads to lower depths indicated ISMapper was reliable for average genome-wide read
41 depths >20x. All ISAbal insertions identified by ISMapper in the *A. baumannii*
42 genomes were confirmed by PCR. In each *A. baumannii* genome, ISMapper
43 successfully identified an IS insertion upstream of the *ampC* beta-lactamase that could
44 explain phenotypic resistance to third-generation cephalosporins. The utility of
45 ISMapper was further demonstrated by profiling genome-wide IS6110 insertions in

46 138 publicly available *Mycobacterium tuberculosis* genomes, revealing lineage-
47 specific insertions and multiple insertion hotspots.

48 **Conclusions**

49 ISMapper provides a rapid and robust method for identifying IS insertion sites direct
50 from short read data, with a high degree of accuracy demonstrated across a wide
51 range of bacteria.

52 **Keywords**

53 insertion sequence (IS), bacteria, genomics, short read analysis, tuberculosis,
54 antimicrobial resistance

55 **Background**

56 Insertions sequences (IS) are small transposable elements that encode the proteins
57 required for their own transposition. The ISfinder database [1] currently contains over
58 500 distinct IS. During transposition some IS create direct repeats, or target site
59 duplications, in the sequences into which they are integrating. The presence and
60 length of these duplications vary widely between IS and are characteristic of
61 individual IS [2]. Rates of transposition vary between IS and host species, but are
62 frequently in the order of the rate of nucleotide substitutions, making IS activity one
63 of the more dynamic evolutionary forces at play in many bacterial genomes. The
64 movement of IS can also have functional consequences for bacterial genomes. IS have
65 been implicated in large changes to genome structure, by expanding in copy number
66 in microbial genomes, with subsequent loss of IS resulting in inactivation of genes,
67 pseudogene formation, mediating deletion of intervening sequences between two
68 copies of the IS, or rearrangements of the genome [3].

69

70 In addition, IS insertions upstream of protein coding sequences can result in their
71 enhanced expression, leading to different phenotypes depending on the function of the
72 over-expressed gene. There are several known examples of IS-mediated gene
73 expression leading to clinically important increases in antimicrobial resistance. For
74 example, increased resistance to fluoroquinolones such as ciprofloxacin can result
75 from the insertion of *IS1* or *IS10* upstream of the *acrEF* efflux pump in *Salmonella*
76 Typhimurium [4], or the insertion of *IS186* upstream of the *acrAB* efflux pump in
77 *Escherichia coli* [5]. In *Acinetobacter baumannii*, insertion of *ISAbal* or *ISAbal25*
78 upstream of the intrinsic beta-lactamase *ampC* can cause resistance to third generation
79 cephalosporins including ceftazidime and cefotaxime [6, 7]. Insertions of the same IS
80 in nearby locations can generate a composite transposon, capable of mobilizing the
81 intervening sequence and transferring it to new genomic locations. For example, the
82 composite transposon *Tn6168* was generated spontaneously via insertions of *ISAbal*
83 on either side of *ampC*, including one copy of *ISAbal* that upregulates *ampC*
84 expression [8]. *Tn6168* has then transferred into different *A. baumannii* backgrounds,
85 conferring horizontally-acquired resistance to third generation cephalosporins [8].
86
87 IS insertions can also result in the upregulation of virulence genes in clinically
88 important human pathogens. For example, an outbreak of tuberculosis in Spain in the
89 1990s was associated with the B strain of *Mycobacterium bovis* carrying an insertion
90 of *IS6110* in the promoter region of the virulence gene *phoP*, resulting in its
91 upregulation [9]. In *Neisseria meningitidis*, insertion of *IS1301* in the middle of the
92 capsule locus has been shown to cause increased expression of operons on either side
93 of the IS, contributing to protection from the human immune system and enhanced
94 pathogenicity [10]. IS have also been shown to enhance niche adaptation in bacteria,

95 for example IS1247 insertion upstream of *dhlB* in *Xanthobacter autotrophicus* results
96 in increased resistance to bromoacetate [11]. This region has also been mobilised by
97 the IS and transferred to a plasmid [11]. In *E. coli*, IS3 has been shown to up-regulate
98 threonine expression, allowing the bacteria to adapt to a low-carbon environment and
99 utilise threonine as its sole carbon source [12].

100

101 The profiling of IS insertion patterns has been used for typing purposes in numerous
102 bacterial species of importance to human health. For example, copy number and
103 position of IS200 in *Salmonella enterica* [13], IS6110 in *Mycobacterium tuberculosis*
104 [14], IS1004 in *Vibrio cholerae* [15] and ISAbal in *A. baumannii* [16] has been used
105 to profile these bacterial pathogens, allowing the identification and tracking of distinct
106 subtypes. To date, IS-based typing schemes for various bacteria have relied on
107 digesting the genome followed by either hybridizing IS probes to fragments in a gel or
108 PCR probing [13-15]. The detection of precise insertion sites can be achieved using
109 PCR, and may be done for typing purposes [17] or for the detection of functionally
110 important insertions [7, 9].

111

112 With the advent of cheap high-throughput short-read sequencing, whole genome
113 sequencing (WGS) of bacteria is increasingly common and is replacing traditional
114 methods for characterizing and typing bacterial genomes. Unfortunately the detection
115 of IS is complicated wherever read lengths are shorter than the length of the IS, as is
116 the case for platforms that are currently most widely used – Illumina and Ion Torrent.
117 IS insertion sites can readily be identified in finished bacterial genomes or in draft
118 assemblies of genomes with single-copy IS, using tools such as nucleotide BLAST or
119 ISfinder [1]. However where multiple copies of the same IS are present within a

120 single genome (including on the chromosome and/or plasmids), this complicates
121 assembly of short-read data and makes IS insertion sites difficult to identify reliably.
122 The IS detection problem can be resolved using long-read sequencing technologies
123 such as the SMRT Cell (Pacific Biosciences) or MinION (Oxford Nanopore)
124 platforms; however given the relative cost efficiency and reliability of short-read
125 sequencing, together with the current widespread use of Illumina for bacterial WGS
126 and wealth of available short-read data for clinically important bacteria, there remains
127 a need for a simple tool to identify IS insertion sites from short-read data.

128

129 Several studies report the use of mapping-based approaches to identify IS insertion
130 sites from bacterial short-read data [18, 19], however none provide software code or
131 validation of the approach used. There are tools available for detecting transposons or
132 structural variation in genomes, for example MindTheGap [20] and BreakDancer
133 [21], however these do not perform well in the identification of IS in bacterial
134 genomes nor were they designed to do so. Some programs could potentially be used
135 for this purpose, such as RelocaTE [22] and RetroSeq [23], however these require
136 additional input or prior knowledge about the IS which may not always be available.
137 TIF (Transposon Insertion Finder) [24] and *breseq* [25] could potentially be used for
138 the detection of IS insertion sites in bacterial genomes, however they were not
139 designed specifically for this purpose and did not perform well on our data sets (see
140 Results).

141

142 Here we present a rapid and robust tool for accurate detection of IS sequences,
143 including insertion site and orientation, direct from short-read data. The method is
144 freely available in the form of open-source code called ISMapper, and here we

145 validate its use via analysis of simulated and real short-read data from a range of IS
146 and bacterial species. ISMapper requires short reads and query IS sequences as input,
147 and can be used either for typing against a reference genome or to assist with manual
148 resolution of complex short-read assemblies.

149 **Implementation**

150 An overview of the ISMapper workflow is shown in **Figure 1**. ISMapper takes as
151 input: (i) a set of paired end Illumina reads for an isolate of interest, (ii) an IS query
152 sequence in fasta format, and (iii) either a reference genome (for typing) or an
153 assembly of the read set (for assembly improvement), in GenBank or FASTA format
154 (**Figure 1a**). Paired end Illumina reads are mapped to the IS query sequence using
155 BWA-MEM (v0.7.5a or later) [26]. From the resulting alignment file (SAM format),
156 unmapped reads whose pairs map to the end of the IS query sequence (that is, reads
157 representing the sequences directly flanking the IS) are extracted using SAMtools
158 view (v0.1.19 or later) [27] to retrieve reads based on SAM flags (**Figure 1b**).
159 Specifically, left flanking reads (taking input sequence as left to right) are extracted
160 using flag ‘-f 36’ and right flanking reads are extracted using flag ‘-F 40 -f 4’ and
161 stored in separate BAM files, which are then converted to FASTQ format using
162 BedTools (v2.20.1) [28]. In addition, Samblaster (v0.1.21) [29] is used to extract from
163 the SAM file any reads that map to the end of the IS and extend into the neighbouring
164 sequence (i.e. “soft clipped” reads, **Figure 1b**). The resulting FASTQ file is filtered
165 using BioPython to extract the soft clipped portion of reads, where those sequences fit
166 a specified size range (default 5- 30 bp). The resulting sequences are sorted into left
167 and right flanking sequences; these are each mapped separately to the reference
168 genome or assembly using BWA-MEM, to identify the location(s) of the query IS in

169 the genome under analysis (**Figure 1c**). Insertion site information is extracted from
170 the resulting alignments using BedTools (coverage command) to summarise coverage
171 of the reference by left and right flanking reads; these are filtered by read depth
172 (default, minimum read depth $\geq 6x$) to minimize false positive hits, and regions that
173 overlap or are separated by a short distance (default, ≤ 100 bp) are merged using
174 BedTools (merge command). Pairs of left and right flanking regions that likely
175 represent either side of the same IS insertion are identified on the basis of positional
176 information, using BedTools (intersect and closest commands). Left and right regions
177 that overlap are considered to indicate a novel IS insertion not present in the
178 reference, with the overlap resulting from target site duplication arising during IS
179 transposition (**Figure 1c**, novel site). Where left and right regions are separated by a
180 sequence that is approximately the length of the IS query, the intervening sequence is
181 extracted and compared to the IS query using nucleotide BLAST+ (v2.2.25 or later)
182 [30] to confirm whether this is a known insertion site that is present in the reference
183 (**Figure 1c**, known site).

184

185 Extensive testing of ISMapper revealed that it was sometimes unable to resolve IS
186 positions that were adjacent to a repeat region (segments of DNA that were repeated
187 multiple times around the genome; see Results). This is because when the IS-flanking
188 reads were mapped back to the reference genome, those that belonged to the
189 neighbouring multi-copy sequence were randomly assigned by BWA-MEM to the
190 various locations of the repeat sequence, resulting in low read depth at the ‘true’ IS-
191 adjacent copy of the multi-copy sequence, which can fall below the minimum depth
192 filter (**Figure 2**). In such cases, the sequence on the other side of the IS is usually not
193 a multi-copy sequence and thus does not suffer the same problem, and so is usually

194 identified as a confident IS-flanking region without a corresponding partner region
195 (**Figure 2**, purple block). Therefore, when ISMapper identifies an un-partnered IS-
196 flanking region, it checks the original alignments for evidence of a nearby low-
197 coverage partner region that failed to pass the depth filter and returns this as a
198 potential but uncertain IS location, indicated by a ‘?’ character in the results table
199 (**Figure 2**, green block).

200

201 ISMapper generates two main output files summarizing the results: (i) a GenBank file
202 of the reference sequence, annotated with the IS-flanking regions and (ii) a table
203 indicating the locations and characteristics of each IS-flanking region identified
204 (**Figure 1d**). The table includes details of the location of the IS insertions; the
205 distance between the left and right flanking regions (where a negative number
206 indicates an overlap of left and right regions, indicating the size and sequence of the
207 target site duplication); a call as to whether the insertion is present in the reference or
208 is a novel insertion site (and, where the insertion site is present in the reference, the
209 percent coverage and sequence homology with the IS query); and details of the
210 gene(s) closest to the IS insertion site (including locus tag, product, gene name and
211 distance from the IS to the start codon). Insertions are also marked to indicate less
212 confident calls. A ‘*’ indicates an imprecise hit; i.e., where the gap between left and
213 right regions is larger than expected for a novel insertion, but is not consistent with an
214 IS insertion at that location in the reference. A ‘?’ indicates an uncertain hit, where
215 only one end (left or right of the predicted insertion) passes the minimum read depth
216 threshold; this often occurs when the IS is inserted within or adjacent to a multi-copy
217 sequence, as described above (**Figure 2**). When run in assembly improvement mode,
218 the table produced is simpler and indicates which contigs are predicted to end

219 adjacent to the IS (indicating left or right orientation), assisting the user to decide
220 whether some contigs could be joined together based on the available IS evidence.
221
222 ISMapper is lightweight code – a test run on a laptop computer (MacBook Air) with
223 8GB of RAM and a 1.3GHz i5 processor was able to analyse a read set comprising
224 2.5 million 100 bp paired-end reads in approximately ten minutes for a single IS
225 query. Because ISMapper analyses each read set and query IS independently,
226 screening of multiple read sets and query IS can be easily performed in parallel across
227 multiple cores. To facilitate easy compilation of results from multiple jobs, ISMapper
228 includes a Python script to cross-tabulate results from multiple read sets, generating a
229 single summary table per query IS (script ‘compiled_table.py’).

230 **Results and Discussion**

231

232 **Validation of IS detection using simulated reads**

233 Nine publicly available genomes from a variety of bacterial genera, and including
234 both chromosomes and plasmids, were downloaded from NCBI (**Table 1**). ISfinder
235 [1] was used to identify the IS present in each genome sequence. All sequences that
236 had >50% identity to a sequence in ISfinder and were present in at least two copies
237 were extracted as query IS for testing with ISMapper. Nucleotide BLAST+ was used
238 to confirm the precise locations and orientations for each query IS in all genomes
239 (total 251 insertions of 17 distinct IS, see **Table 1**). Short reads (100 bp) were
240 simulated from each genome sequence using the wgsim command in SAMtools
241 (v0.1.19), with default parameter settings.

242

243 ISMapper was run with default parameter settings on each combination of genome,
244 query IS and simulated reads. ISMapper was able to accurately locate each IS position
245 and its orientation (ranging between 2 and 61 positions per genome) for the majority
246 of genomes (**Table 1**). In total, 96.4% of IS insertions were correctly detected. The
247 exceptions occurred in three genomes (*K. pneumoniae* plasmid pNDMAR, *Y. pestis*
248 CO92 and *E. coli* O157:H7), in which ISMapper correctly identified 151 IS insertion
249 sites and failed to identify nine (94% detection). Closer inspection revealed that the
250 missed IS were each located next to multi-copy repeat sequences, complicating the
251 second mapping step as discussed above and outlined in **Figure 2**. Switching on
252 reporting of all alignments above a mapping score threshold of 30 (-a and -T 30 in
253 BWA-MEM) enabled the detection of a further IS100 site in *Y. pestis*. By default this
254 option is turned off in ISMapper as it tends to create noise in the mapping, making it
255 more difficult to distinguish true and false positives; however this can be useful if an
256 IS site of interest is known or suspected to be flanked by further repeats.

257

258 **Validation of IS detection using real Illumina read sets derived from isolates** 259 **with finished genomes**

260 Next we validated ISMapper using six genomes for which both Illumina read data and
261 finished genomes were publicly available (Table 2). Each finished genome sequence
262 was analysed with ISfinder [1] to identify query IS for testing as described above, and
263 nucleotide BLAST was used to confirm the precise locations and orientations of each
264 IS in each genome. The resulting test set comprised 106 insertions of 14 query IS.
265 Using default settings, ISMapper was able to accurately identify each IS insertion site
266 and its orientation, between 2 and 26 per genome, for the majority of genomes (Table
267 2). In total, 104 (98%) IS insertions were correctly detected by ISMapper. Three of

268 four IS431mec insertions in *Staphylococcus aureus* TW20 were correctly detected,
269 however the fourth was missed by ISMapper as it was flanked by another IS431mec
270 and further repeat sequences. Two of three IS1 insertions in *Salmonella* Typhi CT18
271 were correctly detected however a third, located between *tviE* and *tviD*, was
272 problematic. ISMapper was able to identify the region flanking the IS at *tviE*, but was
273 unable to detect the corresponding region in *tviD*. Inspection of a BWA-derived
274 alignment of the full Illumina read set to the CT18 chromosome showed that the
275 entire region spanning from *tviD* to *tviA* was devoid of sequence reads, suggesting
276 that this region may have been deleted during culture in the laboratory prior to the
277 extraction of DNA for Illumina sequencing. This region encodes the biosynthesis of
278 the Vi capsule of *S. Typhi*, and is known to be lost sporadically during culture [31].
279 This illustrates that situations where one end of the IS is detected but the other is not
280 can often be ‘accurate’ in the sense that the result reflects underlying structural
281 variation in the genome, including potentially IS-mediated deletions.

282

283 **Detection of antibiotic resistance-mediating IS insertions in *Acinetobacter***
284 ***baumannii*, confirmed by PCR**

285 The genomes of seven ceftazidime resistant *A. baumannii* isolates, belonging to
286 global clone (GC) 1, were sequenced via Illumina HiSeq to generate 100 bp paired
287 end reads. Resistance gene screening of the Illumina data using SRST2 [32] and the
288 ARG-Annot database [33] confirmed earlier PCR data indicating that none of these
289 isolates carried acquired extended spectrum beta-lactamase (ESBL) genes that can
290 confer resistance to third-generation cephalosporins. However, it is known that the
291 insertion of ISAbal1 upstream of the intrinsic chromosomally encoded *ampC* beta-

292 lactamase gene can cause increased resistance to third-generation cephalosporins in *A.*
293 *baumannii* [6].

294

295 We used ISMapper to screen for the ISAbal query sequence (accession AY758396),
296 sourced from ISfinder [1]. Using default parameters, ISMapper identified ISAbal
297 insertions in all seven GC1 genomes. IS positions were assessed relative to the
298 genome sequence of *A. baumannii* GC1 reference A1 (accession CP010781).

299 ISMapper found between 3 and 5 ISAbal insertions in each GC1 isolate, including an
300 insertion upstream of *ampC* in all 7 genomes that was in the orientation required to
301 induce upregulation and explain the observed cephalosporin resistance phenotype
302 (**Figure 3**). In addition, out of 29 total ISAbal insertions, ISMapper was able to
303 correctly identify 26 target site duplications (9 bp in the case of ISAbal). All ISAbal
304 insertions were novel compared to the reference genome A1 (**Figure 3**) and were
305 confirmed using PCR, as described in [6].

306

307 **Impact of read depth on ISMapper performance**

308 To test the effect of read depth on the performance of ISMapper, each of the seven
309 GC1 *A. baumannii* read sets were randomly subsampled to depths of approximately
310 10x, 15x, 20x, 25x, 50x, 75x and 100x, with ten replicates per depth level per read set.
311 ISMapper was then run using default settings to screen for ISAbal insertions. The
312 results indicated that at mean genome-wide read depths of approximately 20x,
313 ISMapper was able to identify 95% of insertions correctly (**Figure 4**). However, all of
314 these calls were either imprecise (gap size larger than expected) or uncertain (high
315 coverage end paired with a low coverage end). An average genome-wide read depth
316 of ~50x was required to find all insertions, with confident calls for >60%, however

317 there was clearly some variation depending on read quality (**Figure 4**). To achieve
318 100% detection with high confidence, average genome-wide read depths of >75x
319 were required (**Figure 4**).

320

321 **Comparison of ISMapper with TIF and *breseq***

322 The seven *A. baumannii* GC1 genomes were used to test both *breseq* [25] and TIF
323 (Transposon Insertion Finder) [24]. *breseq* uses split read mapping to a reference
324 genome along with statistical models to determine new junctions and deletions in the
325 isolates of interest. As input, *breseq* takes paired end reads in FASTQ format, and a
326 reference genome in Genbank format. The *breseq* manual indicates that new
327 insertions of mobile elements can be determined by looking for ‘JC JC’ evidence types
328 in the final html output. All seven *A. baumannii* isolates were screened using default
329 parameters and the reference genome A1 (accession CP010781). In all cases, *breseq*
330 was unable to identify any mobile element insertions, including no structural variation
331 at the known ISAbal insertion sites, although many other types of structural variation
332 were detected.

333

334 TIF requires as input paired end reads (FASTQ format), the head and tail sequences
335 (approximately 17 bp) of the IS of interest as well as the size of the target site
336 duplication the IS makes during transposition. TIF uses regular expressions to search
337 for the head and tail sequences in the reads, and these reads are then extracted and
338 grouped by their target site duplications. Unfortunately, following communication
339 with the authors, we were unable to get TIF to output any results using our data. Other
340 disadvantages of TIF are the requirements to (i) specify the size of target site
341 duplications (which not all IS make and is not always known), (ii) manually extract

342 subsequences of the IS rather than inputting the complete sequence, and (iii) manually
343 edit a Perl script in order to specify inputs to the program.

344

345 **Example use case: Exploration of IS6110 insertions in *Mycobacterium***

346 ***tuberculosis***

347 While IS insertions are thought to be important for shaping the evolution of bacteria
348 in a variety of ways, high-resolution comparative genomic studies of bacterial
349 pathogens have largely ignored IS due to the difficulties associated with accurate
350 detection of insertion sites from high-throughput short read data. An important
351 example is IS6110 in *M. tuberculosis* [34]. Profiling of IS6110 insertions using PCR
352 and restriction fragment based polymorphism (RFLP) based methods has been
353 reported for typing purposes [35], and specific insertions have been linked to
354 clinically relevant changes in function including in outbreak strains [9, 36, 37]
355 However while numerous studies have reported the genomic analysis of hundreds of
356 *M. tuberculosis* isolates sequenced on the Illumina platform, these have not included
357 analysis of IS6110 insertions. Thus, to demonstrate the utility of ISMapper for
358 comparative profiling of IS in an important bacterial pathogen, we analysed the
359 distribution of IS6110 within 138 publicly available genomes representing the major
360 lineages of *M. tuberculosis* [38]. Paired-end Illumina reads were downloaded from
361 NCBI (ERP001731). A core genome phylogeny was generated from these reads by
362 SNP (single nucleotide polymorphism) calling against reference genome H37Rv
363 (accession NC_000962) (methods as described in [39]), followed by maximum
364 likelihood phylogenetic inference on the SNP alignment using RAxML (GTR+G
365 substitution model, 1000 bootstraps) to build a genome-wide phylogenetic tree.

366 ISMapper was run with default settings to screen for insertions of *IS6110* (accession
367 X17348) in each read set, relative to reference genome H37Rv.
368
369 A total of 392 unique *IS6110* insertion sites were identified by ISMapper,
370 approximately one per 10 kbp of the 4.4 Mbp reference genome. The frequency of
371 each insertion within each of the six main global lineages is shown in **Figure 5b**. The
372 data indicate multiple lineage-specific *IS6110* insertions in lineages 2-6, but none that
373 were shared by multiple lineages, suggesting that *IS6110* insertions began to
374 accumulate only after *M. tuberculosis* diverged into these distinct lineages. Isolates in
375 the “modern” lineages 2-4 and in the West African lineage 5 had more *IS6110*
376 insertions overall, with far fewer insertions observed in the “ancient” East and West
377 African lineages 1 and 6 (**Figure 5c**). Lineage 2, which includes the highly successful
378 Beijing sublineage, had the highest number of *IS6110* although it was not the most
379 common lineage in the collection (n=23); it could be that these insertions contribute to
380 the adaptive fitness of the Beijing lineage.

381
382 The spatial distribution of unique *IS6110* insertions within the *M. tuberculosis*
383 genome (**Figure 5d**) revealed several clusters of insertions detected by ISMapper.
384 Many of these clusters comprised multiple independent insertions into PE/PPE genes
385 (which are surface-associated and interact with the host immune system), as well as
386 the membrane associated proteins *mmpS1* and *mmpL12*. There was substantial
387 clustering of *IS6110* insertions interrupting genes encoding the CRISPR machinery,
388 which is involved in immunity to bacteriophage and other foreign DNA. Further, all
389 three phospholipase genes, which are involved in virulence by inducing cell death in
390 macrophages [40] and are encoded by the *plcABC* operon, contained multiple *IS6110*

391 insertions detected by ISMapper. This locus is a known hotspot for IS6110 insertions
392 and has been shown to mediate deletions of segments of this region [41]. IS6110
393 insertions upstream of *phoP*, which have been associated with upregulation and
394 enhanced virulence in *M. tuberculosis* [9], were identified in multiple lineages (1
395 insertion in 6 lineage 2 genomes; singular insertions in one genome each in lineage 3
396 and 5) and may be indicative of positive selection for enhanced *phoP* expression and
397 virulence. These findings from ISMapper analysis are consistent with those reported
398 from PCR-based screens of smaller sets of isolates, but provide a more
399 comprehensive picture of IS dynamics in *M. tuberculosis* that could be extended to
400 much larger genomic data sets and other important pathogens.

401 **Conclusions**

402 ISMapper is a lightweight and reliable tool for the detection of IS insertion sites in
403 bacterial genomes using high-throughput short-read sequencing data, which is now
404 ubiquitous in microbial research and clinical investigations. ISMapper performed well
405 on real and simulated data from 32 different IS and 13 bacterial species, detecting all
406 but the most complex instances involving multiple neighbouring IS insertions or other
407 repeated sequences. ISMapper was able to detect antimicrobial resistance-associated
408 ISAba1 insertions in *A. baumannii*, with all sites detected by the program being
409 subsequently confirmed by PCR. Compared to other tools such as *breseq* and TIF,
410 ISMapper is ideal for detecting new positions for known IS in bacterial genomes. In
411 addition, ISMapper was able to rapidly produce a wealth of data on IS6110 insertions
412 in *M. tuberculosis*, allowing quick identification of lineage-specific insertions and
413 specific regions enriched for insertions that may be functionally significant.

414 **Availability and Requirements**

- 415 • **Project name:** ISMapper
- 416 • **Project home page:** https://github.com/jhawkey/IS_mapper
- 417 • **Programming language:** Python v2.7.5
- 418 • **Operating system(s):** platform independent, requires Python 2.7 and
419 dependencies
- 420 • **Other requirements:** BioPython v1.63, BWA v0.7.12, SAMtools v1.1,
421 Bedtools v2.20.1, BLAST+ v2.2.28, Samblaster v0.1.21
- 422 • **License:** Modified BSD

423 **Competing Interests**

424 The authors declare no competing interests.

425 **Authors' Contributions**

426 JH developed the code, analysed data and wrote the paper. KEH conceived the study
427 and helped to draft the manuscript. HBJ participated in design and coordination of the
428 study and contributed to data interpretation. RRW and DJE developed code. MH
429 performed PCR and sequence analysis. RMH provided sequence data and isolates for
430 validation and contributed to data interpretation. All authors read and approved the
431 final manuscript.

432 Acknowledgements

433 This work was supported by the National Health and Medical Research Council of
434 Australia (Fellowship #1061409 to KEH; Project Grant #1043830 to KEH and RMH)
435 and the Victorian Life Sciences Computation Initiative (VLSCI, #VR0082).

436 References

- 437 1. Siguier P, Perochon J, Lestrade L, Mahillon J, Chandler M: **ISfinder: the**
438 **reference centre for bacterial insertion sequences.** *Nucleic Acids Res* 2006,
439 **34**(Database issue):D32–6.
- 440 2. Mahillon J, Chandler M: **Insertion Sequences.** *Microbiol Mol Biol Rev* 1998,
441 **62**:725–774.
- 442 3. Siguier P, Goureyre E, Chandler M: **Bacterial insertion sequences: their**
443 **genomic impact and diversity.** *FEMS Microbiol Rev* 2014, **38**:865–891.
- 444 4. Olliver A, Vallé M, Chaslus-Dancla E, Cloeckeaert A: **Overexpression of the**
445 **multidrug efflux operon *acrEF* by insertional activation with *IS1* or *IS10***
446 **elements in *Salmonella enterica* serovar *typhimurium* DT204 *acrB* mutants**
447 **selected with fluoroquinolones.** *Antimicrob Agents Chemother* 2005, **49**:289–301.
- 448 5. Jellen-Ritter AS, Kern WV: **Enhanced expression of the multidrug efflux pumps**
449 **AcrAB and AcrEF associated with insertion element transposition in *Escherichia***
450 ***coli* mutants selected with a fluoroquinolone.** *Antimicrob Agents Chemother* 2001,
451 **45**:1467–1472.
- 452 6. Hamidian M, Hall RM: **ISAbal targets a specific position upstream of the**
453 **intrinsic *ampC* gene of *Acinetobacter baumannii* leading to cephalosporin**
454 **resistance.** *J Antimicrob Chemother* 2013, **68**:2682–2683.
- 455 7. Hamidian M, Hancock DP, Hall RM: **Horizontal transfer of an ISAbal25-**
456 **activated *ampC* gene between *Acinetobacter baumannii* strains leading to**
457 **cephalosporin resistance.** *J Antimicrob Chemother* 2013, **68**:244–245.
- 458 8. Hamidian M, Hall RM: **Tn6168, a transposon carrying an ISAbal-activated**
459 ***ampC* gene and conferring cephalosporin resistance in *Acinetobacter baumannii*.**
460 *J Antimicrob Chemother* 2014, **69**:77–80.
- 461 9. Soto CY, Menéndez MC, Pérez E, Samper S, Gómez AB, García MJ, Martín C:
462 **IS6110 mediates increased transcription of the *phoP* virulence gene in a**
463 **multidrug-resistant clinical isolate responsible for tuberculosis outbreaks.** *J Clin*
464 *Microbiol* 2004, **42**:212–219.
- 465 10. Uria MJ, Zhang Q, Li Y, Chan A, Exley RM, Gollan B, Chan H, Feavers I,

- 466 Yarwood A, Abad R, Borrow R, Fleck RA, Mulloy B, Vazquez JA, Tang CM: A
467 **generic mechanism in *Neisseria meningitidis* for enhanced resistance against**
468 **bactericidal antibodies.** *J Exp Med* 2008, **205**:1423–1434.
- 469 11. Van Der Ploeg J, Willemsen M, Van Hall G, Janssen DB: **Adaptation of**
470 ***Xanthobacter autotrophicus* GJ10 to bromoacetate due to activation and**
471 **mobilization of the haloacetate dehalogenase gene by insertion element IS1247.** *J*
472 *Bacteriol* 1995, **177**:1348–1356.
- 473 12. Aronson BD, Levinthal M, Somerville RL: **Activation of a cryptic pathway for**
474 **threonine metabolism via specific IS3-mediated alteration of promoter structure**
475 **in *Escherichia coli*.** *J Bacteriol* 1989, **171**:5503–5511.
- 476 13. Soria G, Barbé J, Gibert I: **Molecular fingerprinting of *Salmonella***
477 ***typhimurium* by IS200-typing as a tool for epidemiological and evolutionary**
478 **studies.** *Microbiologia* 1994, **10**:57–68.
- 479 14. Das S, Paramasivan CN, Lowrie DB, Prabhakar R, Narayanan PR: **IS6110**
480 **restriction fragment length polymorphism typing of clinical isolates of**
481 ***Mycobacterium tuberculosis* from patients with pulmonary tuberculosis in**
482 **Madras, South India.** *Tuber Lung Dis* 1995, **76**:550–554.
- 483 15. Bik EM, Gouw RD, Mooi FR: **DNA fingerprinting of *Vibrio cholerae* strains**
484 **with a novel insertion sequence element: a tool to identify epidemic strains.** *J Clin*
485 *Microbiol* 1996, **34**:1453–1461.
- 486 16. Adams MD, Chan ER, Molyneaux ND, Bonomo RA: **Genomewide Analysis of**
487 **Divergence of Antibiotic Resistance Determinants in Closely Related Isolates of**
488 ***Acinetobacter baumannii*.** *Antimicrob Agents Chemother* 2010, **54**:3569–3577.
- 489 17. Suzuki M, Matsumoto M, Hata M, Takahashi M, Sakae K: **Development of a**
490 **rapid PCR method using the insertion sequence IS1203 for genotyping Shiga**
491 **toxin-producing *Escherichia coli* O157.** *J Clin Microbiol* 2004, **42**:5462–5466.
- 492 18. Doig KD, Holt KE, Fyfe JA, Lavender CJ, Eddyani M, Portaels F, Yeboah-Manu
493 D, Pluschke G, Seemann T, Stinear TP: **On the origin of *Mycobacterium ulcerans*,**
494 **the causative agent of Buruli ulcer.** *BMC Genomics* 2012, **13**:258.
- 495 19. Doughty EL, Sergeant MJ, Adetifa I, Antonio M, Pallen MJ: **Culture-**
496 **independent detection and characterisation of *Mycobacterium tuberculosis* and**
497 ***M. africanum* in sputum samples using shotgun metagenomics on a benchtop**
498 **sequencer.** *PeerJ* 2014, **2**:e585.
- 499 20. Rizk G, Gouin A, Chikhi R, Lemaitre C: **MindTheGap: integrated detection**
500 **and assembly of short and long insertions.** *Bioinformatics* 2014, **30**:3451–3457.
- 501 21. Chen K, Wallis JW, McLellan MD, Larson DE, Kalicki JM, Pohl CS, McGrath
502 SD, Wendl MC, Zhang Q, Locke DP, Shi X, Fulton RS, Ley TJ, Wilson RK, Ding L,
503 Mardis ER: **BreakDancer: an algorithm for high-resolution mapping of genomic**
504 **structural variation.** *Nat Meth* 2009, **6**:677–681.
- 505 22. Robb SMC, Lu L, Valencia E, Burnette JM, Okumoto Y, Wessler SR, Stajich JE:

- 506 **The use of RelocaTE and unassembled short reads to produce high-resolution**
507 **snapshots of transposable element generated diversity in rice.** *G3* 2013, **3**:949–
508 957.
- 509 23. Keane TM, Wong K, Adams DJ: **RetroSeq: transposable element discovery**
510 **from next-generation sequencing data.** *Bioinformatics* 2013, **29**:389–390.
- 511 24. Nakagome M, Solovieva E, Takahashi A, Yasue H, Hirochika H, Miyao A:
512 **Transposon Insertion Finder (TIF): a novel program for detection of de novo**
513 **transpositions of transposable elements.** *BMC Bioinformatics* 2014, **15**:71.
- 514 25. Barrick JE, Colburn G, Deatherage DE, Traverse CC, Strand MD, Borges JJ,
515 Knoester DB, Reba A, Meyer AG: **Identifying structural variation in haploid**
516 **microbial genomes from short-read resequencing data using breseq.** *BMC*
517 *Genomics* 2014, **15**:1039.
- 518 26. Li H, Durbin R: **Fast and accurate short read alignment with Burrows-**
519 **Wheeler transform.** *Bioinformatics* 2009, **25**:1754–1760.
- 520 27. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis
521 G, Durbin R, 1000 Genome Project Data Processing Subgroup: **The Sequence**
522 **Alignment/Map format and SAMtools.** *Bioinformatics* 2009, **25**:2078–2079.
- 523 28. Quinlan AR, Hall IM: **BEDTools: a flexible suite of utilities for comparing**
524 **genomic features.** *Bioinformatics* 2010, **26**:841–842.
- 525 29. Faust GG, Hall IM: **SAMBLASTER: fast duplicate marking and structural**
526 **variant read extraction.** *Bioinformatics* 2014, **30**:2503–2505.
- 527 30. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden
528 TL: **BLAST+: architecture and applications.** *BMC Bioinformatics* 2009:421.
- 529 31. Holt KE, Parkhill J, Mazzoni CJ, Roumagnac P, Weill F-X, Goodhead I, Rance R,
530 Baker S, Maskell DJ, Wain J, Dolecek C, Achtman M, Dougan G: **High-throughput**
531 **sequencing provides insights into genome variation and evolution in *Salmonella***
532 **Typhi.** *Nat Genet* 2008, **40**:987–993.
- 533 32. Inouye M, Dashnow H, Raven L-A, Schultz MB, Pope BJ, Tomita T, Zobel J,
534 Holt KE: **SRST2: Rapid genomic surveillance for public health and hospital**
535 **microbiology labs.** *Genome Med* 2014, **6**:90.
- 536 33. **ARG-ANNOT - Antibiotic Resistance Gene-ANNOTation**
537 [<http://en.mediterranee-infection.com/article.php?laref=283&titre=arg-annot>]
- 538 34. McEvoy CRE, Falmer AA, van Pittius NCG, Victor TC, van Helden PD, Warren
539 RM: **The role of IS6110 in the evolution of *Mycobacterium tuberculosis*.**
540 *Tuberculosis* 2007, **87**:393–404.
- 541 35. van Embden JD, Cave MD, Crawford JT, Dale JW, Eisenach KD, Gicquel B,
542 Hermans P, Martín C, McAdam R, Shinnick TM: **Strain identification of**
543 ***Mycobacterium tuberculosis* by DNA fingerprinting: recommendations for a**
544 **standardized methodology.** *J Clin Microbiol* 1993, **31**:406–409.

- 545 36. Beggs ML, Eisenach KD, Cave MD: **Mapping of IS6110 insertion sites in two**
546 **epidemic strains of *Mycobacterium tuberculosis***. *J Clin Microbiol* 2000, **38**:2923–
547 2928.
- 548 37. Alonso H, Aguilo JI, Samper S, Caminero JA, Campos-Herrero MI, Gicquel B,
549 Brosch R, Martín C, Otal I: **Deciphering the role of IS6110 in a highly**
550 **transmissible *Mycobacterium tuberculosis* Beijing strain, GC1237**. *Tuberculosis*
551 2011, **91**:117–126.
- 552 38. Comas I, Coscolla M, Luo T, Borrell S, Holt KE, Kato-Maeda M, Parkhill J,
553 Malla B, Berg S, Thwaites G, Yeboah-Manu D, Bothamley G, Mei J, Wei L, Bentley
554 S, Harris SR, Niemann S, Diel R, Aseffa A, Gao Q, Young D, Gagneux S: **Out-of-**
555 **Africa migration and Neolithic coexpansion of *Mycobacterium tuberculosis* with**
556 **modern humans**. *Nat Genet* 2013, **45**:1176–1182.
- 557 39. Holt KE, Thieu Nga TV, Thanh DP, Vinh H, Kim DW, Vu Tra MP, Campbell JI,
558 Hoang NVM, Vinh NT, Minh PV, Thuy CT, Nga TTT, Thompson C, Dung TTN,
559 Nhu NTK, Vinh PV, Tuyet PTN, Phuc HL, Lien NTN, Phu BD, Ai NTT, Tien NM,
560 Dong N, Parry CM, Hien TT, Farrar JJ, Parkhill J, Dougan G, Thomson NR, Baker S:
561 **Tracking the establishment of local endemic populations of an emergent enteric**
562 **pathogen**. *Proc Natl Acad Sci U S A* 2013, **110**:17522–17527.
- 563 40. Assis PA, Espíndola MS, Paula-Silva FW, Rios WM, Pereira PA, Leão SC, Silva
564 CL, Faccioli LH: ***Mycobacterium tuberculosis* expressing phospholipase C**
565 **subverts PGE₂ synthesis and induces necrosis and alevolar macrophages**. *BMC*
566 *Microbiol* 2014, **14**:128.
- 567 41. Vera-Cabrera L, Hernández-Vera MA, Welsh O, Johnson WM, Castro-Garza J:
568 **Phospholipase region of *Mycobacterium tuberculosis* is a preferential locus for**
569 **IS6110 transposition**. *J Clin Microbiol* 2001, **39**:3499–3504.

570

571 **Figure Legends**

572 **Figure 1: Workflow for ISMapper.** (a) Inputs are reads (fastq format) and an IS
573 query (fasta), as well as either a reference genome to compare to or an assembly of
574 the reads (fasta or genbank). (b) Reads are mapped to the IS query using BWA
575 (dashed lines) and their pairs (i.e. flanking reads, solid lines) are retrieved from the
576 resulting SAM file using mapping flags (SAMTools). The unmapped component of
577 soft clipped reads (solid+dashed lines), identified from the SAM file using
578 Samblaster) are also retrieved using BioPython. (c) Flanking and softclipped read
579 sequences are then mapped to either the reference genome or the read assembly
580 (BWA), and the final mapping output is then filtered for depth to remove low
581 coverage regions. Left and right end blocks are extracted from the resulting BAM file
582 (using BedTools) and compared to one another to find either intersecting regions,
583 indicating a novel IS position, or close regions, indicating a known IS position in the
584 reference. In a novel position, overlapping sequences from the left and right ends
585 most likely indicate the target site duplication generated during IS transposition
586 (zoomed in section). (d) Results are tabulated, indicating the position and orientation
587 of each site, whether it is novel or known, and information about the genes flanking
588 the insertion site.

589 Figure 2: Issues can arise in the second mapping step of ISMapper if the IS insertion
590 is next to a region that occurs more than once in the reference genome. In this
591 example, the left flanking reads (green arrows) have mapped to all possible copies of
592 the repeat sequence (dark red boxes) and are each randomly assigned to a repeat copy
593 by BWA, resulting in each copy falling below the read depth cut-off (default 6x). To
594 overcome this, where an unpaired flanking region is identified (purple box),
595 ISMapper searches for a potential low-depth partner flanking region that is close to or

596 intersecting the confident region (green box). This is recorded as an uncertain call,
597 labelled in the output with the ‘?’ character.

598 **Figure 3: ISAba1 positions detected in seven *Acinetobacter baumannii* genomes.**

599 The A1 reference genome is indicated by the outer black circle, tick marks indicate
600 genome coordinates (in Mbp), genes targeted in the two available MLST schemes are
601 marked in blue for orientation purposes. ISAba1 detected within the individual test
602 genomes are shown on inner rings (purple), orientation of the IS is indicated by
603 shading (dark, left end; light, right end). Novel IS insertions were identified upstream
604 of the beta-lactamase *ampC* (green) in all isolates, explaining the observed phenotypic
605 resistance to third generation cephalosporins.

606 **Figure 4: ISAba1 detection rate as a function of read depth, for seven *A.***

607 ***baumannii* GC1 genomes sequenced with Illumina.** Mean (lines) and standard
608 deviation (shaded areas) proportions of IS insertions correctly detected per genome,
609 amongst 70 replicate read sets sampled at each depth level (7 read sets x 10 replicates
610 each at read depths 10x, 15x, 20x, 25x, 50x, 75x, 100x). Blue, IS insertion site
611 detected; red, detected with high confidence; purple, detected with low precision
612 (larger than expected gap size between left and right flanking regions, indicated with
613 ‘*’ in output); green, detected with low confidence (high-depth evidence for one side
614 only, low-depth evidence for the other, indicated with ‘?’ in output).

615 **Figure 5: Analysis of IS6110 insertions in a diverse set of *M. tuberculosis***

616 **genomes.** Analyses are based on ISMapper analysis of publicly available Illumina
617 paired end data from 138 genomes. (a) Phylogenetic tree for *M. tuberculosis* based on
618 genome-wide SNP calls, with midpoint rooting and collapsed to lineage level.
619 Number of isolates analysed in each lineage (L) is indicated in brackets. (b) Heat map

620 indicating the frequency of each *IS6110* insertion site (columns) detected in each
621 lineage (rows), according to inset legend (i.e. a black cell indicates that the IS
622 insertion site was detected in all isolates of the given lineage, white cell indicates that
623 the insertion site was not detected at all in that lineage). **(c)** Boxplots show number of
624 *IS6110* insertions detected per genome, for each lineage. Black line, median; boxes,
625 interquartile range; whiskers, minimum and maximum values. **(d)** Histogram of
626 *IS6110* insertions in 1,000 bp windows along the *M. tuberculosis* chromosome.
627 Dashed line, threshold for defining insertion hotspots.

628

629 **Tables**

630 **Table 1:** Validation of ISMapper using simulated reads. * indicates ISMapper was

631 unable to resolve some IS positions due to repeat regions.

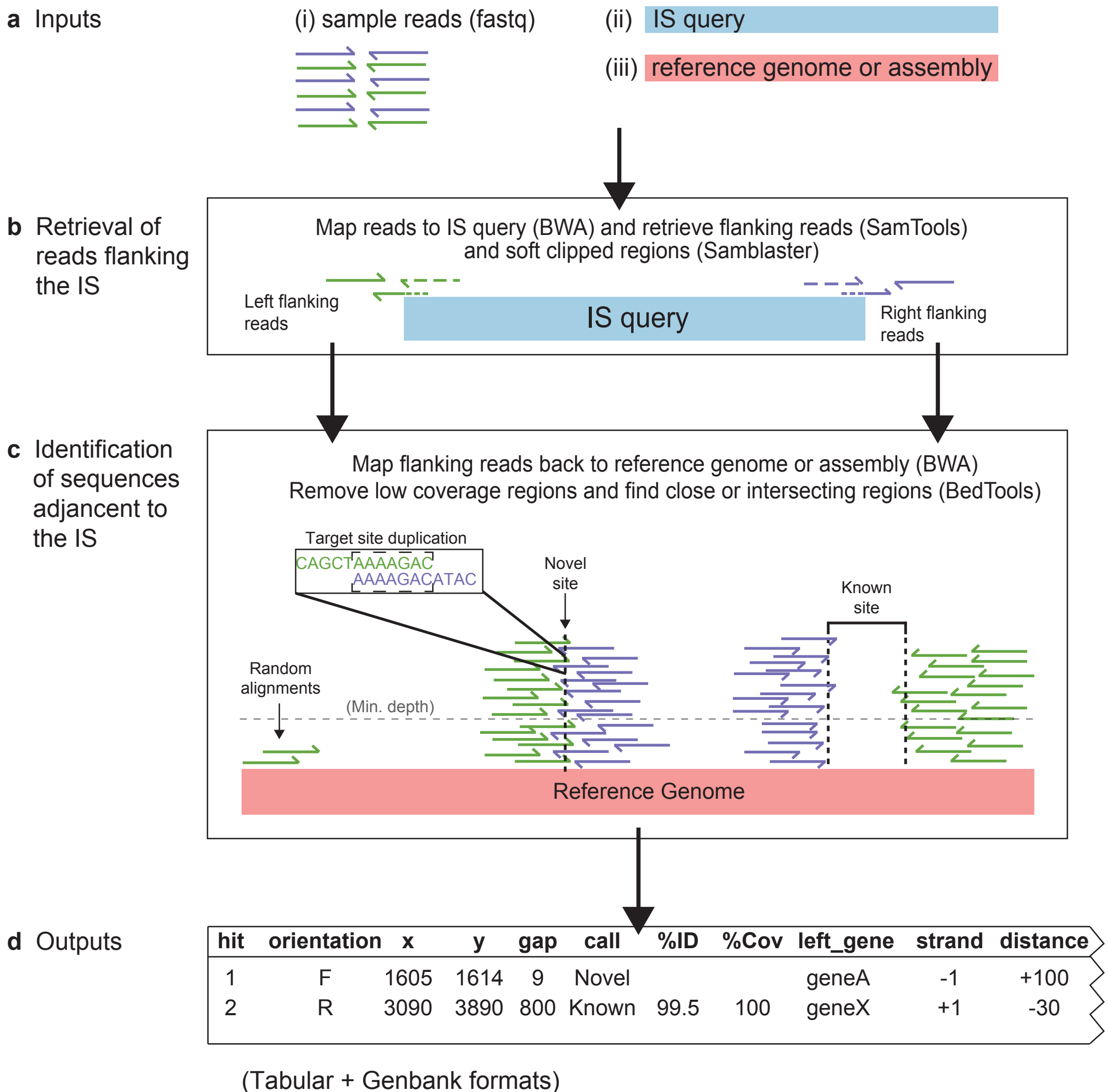
Isolate	Accession	IS	Found	Orientation
<i>S. Typhi</i> CT18	NC_003198	IS200	26/26	26/26
		IS1	3/3	3/3
<i>S. Typhimurium</i> LT2	NC_003197	IS200	6/6	6/6
		ISSty2	2/2	2/2
		IS1351	2/2	2/2
<i>S. Typhi</i> plasmid pHCM1	NC_003384	IS26	4/4	4/4
		ISVsa5	2/2	2/2
<i>S. Paratyphi</i> plasmid pAKU_1	AM412236	IS1	5/5	5/5
		IS26	6/6	6/6
		ISVsa5	2/2	2/2
<i>K. pneumoniae</i> plasmid pNDMAR	JN420336	IS1	7/7	7/7
		IS3000	3/3	3/3
		IS26*	3/5	3/5
<i>Yersinia pestis</i> C092	NC_003143	ISEcp1	2/2	2/2
		IS100*	43/44	43/44
		IS1661*	7/8	7/8
<i>Escherichia coli</i> O104:H4	NC_018658	IS1541*	61/64	61/64
		IS1	10/10	10/10
		IS421	4/4	4/4
		IS609	4/4	4/4
		ISEc1	4/4	4/4
<i>E. coli</i> O157:H7	NC_002695	ISKpn26	4/4	4/4
		IS629*	17/18	17/18
		IS609	2/2	2/2
		ISEc1	4/4	4/4
		ISEc8*	9/10	9/10

632

633 **Table 2:** Validation of ISMapper using real Illumina reads for which finished
 634 genomes were also available. * indicates ISMapper was unable to resolve some
 635 positions due to repeat regions.

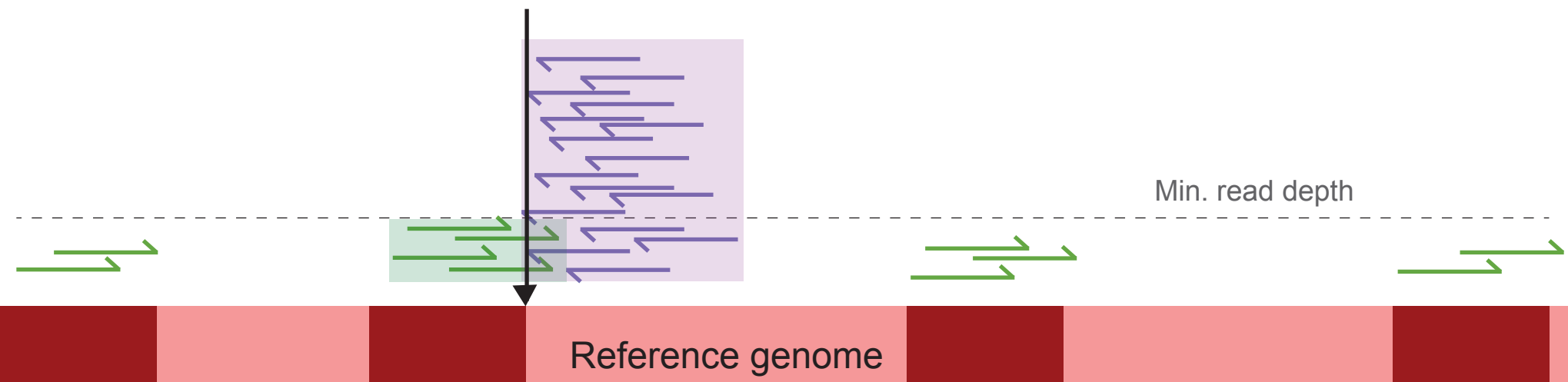
Isolate	Genome	FastQ	IS	Found	Orientation
<i>Streptococcus suis</i> P1/7	AM946016	ERR225612	ISSu3	4/4	4/4
			ISSu4	2/2	2/2
<i>Staphylococcus aureus</i> TW20	NC_017331	ERR043367	ISSep3	3/3	3/3
			IS256	8/8	8/8
			IS431mec*	3/4	3/4
			IS1181	2/2	2/2
<i>Klebsiella pneumoniae</i> NJST258_1	CP006923	SRR1166975	ISKpn1	6/6	6/6
			IS5B	8/8	8/8
			IS903B	2/2	2/2
			IS1294	3/3	3/3
			ISKpn18	2/2	2/2
			ISKpn26	7/7	7/7
<i>S. Typhi</i> CT18	NC_003198	ERR343331	IS200	26/26	26/26
			IS1*	2/3	2/3
<i>S. Typhi</i> Ty2	AE014613	ERR343332	IS200	26/26	26/26

636



- ← left flanking read
- ← right flanking read
- repeat copies

IS insertion site



■ ISAba1
■ MLST genes

