

MOBFinder: A tool for mobilization typing for plasmid metagenomic fragments based on a language model

--Manuscript Draft--

Manuscript Number:	GIGA-D-24-00070R1	
Full Title:	MOBFinder: A tool for mobilization typing for plasmid metagenomic fragments based on a language model	
Article Type:	Technical Note	
Funding Information:	National Key Research and Development Program of China (2022YFA0806400)	Prof. Hongwei Zhou
	National Natural Science Foundation of China (82102508)	Dr. Zhencheng Fang
	National Natural Science Foundation of China (81925026)	Prof. Hongwei Zhou
Abstract:	<p>Background Mobilization typing (MOB) is a classification scheme for plasmid genomes based on their relaxase gene. The host ranges of plasmids of different MOB categories are diverse and MOB is crucial for investigating plasmid mobilization, especially the transmission of resistance genes and virulence factors. However, MOB typing of plasmid metagenomic data is challenging due to the highly fragmented characteristics of metagenomic contigs.</p> <p>Results We developed MOBFinder, an 11-class classifier, for categorizing plasmid fragments into 10 MOB types and a non-mobilizable category. We first performed MOB typing to classify complete plasmid genomes according to relaxase information and then constructed an artificial benchmark dataset of plasmid metagenomic fragments (PMFs) from those complete plasmid genomes whose MOB types are well annotated. Next, based on natural language models, we used word vectors to characterize the PMFs. Several random forest classification models were trained and integrated to predict fragments of different lengths. Evaluating the tool using the benchmark dataset, we found that MOBFinder outperforms previous tools such as MOBscan and MOB-suite, with an overall accuracy approximately 59% higher than that of MOB-suite. Moreover, the balanced accuracy, harmonic mean, and F1-score reached up to 99% for some MOB types. When applied to a cohort of patients with type II diabetes (T2D), MOBFinder offered insights suggesting that the MOB type plasmid, which is widely present in Escherichia and Klebsiella, and the MOBQ type plasmid, might accelerate antibiotic resistance transmission in patients suffering from T2D.</p> <p>Conclusions To the best of our knowledge, MOBFinder is the first tool for MOB typing of PMFs. The tool is freely available at https://github.com/FengTaoSMU/MOBFinder.</p>	
Corresponding Author:	Zhencheng Fang Zhujiang Hospital of Southern Medical University Guangzhou, CHINA	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:	Zhujiang Hospital of Southern Medical University	
Corresponding Author's Secondary Institution:		
First Author:	Tao Feng	
First Author Secondary Information:		
Order of Authors:	Tao Feng Shufang Wu	

	Hongwei Zhou
	Zhencheng Fang
Order of Authors Secondary Information:	
Response to Reviewers:	The response to specific reviewer and editor comments has been uploaded in the cover letter.
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
Experimental design and statistics	Yes
<p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	
Resources	Yes
<p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	
Availability of data and materials	Yes
<p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories</p>	

(where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.

Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist?](#)

46

47

48 **Abstract**

49 **Background**

50 Mobilization typing (MOB) is a classification scheme for plasmid genomes based on their relaxase
51 gene. The host ranges of plasmids of different MOB categories are diverse and MOB is crucial for
52 investigating plasmid mobilization, especially the transmission of resistance genes and virulence
53 factors. However, MOB typing of plasmid metagenomic data is challenging due to the highly
54 fragmented characteristics of metagenomic contigs.

55 **Results**

56 We developed MOBFinder, an 11-class classifier, for categorizing plasmid fragments into 10 MOB
57 types and a non-mobilizable category. We first performed MOB typing to classify complete plasmid
58 genomes according to relaxase information and then constructed an artificial benchmark dataset of
59 plasmid metagenomic fragments (PMFs) from those complete plasmid genomes whose MOB types
60 are well annotated. Next, based on natural language models, we used word vectors to characterize
61 the PMFs. Several random forest classification models were trained and integrated to predict
62 fragments of different lengths. Evaluating the tool using the benchmark dataset, we found that
63 MOBFinder outperforms previous tools such as MOBscan and MOB-suite, with an overall *accuracy*
64 approximately 59% higher than that of MOB-suite. Moreover, the *balanced accuracy*, *harmonic*
65 *mean*, and *F1-score* reached up to 99% for some MOB types. When applied to a cohort of patients
66 with type II diabetes (T2D), MOBFinder offered insights suggesting that the MOBF type plasmid,
67 which is widely present in *Escherichia* and *Klebsiella*, and the MOBQ type plasmid, might
68 accelerate antibiotic resistance transmission in patients suffering from T2D.

69 **Conclusions**

70 To the best of our knowledge, MOBFinder is the first tool for MOB typing of PMFs. The tool is
71 freely available at <https://github.com/FengTaoSMU/MOBFinder>.

72

73 **Keywords:** MOB typing; language model; metagenomic sequencing; plasmid; random forest

74

75 **1. Introduction**

76 Plasmids are usually small, double-stranded, and circular DNA molecules found within bacterial
77 cells [1]. Being separate from the bacterial chromosome, plasmids have the ability to replicate
78 independently and can be transferred between bacteria through conjugation [2]. Bacteria,
79 specifically pathogenic strains, can acquire antibiotic resistance genes or virulence factors via
80 plasmid-mediated horizontal gene transfer, aiding their ability to adapt to various environments [3].
81

82 Plasmid classification is important for investigating multiple properties of plasmids, such as host
83 range, replication patterns, and mobilization mechanisms [4]. Many classification schemes have
84 been developed according to the distinct characteristics of plasmids, including taxonomic
85 classification, replicon typing (Rep), incompatibility typing (Inc), mate-pair formation typing
86 (MPF), and mobilization typing (MOB). In taxonomic classification, plasmids are categorized based
87 on their host bacteria [5]. Rep typing classifies plasmids according to genes controlling their
88 replication, known as replication initiation genes [4, 6]. Inc typing takes advantage of the fact that
89 plasmids with similar replication or partition systems are incompatible within the same cell,
90 categorizing plasmids based on compatibility [6]. MPF typing is based on genes encoding the MPF
91 system, which consists of proteins that mediate contact and DNA exchange between donor and
92 recipient cells during conjugation [4, 7]. Finally, MOB typing classifies plasmids based on the
93 relaxase gene, which is present in all transmissible plasmids [8-10]. And plasmids with different
94 relaxase types are categorized as different MOB types, each of possesses a distinct transmission
95 mechanism that determines its taxonomic host range [4, 11]. This variation among different MOB
96 types is critical in researching the spread of virulence traits, the emergence of antibiotic resistance,
97 and the adaptation and evolution of bacteria. Moreover, MOB typing has been found to be effective
98 for identifying novel mobilizable plasmids that were previously unassigned to any Rep or Inc types,
99 and for investigating the mobilization characteristics of plasmids with similar mobilization systems
100 [12, 13].

101

102 Recently, many experimental and computational schemes have been devised for plasmid typing, as
103 well as to explore the diversity and functionality of plasmids (Table 1). For example, plasmid
104 taxonomic PCR (PlasTax-PCR) [14], PCR-based replicon typing (PBRT) [15], and degenerate
105 primer MOB typing (DPMT) [12] are multiplex PCR methods for identifying plasmids with

106 analogous replication or mobilization systems. PlasTrans, based on deep learning, identifies
107 mobilizable metagenomic plasmid fragments [16]. Web servers such as PlasmidFinder [6], pMLST,
108 and oriTfinder [17] were established based on collected maker gene databases and alignment-based
109 methods to facilitate Rep, Inc, and MOB typing. COPLA [5], based on average nucleotide identity,
110 performs taxonomic classifications of complete plasmid genomes with an overall accuracy of 41%.
111 For the MOB typing, MOBscan [18] uses the HMMER model to annotate relaxase genes and
112 classify plasmids accordingly. MOB-suite [19, 20] performs plasmid typing for plasmid assemblies.
113 First, it uses Mash distance to cluster plasmid assemblies into clusters; then, it uses marker gene
114 databases to annotate them.

115

116 **Table 1.** Experimental and computational schemes developed for plasmid classification.

117

118 Metagenomic sequencing makes it possible to obtain all plasmid DNA from microbial communities
119 at once, and a number of computational tools for identifying plasmid fragments from metagenomic
120 data have been developed, such as PlasFlow [21], PlasmidSeeker [22], PlasClass [23], PPR-Meta
121 [24] and PlasForest [25]. As DNA fragments of plasmids and bacteria are intermingled in
122 metagenomic data [26], recognizing the transmission mechanisms and host ranges of plasmids can
123 be challenging. To this end, it is crucial to annotate MOB types of metagenomic plasmid fragments.
124 However, this is difficult when plasmid assembly fragments are incomplete and essential genes for
125 annotation are lacking. Therefore, it is worthwhile to consider alternative methods. Given that
126 plasmids of the same MOB type have similar transmission mechanisms and host ranges, their
127 genomic signatures (e.g., GC content and codon usage) tend to also be alike, not only relaxase [4,
128 27]. In this context, neural networks, which have demonstrated strong performance in the
129 classification and identification of biological sequences [28, 29] could be useful. Furthermore,
130 language models [30, 31] derived from such neural networks have also showcased their impressive
131 ability to characterize sequence features [32, 33]. In this methodology, short sequences of
132 nucleotides (referred to as k -mers) or amino acids are analogous to “words”, and the longer
133 sequences of DNA or proteins are analogous to “sentences”. Through the application of
134 unsupervised learning on large datasets, each “word” is linked to a feature vector that captures its
135 context, offering a more sophisticated analysis than the traditional k -mer frequency method, which

136 simply counts the occurrence of nucleotide sequences without acknowledging their biochemical
137 characteristics. Unlike the conventional method, this language model-based approach assesses
138 sequences based on their contextual importance across different genetic environments, positioning
139 contextually similar sequences close together in a multidimensional space. This technique provides
140 deeper insights into the biochemical complexities of nucleotide sequences, thereby furnishing a
141 more comprehensive understanding of an organism's functional biology [34]. To characterize the
142 features of plasmids within the same MOB type, we employed language models to perform the
143 MOB annotation. In addition to the relaxase-coding gene, language models exhibit the ability to
144 capture more biological features and associations within comparable mobilization systems, making
145 it possible to perform MOB annotation for metagenomic plasmid assemblies.

146

147 Thus, we presented MOBFinder, a tool for annotating MOB types from plasmid metagenomic
148 fragments (PMFs). MOBFinder can process single or multiple plasmid DNA sequences, and
149 provides predicted MOB types for each input fragment, including MOBB, MOBC, MOBF, MOBH,
150 MOBL, MOBM, MOBP, MOBQ, MOBT, MOBV and non-MOB. Moreover, it provides the option
151 to annotate plasmid bins from metagenomics data.

152

153 An overview of this work is shown in Figure 1A, and the development of MOBFinder involved the
154 following steps: (1) Benchmark dataset construction. Plasmid complete genomes obtained from the
155 National Center for Biotechnology Information (NCBI) were classified into different MOB types
156 based on relaxase databases. Then, to simulate plasmid fragments in metagenomic data, an artificial
157 benchmark dataset of varying lengths is generated. (2) Word embeddings. Numerical word vectors
158 were generated using skip-gram to characterize the sequence features of different MOB categories.
159 (3) Classification model ensemble and optimization. Several classification models, specifically
160 designed for different lengths, were trained and integrated to predict fragments of different lengths.
161 Evaluations against a test dataset demonstrated that MOBFinder is a powerful tool for MOB typing
162 of plasmid fragments and bins. Its application to a cohort of patients with type II diabetes (T2D)
163 revealed a potential correlation between some MOB types and the spread of antibiotic resistance
164 genes among T2D patients. This suggests that MOBFinder is an effective data analysis approach for
165 investigating plasmid-mediated horizontal gene transfer within microbial communities.

166

167 **2. Materials and methods**

168 **2.1. The workflow of MOBFinder**

169 To annotate the MOB type of plasmid fragments in metagenomics, we designed MOBFinder (Figure
170 1). As MOB-suite [19, 20] didn't offer a quantitative likelihood score for the outcomes and some
171 plasmids would be classified into multiple MOB types (Figure S1), we constructed a benchmark
172 dataset using a high-resolution MOB typing strategy for categorizing complete plasmid genomes
173 (Figure 1B, 1C). Then, based on a language model and random forest, we designed an algorithm to
174 perform MOB typing for PMFs (Figure 1D, 1E).

175

176 **Figure 1. Flowchart of the technical approach utilized in this study.** (A) General workflow of
177 the development and testing of MOBFinder. (B) Using plasmid relaxases with known MOB types
178 as reference sequences, we developed a database of relaxases from the non-redundant (NR) database
179 representing different MOB types. (C) Utilizing the relaxase database, complete plasmid genomes
180 from the NCBI were subjected to MOB typing. (D) Those complete genomes were also used to train
181 a 4-mer language model using the skip-gram algorithm, allowing each 4-mer to be represented by a
182 100-dimensional word vector. For a DNA fragment, the average word vector of all 4-mers on its
183 sequence serves as the feature vector for that DNA. (E) We constructed simulated metagenomic
184 contigs from the complete genomes that had been MOB typed as a benchmark and encoded these
185 contigs into word vectors. Then these word vectors were used to train a random forest algorithm.
186 Then the trained model, with metagenomic DNA fragments as input, was used to predict the MOB
187 typing of the corresponding DNA fragment based on its word vectors.

188

189 **2.2. MOB typing of complete plasmid genomes**

190 Traditionally, plasmid MOB typing of complete plasmid genomes has been a bioinformatics task
191 based on the analysis of relaxase sequence similarity. The practice of annotating MOB types through
192 BLAST similarity searches using representative sequences of different MOB type relaxases has
193 gradually evolved into the standard method for MOB typing [4, 19, 20]. In this work, we constructed
194 a benchmark dataset of simulated metagenomic contigs based on complete plasmid genomes with
195 known MOB types. Previous studies have included a relatively small number of plasmids in their

196 analyses. To further expand the MOB typing training dataset, we annotated the newly collected
197 plasmid complete genome data for MOB typing according to relaxase information.

198

199 Ten validated MOB relaxase protein families were collected, including MOBB, MOBC, MOBF,
200 MOBH, MOBL, MOBM, MOBP, MOBQ, MOBT and MOBV [7-10, 35, 36] (Figure 1B). For each
201 MOB category, blastp (RRID:SCR_001010) [37] was used to search homologous protein sequences
202 against the NCBI non-redundant protein sequence database, with an *e-value* threshold of 1e-10, a
203 *query coverage* threshold of 70%, and an *identity* threshold of 70%. A previous study applied an *e-*
204 *value* threshold of 1e-5, and minimum requirements for *query coverage* and *identity* set at 50% [4].
205 However, employing these criteria, we observed that some relaxases were annotated as belonging
206 to multiple MOB types. To eliminate ambiguous annotations and construct a more reliable dataset
207 for the training of MOBfinder, we imposed the stricter criteria mentioned above. After the
208 expansion of protein sequences, local relaxase databases were built using the ‘makeblastdb’
209 command for MOB typing of plasmid genomes.

210

211 Plasmid genomes were retrieved from the NCBI nucleotide database using the keywords ‘complete’
212 and ‘plasmid,’ and incomplete fragments were removed manually for further analysis. The accession
213 list of these plasmids is provided in Supplementary Table 1. For each plasmid genome, coding
214 sequences were extracted from the genbank file, and blastp [37] was employed to search for the
215 best alignment of local relaxase databases. Here, we defined the *mob_score* to measure the
216 likelihood of homology:

$$217 \quad \text{mob_score} = \sqrt{0.01 * \text{qcov_max} * (1 - 1/\log_{10}(\text{bitscore_max}))}$$

218 where *qcov_max* and *bitscore_max* represent the *query coverage* and *bitscore* corresponding to the
219 match with the highest bit score, respectively. To identify plasmid genomes encoding known
220 relaxase families, we set a *mob_score* threshold of 0.5, which was established in conjunction with a
221 minimum *query coverage* of 50% and a minimum *bitscore* of 100. To further enhance the reliability
222 of our classification, we introduced an *e-value* cutoff, conservatively set at 1e-10, to complete the
223 plasmid genome classification (Figure 1C). In instances where plasmid genomes yielded no blast
224 results or exhibited an *e-value* exceeding 0.01, we categorized them as non-MOB.

225

226 **2.3. Word embeddings using a language model**

227 To characterize the features and patterns within each MOB category and use numerical word vectors
228 to represent them, we utilized a skip-gram language model [30, 31] to learn from plasmid genomes.
229 Using a sliding window, the model calculated the likelihood between segmented words and
230 outputted a probability distribution over the context words. The training steps were as follows
231 (Figure 1D):

232

233 (1) Word generation. Since DNA sequences are composed of different nucleotide characters, we
234 used a k -mer sliding window to generate overlapping input words. For example, with $k=4$,
235 ‘ATCGCTGA’ would be segmented into ‘ATCG,’ ‘TCGC,’ ‘CGCT,’ ‘GCTG,’ and ‘CTGA’. In this
236 step, unique words were generated.

237

238 (2) Word encoding initialization. Each word was initially assigned a random vector.

239

240 (3) Skip-gram model. We employed a standard skip-gram model as described in previous studies
241 [30, 31] to generate word vectors through the dna2vec module [31]. A two-layer neural network was
242 used to construct the skip-gram model. The initialized vectors were used as input, and the output
243 was a probability distribution over the input words. Layer 1 was a hidden layer to convert the
244 initialized vectors into a 100-dimensional word vector representation as predefined by Ng [31].
245 Layer 2 was used to compute and maximize the probability of the correct context words using the
246 negative sampling function, with the size of context words set to 20 (10 words for upstream and
247 downstream, respectively) as pre-set by Ng [31].

248

249 (4) Model training. For each input plasmid genome, we used an optimization algorithm to minimize
250 the loss function. Then, using the default settings, we used backpropagation to update the neural
251 network parameters (word vectors) for 10 epochs.

252

253 (5) Word vector extraction. After the training process, the word vectors in the hidden layer were
254 extracted to characterize the plasmid fragments.

255

256 **2.4. Benchmark dataset construction**

257 Because there are no real metagenomic data to serve as a benchmark, using simulated data as a
258 benchmark dataset is a common approach when developing bioinformatics tools [16, 24]. Therefore,
259 in the development of MOBFinder, we artificially generated simulated datasets through the
260 following steps:

261

262 (1) For classified plasmid genomes in each MOB category, we randomly split them at a proportion
263 of 70% and 30% to construct the training and test datasets.

264

265 (2) Training dataset. To predict plasmid fragments with different lengths, we generated contigs of
266 different length ranges: 100-400 bp, 401-800 bp, 801-1200 bp, and 1201-1600 bp. For each MOB
267 class in each length range, we randomly generated 90000 artificial contigs. Plasmid fragments
268 longer than 1600 bp were segmented into shorter contigs and predicted using models designed for
269 the corresponding lengths.

270

271 (3) Test dataset. Because some plasmid fragments in real metagenomics datasets were much longer,
272 we generated four length groups to assess the performance of MOBFinder: Group A with a length
273 range of 801-1200 bp, Group B with a length range of 1201-1600 bp, Group C with a length range
274 of 3000-4000 bp, and Group D with a length range of 5000-10000 bp. For each MOB class in these
275 four groups, 500 fragments were randomly extracted.

276

277 **2.5. Classification algorithm**

278 To efficiently handle the training dataset and improve the robustness of MOBFinder, we employed
279 random forest to train four predictive models using the training dataset. The detailed steps are as
280 follows (Figure 1E):

281

282 (1) Word representation calculation. For each contig in the training dataset, we used a 4-mer sliding
283 window to generate overlapping words and transformed them into numerical word vectors using
284 trained word embeddings. To characterize the underlying features and patterns of the input contigs,
285 we summed all the word vectors to compute their average as input of random forest.

286

287 (2) Classification model training. To improve the performance of MOBFinder, we trained four
288 classification models on different lengths in the training dataset: 100-400 bp, 401-800 bp, 801-1200
289 bp, and 1201-1600 bp. The number of trees was set to 500 to generate predictive models.

290

291 (3) Model ensemble. The four trained models were ensembled into MOBFinder to make more
292 accurate predictions. For fragments shorter than 100 bp, we used a model designed for 100-400 bp
293 to predict the MOB type. For those longer than 1600 bp, we segmented them into short contigs and
294 made predictions using the corresponding model. For example, a fragment with a length of 4000 bp
295 would be segmented into three contigs: two with a length of 1600 bp and one of 800 bp. After
296 predicting fragments with the corresponding models, we aggregated and calculated the weighted
297 average scores for each MOB class, and the MOB type with the highest score was selected as the
298 final prediction result for the input fragment.

299

300 (4) Plasmid bin classification. Metagenomic binning is an essential step in the reconstruction of
301 genomes from individual microorganisms. Thus, we designed MOBFinder to perform MOB typing
302 on both plasmid contigs and plasmid bins. If the input is a plasmid bin, MOBFinder predicts the
303 likelihood of each MOB class for fragments within the bin. For each MOB category, MOBFinder
304 aggregates the scores of each sequence within the bin and calculates the weighted average scores
305 based on the sequence length. The MOB category with the maximum score is selected as the
306 prediction result.

307

308 **2.6. Performance validation**

309 A test dataset was used to assess the performance of MOBFinder and compare it to MOB-suite and
310 MOBscan. Because MOBscan can only predict MOB type using plasmid protein sequences rather
311 than DNA sequences, we first annotated the proteins in the plasmid fragments of the test set using
312 Prokka (RRID:SCR_014732) [38] and then used MOBscan to predict the MOB type based on the
313 annotated proteins. We calculated overall *accuracy*, *kappa*, and *run time* by comparing the predicted
314 classes and true classes. We used the online server of MOBscan to perform the MOB annotation,
315 and the calculation of *run time* for MOBScan was confined to the duration spent on preprocessing

316 with Prokka locally. The overall *accuracy* was the proportion of accurate predictions. The *kappa* (a)
 317 was calculated to assess the overall consistency between the predictions and true classes, which took
 318 into account the possibility of random prediction. P_o represented observed accuracy [$P_o = (A_{11} +$
 319 $A_{22} + \dots + A_{nn}) / N$], where A_{11} , A_{22} , and A_{nn} represented the values on the diagonal of the confusion
 320 matrix and n represented the number of MOB categories. N represented the total number of samples.
 321 P_e represented the expected accuracy [$P_e = (E_{11} + E_{22} + \dots + E_{nn}) / N^2$], where E_{11} , E_{22} , and E_{nn}
 322 were the expected values in each cell of the confusion matrix, n was the number of MOB classes,
 323 and N was the total number of samples. The *run time* was recorded using the command ‘time’ in
 324 Linux.

$$325 \quad \quad \quad \text{kappa} = (P_o - P_e) / (1 - P_e) \quad (a)$$

$$326 \quad \quad \quad \text{balanced accuracy} = (TPR + TNR) / 2 \quad (b)$$

$$327 \quad \quad \quad \text{harmonic mean} = 2 * S_n * S_p / (S_n + S_p) \quad (c)$$

$$328 \quad \quad \quad F1 - \text{score} = 2 * \text{precision} * \text{recall} / (\text{precision} + \text{recall}) \quad (d)$$

329 For each MOB category, we also calculated the *balanced accuracy* (b), *harmonic mean* (c) and *F1-*
 330 *score* (d). Considering the class imbalance within the training dataset, *balanced accuracy* was used
 331 to measure the average accuracy of each MOB category, where TPR was the true positive rate [TRP
 332 $= \text{true positives} / (\text{true positives} + \text{false negatives})$] and TNR was the true negative rate [$TNR = \text{true}$
 333 $\text{negatives} / (\text{true negatives} + \text{false positives})$]. The *harmonic mean* provided an overall evaluation
 334 of the model’s performance, where S_n and S_p represented sensitivity [$S_n = \text{true positives} / (\text{true}$
 335 $\text{positives} + \text{false negatives})$] and specificity [$S_p = \text{true negatives} / (\text{true negatives} + \text{false positives})$],
 336 respectively. The *F1-score* combined *precision* and *recall*, providing a balanced measure of the
 337 model’s performance, where *precision* was the number of correct positive predictions out of all
 338 positive predictions [$\text{precision} = \text{true positives} / (\text{true positives} + \text{false positives})$] and *recall* was the
 339 number of correct positive predictions out of all actual positive predictions. [$\text{recall} = \text{true positives}$
 340 $/ (\text{true positives} + \text{false negatives})$].

341

342 A receiver operating characteristic (ROC) curve was used to visualize the performance of
 343 MOBFinder in predicting each MOB category, where the x-axis and y-axis were the false positive
 344 rate (FPR) and true positive rate (TPR). Plots closer to the left and top indicate higher TPR and
 345 lower FPR , which means better performance. For each MOB class, the area under the curve (AUC)

346 value was calculated to quantify the performance of MOBFinder. An AUC value between 0.5 and 1
347 indicates that the model performs better than random chance, and a higher AUC value indicates
348 better prediction capability.

349

350 **2.7. Annotation and analysis of T2D metagenomic data**

351 Metagenomic sequencing data (SRA045646) were retrieved from the NCBI short read archive (SRA)
352 database to investigate whether the plasmids within different MOB classes were associated with
353 antibiotic resistance enrichment in T2D patients, as suggested by previous studies [39, 40]. All
354 metagenomic data were preprocessed using the same protocols. PRINSEQ (RRID:SCR_005454)
355 [41] was used to remove low-quality reads and bowtie2 (RRID:SCR_016368) [42] was used to
356 remove host reads by aligning them to the human GRCH38 reference genome downloaded from the
357 ENSEMBL database. We excluded metagenomic samples that did not pass quality control. Because
358 the abundance of plasmids in metagenomes was much lower than that of bacteria, we only retained
359 samples with more than 10,000,000 paired-end reads for downstream analysis (Supplementary
360 Table 2).

361

362 To improve the efficiency and accuracy of assembly, we used MEGAHIT (RRID:SCR_018551) [43]
363 to generate metagenomic contigs. PPR-Meta (RRID:SCR_016915) [24] was utilized to identify and
364 extract plasmid fragments from the assembled fragments while filtering out bacteria and phage
365 sequences. COCACOLA [44] was employed to cluster plasmid fragments into bins based on
366 sequence similarity and composition. This allowed us to investigate the plasmid fragments from
367 same originate and enabled better annotation and analysis of their functions.

368

369 MOBFinder was applied to annotate the MOB types in each plasmid bin. The average fragments
370 per kilobase per million of each plasmid bin was calculated using bowtie2 to represent its abundance.
371 Next, we analyzed the significance of differences in plasmid bins and various MOB types between
372 healthy and T2D groups using the Wilcoxon rank-sum test. The calculation of p values was adjusted
373 for multiple comparisons using the Benjamini-Hochberg method (denoted as $p.adjust$). ABRicate
374 (RRID:SCR_021093) [45] was utilized to annotate antibiotic resistance genes ($identity > 50\%$ and
375 $qcov > 50\%$) in each plasmid bin, based on four antibiotic resistance gene databases [46-49]. The

376 Tukey's Honest Significant Difference test was performed to compare the identified resistance
377 genes among different MOB classes. All statistical analyses were conducted using R.

378

379 **3. Results**

380 **3.1. MOB typing of plasmid genomes**

381 To construct the benchmark datasets, we obtained 90,395 complete plasmid genomes and
382 categorized them into 11 MOB categories using blast (Table 2). We removed 22,470 of them
383 potentially classified into more than one MOB class, leaving 67,925 classified genomes for the
384 training and optimization of MOBFinder (Figure 2A). Our analysis results revealed significant
385 differences in the number, average length, and GC content of plasmid genomes among MOB types.
386 Notably, non-MOB types included the genomes with the most and longest average length, whereas
387 MOBB and MOBM had the fewest plasmid genomes and shortest average length, respectively. In
388 terms of GC content, MOBL had the lowest and MOBQ had the highest amounts. Moreover,
389 plasmids of different MOB types exhibited diverse host ranges at the genus level (Figure 2B).
390 MOBB was predominantly found in *Bacteroides*, *Hymenobacter*, *Parabacteroides*, *Phocaeicola*
391 and *Spirosoma*. Particularly, *Phocaeicola* has been detected in the human gut and possessed the
392 gene for porphyrin degradation through horizontal gene transfer [50]. MOBC, MOBF, MOBH, and
393 MOBP were all found in *Escherichia* and *Klebsiella*. And *Klebsiella* is a multidrug-resistant
394 bacterium that has demonstrated resistance to multiple antibiotics [51]. MOBL, MOBT, and MOBV
395 were mainly discovered in *Bacillus* and *Enterococcus*. Almost all MOBM type plasmid genomes
396 were present in *Clostridium* and *Enterocloster*, and some species in *Clostridium* could cause various
397 diseases [52]. MOBQ demonstrated a broader host range, including *Acinetobacter*, *Agrobacterium*,
398 *Escherichia*, *Rhizobium*, *Lactiplantibacillus*, and *Staphylococcus*. Non-MOB plasmids were
399 detected in the majority of bacteria. These results illustrate the relationship between different MOB
400 types and their host ranges, and also demonstrate that MOB typing of plasmid fragments is feasible
401 in the absence of relaxases.

402

403 **Table 2.** Number, average length, and GC content of plasmid genomes for each MOB type.

404

405 **Figure 2. Benchmark dataset construction using a high-resolution strategy.** (A) Proportion of

406 classified plasmid genomes. A confidence level of ‘sure’ means that the classified plasmid genomes
407 had a *mob_score* of more than 0.5 and an *e-value* of less than 1e-10, while ‘possible’ did not. Plasmid
408 genomes identified as ‘sure’ were used to generate benchmark datasets. Non-MOB, non-mobilizable
409 plasmid. (B) Host range of the classified plasmid genomes at the genus level. Different colors
410 represent different genera, and genera accounting for less than 5% of the total abundance are
411 grouped under the category ‘other.’

412

413 **3.2. Overall performance of MOBFinder**

414 We evaluated the overall performance of MOBFinder in terms of *accuracy*, *kappa*, and *run time*,
415 and compared the tool to MOBscan and MOB-suite. MOBscan did not perform well, achieving low
416 *accuracy* and *kappa* values across sequences of varying lengths, while MOB-suite exhibited
417 marginally better performance than MOBscan when handling sequences of greater length (Figure
418 3A, 3B). In comparison, the *accuracy* of MOBFinder ranged from 70% to 77%, a significant
419 improvement of at least 59% over MOB-suite (Figure 3A). The *kappa* of MOBFinder ranged
420 between 67% and 75% and was approximately 65% higher than that of MOB-suite (Figure 3B).
421 Moreover, MOBFinder exhibited a shorter *run time* in the test dataset, with a more gradual increase
422 trend (Figure 3C). In general, these results indicate that MOBFinder greatly outperformed the other
423 tools, and consistently improved in accuracy and consistency as the sequence length increased.

424

425 **Figure 3. Overall performance of MOBFinder and comparison to MOB-suite and MOBScan.**

426 Evaluation and comparison in terms of (A) *accuracy*, (B) *kappa*, and (C) *run time* (C). The four
427 fragment length groups in the test dataset were Group A (801-1200 bp), Group B (1201-1600 bp),
428 Group C (3000-4000 bp), and Group D (500-10000 bp). (D) For each MOB type, the *balanced*
429 *accuracy*, *harmonic mean*, and *F1-score* were used to assess the performance of MOBFinder and
430 compared to MOB-suite and MOBscan. Since MOB-suite and MOBscan do not include the
431 prediction of MOBL, only the results of MOBL from MOBFinder are provided. MOBFinder, MOB-
432 suite and MOBscan are represented by blue lines, orange lines and gray lines respectively.

433

434 **3.3. Evaluation by MOB category**

435 Next, to evaluate the discrimination ability of MOBFinder for each MOB type, we calculated the

436 *balanced accuracy*, *harmonic mean*, and *F1-score* using the test dataset (Figure 3D). It
437 demonstrated the highest performance for MOBB and MOBM, while its ability to identify non-
438 MOB types was comparatively low. For MOBM, the *balanced accuracy* and *harmonic mean*
439 reached up to 99% and the *F1-score* exceeded 96% for all length groups. For non-MOB, the
440 *balanced accuracy* was 65%, the *harmonic mean* was 49%, and the *F1-score* was 40%. Compared
441 to MOB-suite, MOBFinder exhibited much better performance in predicting all MOB classes. Even
442 for non-MOB, it showed an approximate 13% improvement over the other tools in terms of
443 *balanced accuracy*, 34% in terms of *harmonic mean*, and 24% in terms of *F1-score*.

444

445 In AUC analyses (Figure 4), all values were greater than 0.8, indicating that the tool effectively
446 distinguished between positive and negative samples in each MOB class. In fact, most values were
447 higher than 0.9, except for MOBT and non-MOB. The performance differences by MOB type might
448 be attributable to the differences in host ranges and sequence features among types. Additionally,
449 the imbalance in the training dataset for each MOB type may also be a primary factor contributing
450 to the performance disparities.

451

452 **Figure 4.** ROC curves and AUC values for MOBFinder. The curves were plotted using the output
453 scores of MOBFinder, and the AUC values were calculated to quantify the performance of the tool
454 for each MOB class.

455

456 **3.4. Application to T2D metagenomic data**

457 In a previous study, enrichment analysis of fecal samples identified antibiotic resistance pathways
458 in patients with T2D [40]. The precise mechanism of this enrichment, however, remained elusive.
459 We used MOBFinder to analyze real T2D metagenomic data [39]. After preprocessing and assembly,
460 2,217,064 metagenomic fragments were generated, and plasmid assemblies were identified using
461 PPR-Meta. Subsequently, the plasmid fragments were clustered into 55 bins and annotated using
462 MOBFinder. By employing MOBFinder, we assigned 2 bins to the MOBF class, 8 bins to MOBL,
463 17 bins to MOBQ, and identified 28 bins as non-MOB (Figure 5A). Furthermore, we detected 15
464 bins that exhibited significant differences between the T2D group and a control group. Among them,
465 1 bin was classified as MOBF, 2 as MOBL, 5 as MOBQ, and 7 as non-MOB (Figure S2). Among

466 above MOB types, MOBQ contains the highest number of bins enriched in T2D, while MOBF is
467 widely present in *Escherichia* and *Klebsiella* (Figure 2B), and that some strains of *Klebsiella* are
468 resistant to multiple antibiotics, including carbapenems [53], these two MOB types might contribute
469 to antibiotic resistance in T2D patients. Indeed, when we compared the average abundance of each
470 MOB type between the T2D group and the control group (Figure 5B), the abundances of MOBF
471 and MOBQ were significantly greater in the T2D group.

472

473 **Figure 5.** Annotation of T2D-related plasmid bins using MOBFinder. (A) Heatmap of plasmid bins
474 between T2D patients and controls. Each column represents a sample, and each row represents a
475 plasmid bin. (B) Comparison of the abundance of the four identified MOB types between T2D
476 patients and controls. The *p-value* was calculated using the Wilcoxon rank-sum test, adjusted using
477 the Benjamini-Hochberg method for multiple comparisons. (**p.adjust* < 0.05, ***p.adjust* < 0.01,
478 and ****p.adjust* < 0.001.)

479

480 In addition, these two MOB types can be transferred among multiple bacterial species. This suggests
481 that an increase in these two MOB types could potentially raise the risk of bacterial infection among
482 individuals with T2D. Subsequently, we used four databases [46-49] to detect drug resistance genes
483 in four MOB types (Figure 6). The number of such genes was significantly higher in MOBF than in
484 the other three MOB types. This suggests that MOBF plasmids may carry more drug resistance
485 genes than the other MOB types. Furthermore, the increase in MOBF and MOBQ plasmids could
486 result in more bacteria acquiring drug resistance genes, thereby leading to more antibiotic resistance
487 pathways in T2D patients. In summary, our results demonstrate the utility of MOBFinder for
488 annotating plasmid fragments in metagenomes, uncovering the potential mechanisms underlying
489 the antibiotic resistance enrichment in metagenomic analysis.

490

491 **Figure 6.** Comparison of resistance genes among different MOB types. Four databases were used
492 to identify antibiotic resistance gene within each MOB type, and the *p-value* was calculated using
493 Tukey's Honest Significant Difference test. The two groups without significance markings indicate
494 no statistical difference. (**p-value* < 0.05, ***p-value* < 0.01 and ****p-value* < 0.001.)

495

496 **3.5. Use of MOBFinder**

497 MOBFinder can predict the MOB type of plasmid fragments and bins in metagenomics. For PMFs,
498 it takes a FASTA file as input. The output file consists of 13 columns. The first column represents
499 the fragment ID, the second column displays the predicted MOB type, and columns 3 to 13 represent
500 the scores for each MOB class, namely MOBB, MOBC, MOBF, MOBH, MOBL, MOBM, MOBP,
501 MOBQ, MOBT, MOBV, and non-MOB.

502

503 For plasmid metagenomic bins, MOBFinder requires two input files: a FASTA file containing the
504 plasmid fragments and a meta table that records the mapping between plasmid fragment IDs and
505 bin IDs. The output results are similar to those of plasmid fragments. The first column is the plasmid
506 bin ID. The second is the predicted MOB class of the plasmid bins. The other columns present the
507 MOB scores of the different MOB types.

508

509 **4. Discussion**

510 We developed MOBFinder based on a language model and the random forest algorithm to classify
511 plasmid fragments and bins from metagenomics data into MOB types. First, using the relaxase-
512 alignment method, plasmid genomes were classified into distinct MOB categories. Analyses
513 revealed substantial differences in parameters such as the number, average length, and GC content
514 of plasmid genomes across MOB types. Additionally, there were noteworthy differences in the host
515 ranges among different MOB classes. These results suggest the potential of utilizing sequence
516 features from different MOB types for PMF MOB typing. To characterize the plasmids within each
517 MOB type, we used the skip-gram model to generate word vectors. Our tool demonstrated superior
518 overall performance compared to other tools. Specifically, for each MOB category, MOBFinder
519 exhibited significant improvements in *balanced accuracy*, *harmonic mean*, and *F1-score*, with
520 values reaching up to 99% for the first two measures in the MOBM category.

521

522 Traditionally, *k*-mer frequency models and one-hot encoding have commonly been employed to
523 digitize biological sequences, extensively applied across various machine learning algorithms [54].

524 However, both models simply mark or count the frequency of various characters in sequences,
525 failing to reflect the biological significance underlying each character. These models may also

526 encounter dimensionality issues [54]. For instance, in the k -mer model, if k is set to 8, the
527 dimensionality of the k -mer vector of each DNA sequence becomes 4^8 , which is problematic in
528 metagenomics where most fragment lengths do not reach this magnitude. This would result in
529 significant noise in the feature vector and cause overfitting. Similarly, in the one-hot model, for a
530 sequence of length L using 4-mers as the base unit, it would require L one-hot vectors each with a
531 dimensionality of 4^4 . In such instances, if the dataset for training is not sufficiently large, this
532 representation method could also lead to overfitting due to high dimensionality. In contrast, word
533 vector models offer a superior solution to these problems. Such models initially perform a random
534 initialization of vectors for each “word.” Taking the skip-gram algorithm utilized in this study as an
535 example, the dimension of a random vector can be 1-of- n , where n represents the size of the
536 vocabulary [30]. Following unsupervised pre-training on large datasets, the algorithm maps
537 characters with similar contexts to similar feature spaces. The dimensions of the coordinates (i.e.,
538 the word vectors) of these feature spaces will be lower than those of the initial random vectors. Thus,
539 through unsupervised pre-training on large datasets, language models can compress high-
540 dimensional initial vectors into lower-dimensional word vectors (e.g., MOBFinder’s word vectors
541 have a dimensionality of 100), enabling the feature vectors to contain more character information
542 while effectively avoiding dimensionality issues during supervised training.

543

544 In a metagenomic sequences classification task, 4-mer is widely used as the basic unit in various
545 bioinformatics tools [55], thus MOBFinder takes this as a “word.” To assess the impact of training
546 word vectors with different k -mer lengths on performance, we compared models with k -mer lengths
547 of 2, 3, 4, 5, 6, 7, and 8 (Figure S3). We observed lower overall *accuracy* and *kappa* values for $k=2$.
548 At $k=4$, the *balanced accuracy*, *harmonic mean*, *F1-score*, and *AUC* values stabilized across
549 different MOB types. Subsequently, as the k -mer length increased, there was no significant
550 improvement in *accuracy* or other metrics, while the *run time* gradually increased. Therefore, we
551 chose a k -mer length of 4 for training word vectors and developing MOBFinder.

552

553 Interestingly, in an analysis of T2D metagenomic sequencing data [39], we noted a significant
554 increase in MOBF and MOBQ type plasmids in T2D patients. Moreover, we found more drug
555 resistance genes in the MOBF class, whose dominant hosts are *Klebsiella* and *Escherichia*, which

556 are associated with the spread of multidrug resistance. Although previous analyses of gut
557 metagenomic data from patients with T2D have reported enrichment of drug resistance pathways
558 [40], our results suggest a potential reason for it: the increased abundance of MOB F and MOB Q
559 type plasmids in the guts of individuals with T2D may disseminate more antibiotic resistance genes,
560 resulting in such enrichment.

561

562 At present, databases contain a large amount of human metagenomic data derived from second-
563 generation sequencing. However, understanding of the functions of numerous disease-linked
564 microbial sequences remains limited, attributable to the incomplete nature of metagenomic
565 fragments. The development of MOB Finder enables MOB annotation for plasmid fragments from
566 metagenomics data and provides a powerful tool for investigating the transmission mechanisms of
567 plasmid-mediated antibiotic resistance genes and virulence factors.

568

569 **5. Conclusions**

570 In summary, MOB Finder is a tool for MOB typing of plasmid fragments and bins from metagenomic
571 data. Analyses of classified plasmid genomes unveiled notable differences in sequence
572 characteristics and host ranges across MOB types. Hence, we employed a language model to extract
573 the sequence features specific to each MOB type and represented them using word vectors.
574 Additionally, we boosted prediction accuracy by training and integrating several random forest
575 classification models. MOB Finder surpassed other tools in performance tests and successfully
576 detected an increase in certain MOB type plasmids in T2D patients. Importantly, these MOB type
577 plasmids harbor potential drug-resistance genes, thus offering an explanation for the observed
578 antibiotic resistance in T2D individuals. This suggests that MOB Finder could potentially aid the
579 formulation of specific medications to curb drug resistance transmission. We anticipate that
580 MOB Finder will be a powerful tool for the analysis of plasmid-mediated transmission.

581

582 **Availability of Source Code and Requirements**

- 583 • Project name: MOB Finder
- 584 • Project homepage: <https://github.com/FengTaoSMU/MOBFinder>
- 585 • Operating system(s): Linux

- 586 • Programming language: Python, R script
- 587 • Other requirements: BLAST, biopython
- 588 • License: GPL-3.0
- 589 • RRID: SCR_024451
- 590 • biotoolsID: MOBFinder

591

592 **Data Availability**

593 Snapshots of our code and other data further supporting this work are openly available in the
594 GigaScience repository, GigaDB [56].

595

596 **Abbreviations**

597 MOB: mobilization typing; Rep: replicon typing; Inc: incompatibility typing; MPF: mate-pair
598 formation typing; non-MOB: non-mobilizable; T2D: type 2 diabetes; TPR: true positive rate; TNR:
599 true negative rate; Sn: sensitivity; Sp: specificity; ROC: the receiver operating characteristic; AUC:
600 the area under the curve; SRA: short read archive; NCBI: National Center for Biotechnology
601 Information; PlasTax-PCR: Plasmid Taxonomic PCR; PBRT: PCR-based replicon typing; DPMT:
602 degenerate primer MOB typing.

603

604 **Competing Interests**

605 The authors declare that they have no competing interests.

606

607 **Funding**

608 This investigation was financially supported by the National Key R&D Program of China
609 (2022YFA0806400) and the National Natural Science Foundation of China (82102508, 81925026).

610

611 **Author Contributions**

612 TF, ZCF, and HWZ proposed and designed this work. TF and ZCF developed and optimized the
613 software. TF, ZCF, SFW, and HWZ wrote and revised the manuscript.

614

615 **Supplementary data**

616

617 **Supplementary Table 1.** Accessions list of classified plasmid genomes.

618 **Supplementary Table 2.** List of metagenomic samples used in our analysis.

619 **Supplementary Figure 1.** MOB typing using MOB-suite. Single-class, plasmid genomes classified
620 into one MOB type; multi-class, plasmid genomes classified into more than one MOB category;
621 non-MOB, non-mobilizable plasmids.

622 **Supplementary Figure 2.** Abundance of each significantly different plasmid bin from various
623 MOB types between patients with type II diabetes and controls.

624 **Supplementary Figure 3.** Comparison results for the development of MOBFinder using word
625 vectors trained with different *k*-mer lengths. (A-C) Overall *accuracy*, *kappa*, and *run time* of the
626 MOB classification model trained with word vectors trained using different lengths of *k*-mers. (D)
627 *Balanced accuracy*, *harmonic mean*, *F1-score*, and AUC of word vectors trained with different *k*-
628 mer lengths across different MOB types.

629

630 **References**

- 631 1. Helinski DR. A Brief History of Plasmids. *EcoSal Plus*. 2022 Dec 15;10(1):eESP00282021.
632 doi: 10.1128/ecosalplus.esp-0028-2021
- 633 2. Garcillán-Barcia MP, Francia MV, de la Cruz F. The diversity of conjugative relaxases and its
634 application in plasmid classification. *FEMS Microbiol Rev*. 2009 May;33(3):657-87. doi:
635 10.1111/j.1574-6976.2009.00168.x
- 636 3. Rodríguez-Beltrán J, DelaFuente J, León-Sampedro R, MacLean RC, San Millán Á. Beyond
637 horizontal gene transfer: the role of plasmids in bacterial evolution. *Nat Rev Microbiol*. 2021
638 Jun;19(6):347-359. doi: 10.1038/s41579-020-00497-1
- 639 4. Shintani M, Sanchez ZK, Kimbara K. Genomics of microbial plasmids: classification and
640 identification based on replication and transfer systems and host taxonomy. *Front Microbiol*.
641 2015 Mar 31;6:242. doi: 10.3389/fmicb.2015.00242
- 642 5. Redondo-Salvo S, Bartomeus-Peñalver R, Vielva L, Tagg KA, et al. COPLA, a taxonomic
643 classifier of plasmids. *BMC Bioinformatics*. 2021 Jul 31;22(1):390. doi: 10.1186/s12859-021-
644 04299-x

- 645 6. Carattoli A, Hasman H. PlasmidFinder and In Silico pMLST: Identification and Typing of
646 Plasmid Replicons in Whole-Genome Sequencing (WGS). *Methods Mol Biol.* 2020;2075:285-
647 294. doi: 10.1007/978-1-4939-9877-7_20
- 648 7. Smillie C, Garcillán-Barcia MP, Francia MV, Rocha EP, de la Cruz F. Mobility of plasmids.
649 *Microbiol Mol Biol Rev.* 2010 Sep;74(3):434-52. doi: 10.1128/MMBR.00020-10
- 650 8. Francia MV, Varsaki A, Garcillán-Barcia MP, Latorre A, Drainas C, de la Cruz F. A
651 classification scheme for mobilization regions of bacterial plasmids. *FEMS Microbiol Rev.*
652 2004 Feb;28(1):79-100. doi: 10.1016/j.femsre.2003.09.001
- 653 9. Garcillán-Barcia MP, Francia MV, de la Cruz F. The diversity of conjugative relaxases and its
654 application in plasmid classification. *FEMS Microbiol Rev.* 2009 May;33(3):657-87. doi:
655 10.1111/j.1574-6976.2009.00168.x
- 656 10. Garcillán-Barcia MP, Alvarado A, de la Cruz F. Identification of bacterial plasmids based on
657 mobility and plasmid population biology. *FEMS Microbiol Rev.* 2011 Sep;35(5):936-56. doi:
658 10.1111/j.1574-6976.2011.00291.x
- 659 11. Bradley P, den Bakker HC, Rocha EPC, McVean G, Iqbal Z. Ultrafast search of all deposited
660 bacterial and viral genomic data. *Nat Biotechnol.* 2019 Feb;37(2):152-159. doi:
661 10.1038/s41587-018-0010-1
- 662 12. Alvarado A, Garcillán-Barcia MP, de la Cruz F. A degenerate primer MOB typing (DPMT)
663 method to classify gamma-proteobacterial plasmids in clinical and environmental settings.
664 *PLoS One.* 2012;7(7):e40438. doi: 10.1371/journal.pone.0040438
- 665 13. Garcillán-Barcia MP, Alvarado A, de la Cruz F. Identification of bacterial plasmids based on
666 mobility and plasmid population biology. *FEMS Microbiol Rev.* 2011 Sep;35(5):936-56. doi:
667 10.1111/j.1574-6976.2011.00291.x
- 668 14. Cuartas R, Coque TM, de la Cruz F, Garcillán-Barcia MP. PLASmid TAXonomic PCR
669 (PlasTax-PCR), a Multiplex Relaxase MOB Typing to Assort Plasmids into Taxonomic Units.
670 *Methods Mol Biol.* 2022;2392:127-142. doi: 10.1007/978-1-0716-1799-1_10
- 671 15. Carattoli A, Bertini A, Villa L, Falbo V, Hopkins KL, Threlfall EJ. Identification of plasmids
672 by PCR-based replicon typing. *J Microbiol Methods.* 2005 Dec;63(3):219-28. doi:
673 10.1016/j.mimet.2005.03.018
- 674 16. Fang Z, Zhou H. Identification of the conjugative and mobilizable plasmid fragments in the

- 675 plasmidome using sequence signatures. *Microb Genom.* 2020 Nov;6(11):mgen000459. doi:
676 10.1099/mgen.0.000459
- 677 17. Li X, Xie Y, Liu M, Tai C, Sun J, Deng Z, Ou HY. oriTfinder: a web-based tool for the
678 identification of origin of transfers in DNA sequences of bacterial mobile genetic elements.
679 *Nucleic Acids Res.* 2018 Jul 2;46(W1):W229-W234. doi: 10.1093/nar/gky352
- 680 18. Garcillán-Barcia MP, Redondo-Salvo S, Vielva L, de la Cruz F. MOBscan: Automated
681 Annotation of MOB Relaxases. *Methods Mol Biol.* 2020;2075:295-308. doi: 10.1007/978-1-
682 4939-9877-7_21
- 683 19. Robertson J, Nash JHE. MOB-suite: software tools for clustering, reconstruction and typing of
684 plasmids from draft assemblies. *Microb Genom.* 2018 Aug;4(8):e000206. doi:
685 10.1099/mgen.0.000206
- 686 20. Robertson J, Bessonov K, Schonfeld J, Nash JHE. Universal whole-sequence-based plasmid
687 typing and its utility to prediction of host range and epidemiological surveillance. *Microb*
688 *Genom.* 2020 Oct;6(10):mgen000435. doi: 10.1099/mgen.0.000435
- 689 21. Krawczyk PS, Lipinski L, Dziembowski A. PlasFlow: predicting plasmid sequences in
690 metagenomic data using genome signatures. *Nucleic Acids Res.* 2018 Apr 6;46(6):e35. doi:
691 10.1093/nar/gkx1321
- 692 22. Roosaare M, Puustusmaa M, Möls M, Vaher M, Remm M. PlasmidSeeker: identification of
693 known plasmids from bacterial whole genome sequencing reads. *PeerJ.* 2018 Apr 2;6:e4588.
694 doi: 10.7717/peerj.4588
- 695 23. Pellow D, Mizrahi I, Shamir R. PlasClass improves plasmid sequence classification. *PLOS*
696 *Comput Biol.* 2020;16:e1007781. doi: 10.1371/journal.pcbi.1007781
- 697 24. Fang Z, Tan J, Wu S, Li M, Xu C, Xie Z, et al. PPR-Meta: a tool for identifying phages and
698 plasmids from metagenomic fragments using deep learning. *GigaScience.* 2019;8:1–14.
699 10.1093/gigascience/giz066
- 700 25. Pradier L, Tissot T, Fiston-Lavier AS, Bedhomme S. PlasForest: a homology-based random
701 forest classifier for plasmid detection in genomic datasets. *BMC Bioinformatics.* 2021 Jun
702 26;22(1):349. doi: 10.1186/s12859-021-04270-w
- 703 26. Sobecky PA, Hazen TH. Horizontal gene transfer and mobile genetic elements in marine
704 systems. *Methods Mol Biol.* 2009;532:435-53. doi: 10.1007/978-1-60327-853-9_25

- 705 27. Suzuki H, Yano H, Brown CJ, Top EM. Predicting plasmid promiscuity based on genomic
706 signature. *J Bacteriol.* 2010 Nov;192(22):6045-55. doi: 10.1128/JB.00277-10
- 707 28. Wu S, Fang Z, Tan J, Li M, Wang C, Guo Q, Xu C, Jiang X, Zhu H. DeePhage: distinguishing
708 virulent and temperate phage-derived sequences in metavirome data with a deep learning
709 approach. *Gigascience.* 2021 Sep 8;10(9):giab056. doi: 10.1093/gigascience/giab056
- 710 29. Fang Z, Feng T, Zhou H, Chen M. DeePVP: Identification and classification of phage virion
711 proteins using deep learning. *Gigascience.* 2022 Aug 11;11:giac076. doi:
712 10.1093/gigascience/giac076
- 713 30. Mikolov T, Chen K, Corrado G, and Dean J. Efficient estimation of word representations in
714 vector space. 2013. arXiv preprint. doi: arXiv:1301.3781.
- 715 31. Patrick Ng. dna2vec: Consistent vector representations of variable-length k-mers. arXiv. doi:
716 10.48550/arXiv.1701.06279
- 717 32. Tsukiyama S, Hasan MM, Fujii S, Kurata H. LSTM-PHV: prediction of human-virus protein-
718 protein interactions by LSTM with word2vec. *Brief Bioinform.* 2021 Nov 5;22(6):bbab228.
719 doi: 10.1093/bib/bbab228
- 720 33. Sharma R, Shrivastava S, Kumar Singh S, Kumar A, Saxena S, Kumar Singh R. Deep-
721 ABPpred: identifying antibacterial peptides in protein sequences using bidirectional LSTM
722 with word2vec. *Brief Bioinform.* 2021 Sep 2;22(5):bbab065. doi: 10.1093/bib/bbab065
- 723 34. Asgari E, Mofrad MR. Continuous Distributed Representation of Biological Sequences for
724 Deep Proteomics and Genomics. *PLoS One.* 2015 Nov 10;10(11):e0141287. doi:
725 10.1371/journal.pone.0141287
- 726 35. Wisniewski JA, Traore DA, Bannam TL, Lyras D, Whisstock JC, Rood JI. TcpM: a novel
727 relaxase that mediates transfer of large conjugative plasmids from *Clostridium perfringens*.
728 *Mol Microbiol.* 2016 Mar;99(5):884-96. doi: 10.1111/mmi.13270
- 729 36. Ramachandran G, Miguel-Arribas A, Abia D, Singh PK, Crespo I, Gago-Córdoba C, Hao JA,
730 Luque-Ortega JR, Alfonso C, Wu LJ, Boer DR, Meijer WJ. Discovery of a new family of
731 relaxases in Firmicutes bacteria. *PLoS Genet.* 2017 Feb 16;13(2):e1006586. doi:
732 10.1371/journal.pgen.1006586
- 733 37. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. BLAST+:
734 architecture and applications. *BMC Bioinformatics.* 2009 Dec 15;10:421. doi: 10.1186/1471-

735 2105-10-421

736 38. Seemann T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics*. 2014 Jul
737 15;30(14):2068-9. doi: 10.1093/bioinformatics/btu153

738 39. Qin J, Li Y, Cai Z, Li S, Zhu J, Zhang F et al. A metagenome-wide association study of gut
739 microbiota in type 2 diabetes. *Nature*. 2012 Oct 4;490(7418):55-60. doi: 10.1038/nature11450

740 40. Wu H, Tremaroli V, Schmidt C, Lundqvist A, Olsson LM, Krämer M, Gummesson A, Perkins
741 R, Bergström G, Bäckhed F. The Gut Microbiota in Prediabetes and Diabetes: A Population-
742 Based Cross-Sectional Study. *Cell Metab*. 2020 Sep 1;32(3):379-390.e3. doi:
743 10.1016/j.cmet.2020.06.011

744 41. Schmieder R, Edwards R. Quality control and preprocessing of metagenomic datasets.
745 *Bioinformatics*. 2011 Mar 15;27(6):863-4. doi: 10.1093/bioinformatics/btr026

746 42. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Methods*. 2012 Mar
747 4;9(4):357-9. doi: 10.1038/nmeth.1923

748 43. Li D, Liu CM, Luo R, Sadakane K, Lam TW. MEGAHIT: an ultra-fast single-node solution
749 for large and complex metagenomics assembly via succinct de Bruijn graph. *Bioinformatics*.
750 2015 May 15;31(10):1674-6. doi: 10.1093/bioinformatics/btv033

751 44. Lu YY, Chen T, Fuhrman JA, Sun F. COCACOLA: binning metagenomic contigs using
752 sequence COmposition, read CoverAge, CO-alignment and paired-end read LinkAge.
753 *Bioinformatics*. 2017 Mar 15;33(6):791-798. doi: 10.1093/bioinformatics/btw290

754 45. Seemann T. Abriicate. Github. <https://github.com/tseemann/abriicate>

755 46. Gupta SK, Padmanabhan BR, Diene SM, Lopez-Rojas R, Kempf M, Landraud L, Rolain JM.
756 ARG-ANNOT, a new bioinformatic tool to discover antibiotic resistance genes in bacterial
757 genomes. *Antimicrob Agents Chemother*. 2014;58(1):212-20. doi: 10.1128/AAC.01310-13

758 47. Jia B, Raphenya AR, Alcock B, Waglechner N, Guo P, Tsang KK, Lago BA, Dave BM, Pereira
759 S, Sharma AN, Doshi S, Courtot M, Lo R, Williams LE, Frye JG, Elsayegh T, Sardar D,
760 Westman EL, Pawlowski AC, Johnson TA, Brinkman FS, Wright GD, McArthur AG. CARD
761 2017: expansion and model-centric curation of the comprehensive antibiotic resistance
762 database. *Nucleic Acids Res*. 2017 Jan 4;45(D1):D566-D573. doi: 10.1093/nar/gkw1004

763 48. Doster E, Lakin SM, Dean CJ, Wolfe C, Young JG, Boucher C, Belk KE, Noyes NR, Morley
764 PS. MEGARes 2.0: a database for classification of antimicrobial drug, biocide and metal

765 resistance determinants in metagenomic sequence data. *Nucleic Acids Res.* 2020 Jan
766 8;48(D1):D561-D569. doi: 10.1093/nar/gkz1010

767 49. Feldgarden M, Brover V, Haft DH, Prasad AB, Slotta DJ, Tolstoy I, Tyson GH, Zhao S, Hsu
768 CH, McDermott PF, Tadesse DA, Morales C, Simmons M, Tillman G, Wasilenko J, Folster JP,
769 Klimke W. Validating the AMRFinder Tool and Resistance Gene Database by Using
770 Antimicrobial Resistance Genotype-Phenotype Correlations in a Collection of Isolates.
771 *Antimicrob Agents Chemother.* 2019 Oct 22;63(11):e00483-19. doi: 10.1128/AAC.00483-19

772 50. Hehemann JH, Correc G, Barbeyron T, Helbert W, Czjzek M, Michel G. Transfer of
773 carbohydrate-active enzymes from marine bacteria to Japanese gut microbiota. *Nature.* 2010
774 Apr 8;464(7290):908-12. doi: 10.1038/nature08937

775 51. Fu S, Wang R, Xu Z, Zhou H, Qiu Z, Shen L and Yang Q. Metagenomic sequencing combined
776 with flow cytometry facilitated a novel microbial risk assessment framework for bacterial
777 pathogens in municipal wastewater without cultivation. *iMeta.*
778 2023;2:e77 .doi:10.1002/imt2.77

779 52. Dieterle MG, Rao K, Young VB. Novel therapies and preventative strategies for primary and
780 recurrent *Clostridium difficile* infections. *Ann N Y Acad Sci.* 2019 Jan;1435(1):110-138. doi:
781 10.1111/nyas.13958

782 53. Yang X, Dong N, Chan EW, Zhang R, Chen S. Carbapenem Resistance-Encoding and
783 Virulence-Encoding Conjugative Plasmids in *Klebsiella pneumoniae*. *Trends Microbiol.* 2021
784 Jan;29(1):65-83. doi: 10.1016/j.tim.2020.04.012

785 54. Jaillard M, Palmieri M, van Belkum A, Mahé P. Interpreting k-mer-based signatures for
786 antibiotic resistance prediction. *Gigascience.* 2020 Oct 17;9(10):giaa110. doi:
787 10.1093/gigascience/giaa110

788 55. Sedlar K, Kupkova K, Provaznik I. Bioinformatics strategies for taxonomy independent
789 binning and visualization of sequences in shotgun metagenomics. *Comput Struct Biotechnol J.*
790 2016 Dec 5;15:48-55. doi: 10.1016/j.csbj.2016.11.005

791 56. Feng T; Wu S; Zhou H; Fang Z. Supporting data for "MOBFinder: A tool for mobilization
792 typing of plasmid metagenomic fragments based on a language model" *GigaScience Database*
793 2024. <https://doi.org/10.5524/102559>

794
795

Table 1. Experimental and computational schemes developed for plasmid classification.

Technology category	Method	Classification scheme	Material	Description
Experimental	DPMT [12]	MOB typing	Plasmid DNA from clinical isolates	Used degenerate primers to hybridize relaxase-coding genes to identify and classify plasmids isolated from clinical isolates
	PlasTax-PCR [14]	Taxonomic typing	Plasmid DNA from clinical isolates	Utilized PCR primers that target conserved segments of the relaxase gene of plasmid taxonomic units (PTUs) to identify specific PTUs of transmissible plasmids
	PBRT [15]	Rep typing or Inc typing	Plasmid DNA from clinical isolates	Used multiplex PCR to amplify DNA fragments of replicons and detect known replicon types of plasmids
Computational	MOBscan [18]	MOB typing	Plasmid protein sequences	Used the HMMER model to annotate the relaxases and further perform MOB typing
	MOB-suite [19, 20]	MOB typing, MPF typing and Rep typing	Complete plasmid genomes or plasmid assembly clusters (Linux)	Utilized collected relaxase, oriT, replicon, and T4SS sequences to construct database, then classified plasmid assembly clusters with BLAST
	PlasTans [16]	transmissible plasmid identification	Plasmid assembly contigs (Linux)	Used the convolutional neural network deep learning algorithm to classify plasmid DNA fragments
	PlasmidFinder [6]	Rep typing or Inc typing	Raw reads or complete plasmid genomes or plasmid assembly contigs (web server)	Utilized collected replicon sequences and BLASTn to perform Rep typing and Inc typing
	pMLST [6]	Rep typing or Inc typing	Raw reads or complete plasmid genomes or plasmid assembly contigs (web server)	Used collected plasmid multilocus sequence typing (pMLST) allele sequences, known sequence type profiles, and BLAST to perform Rep typing and Inc typing
	oriTfinder [17]	MOB typing, MPF typing	Complete plasmid genomes (web server)	Utilized collected oriT, relaxase, T4CP, and T4SS sequences to annotate plasmids with BLAST

COPLA [5]

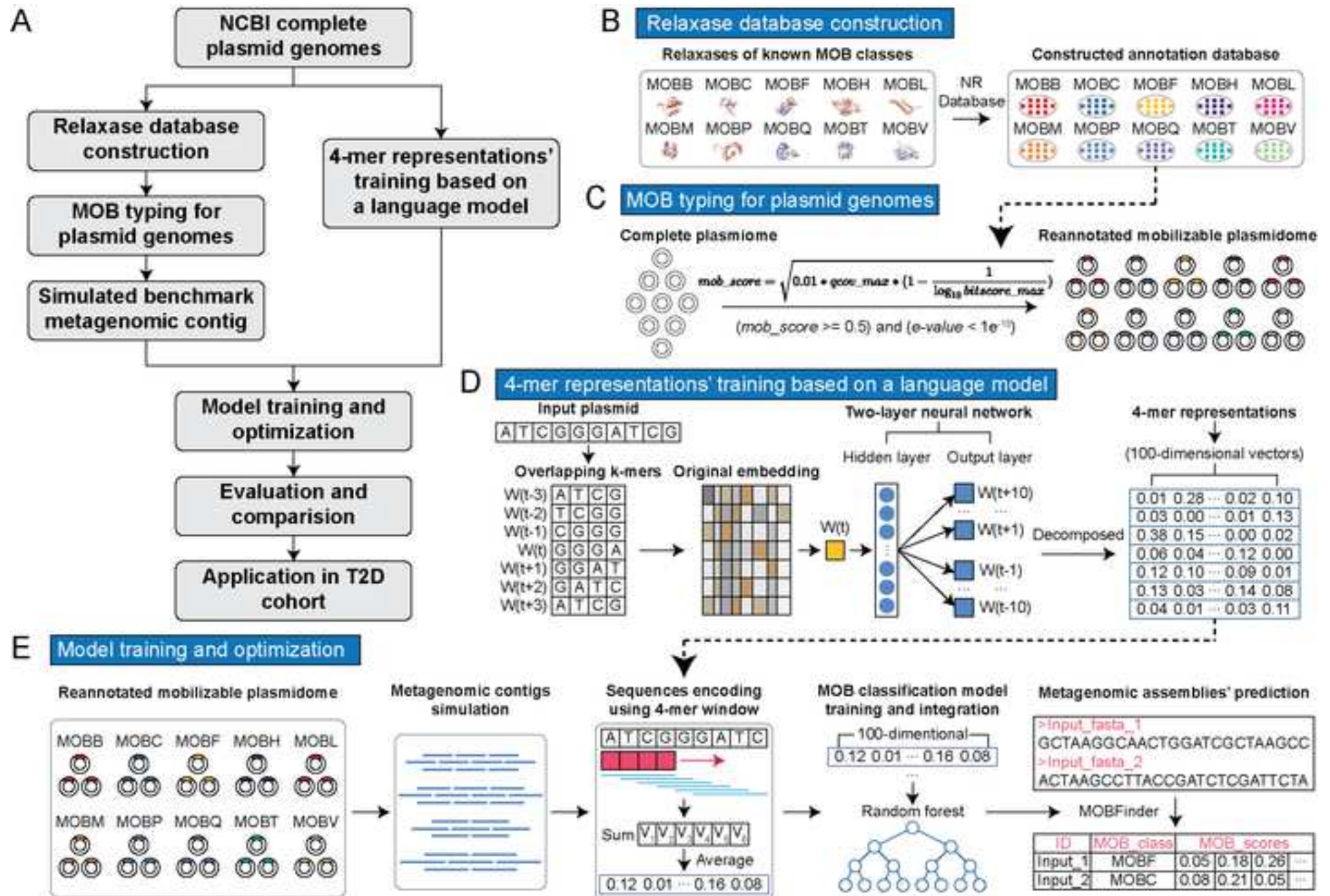
Taxonomic typing

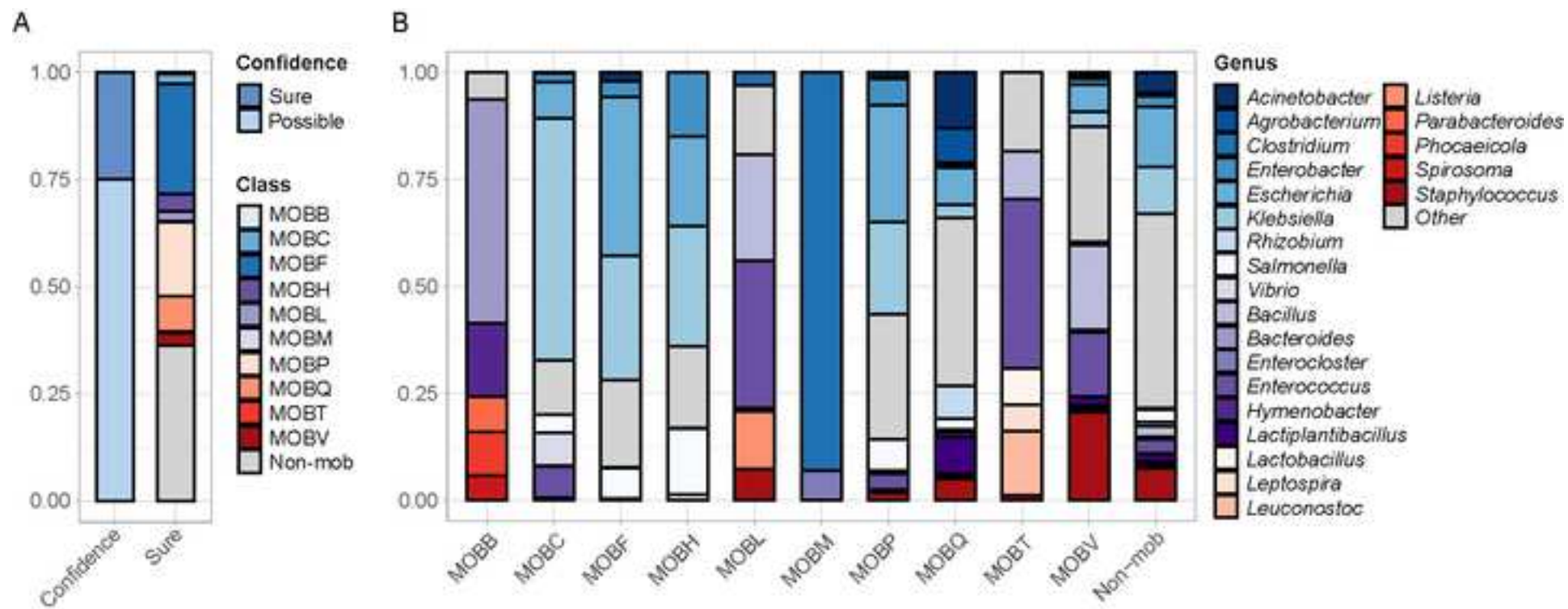
Complete plasmid genomes or
plasmid assembly sets (Linux)

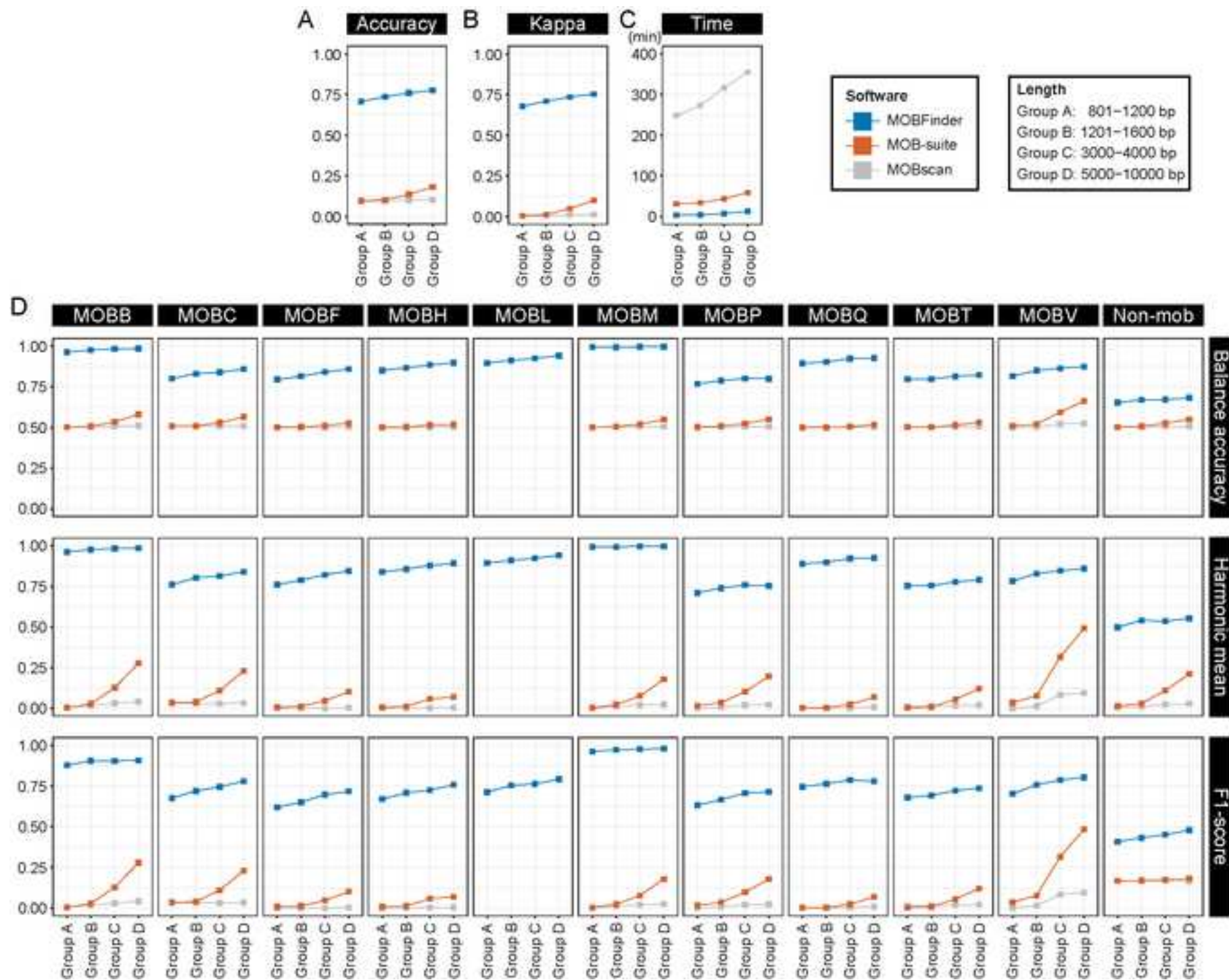
Used average nucleotide identity (ANI) metrics and hierarchical
stochastic block modeling (HSBM) to create plasmid taxonomic
units (PTUs) and predict taxonomic hosts

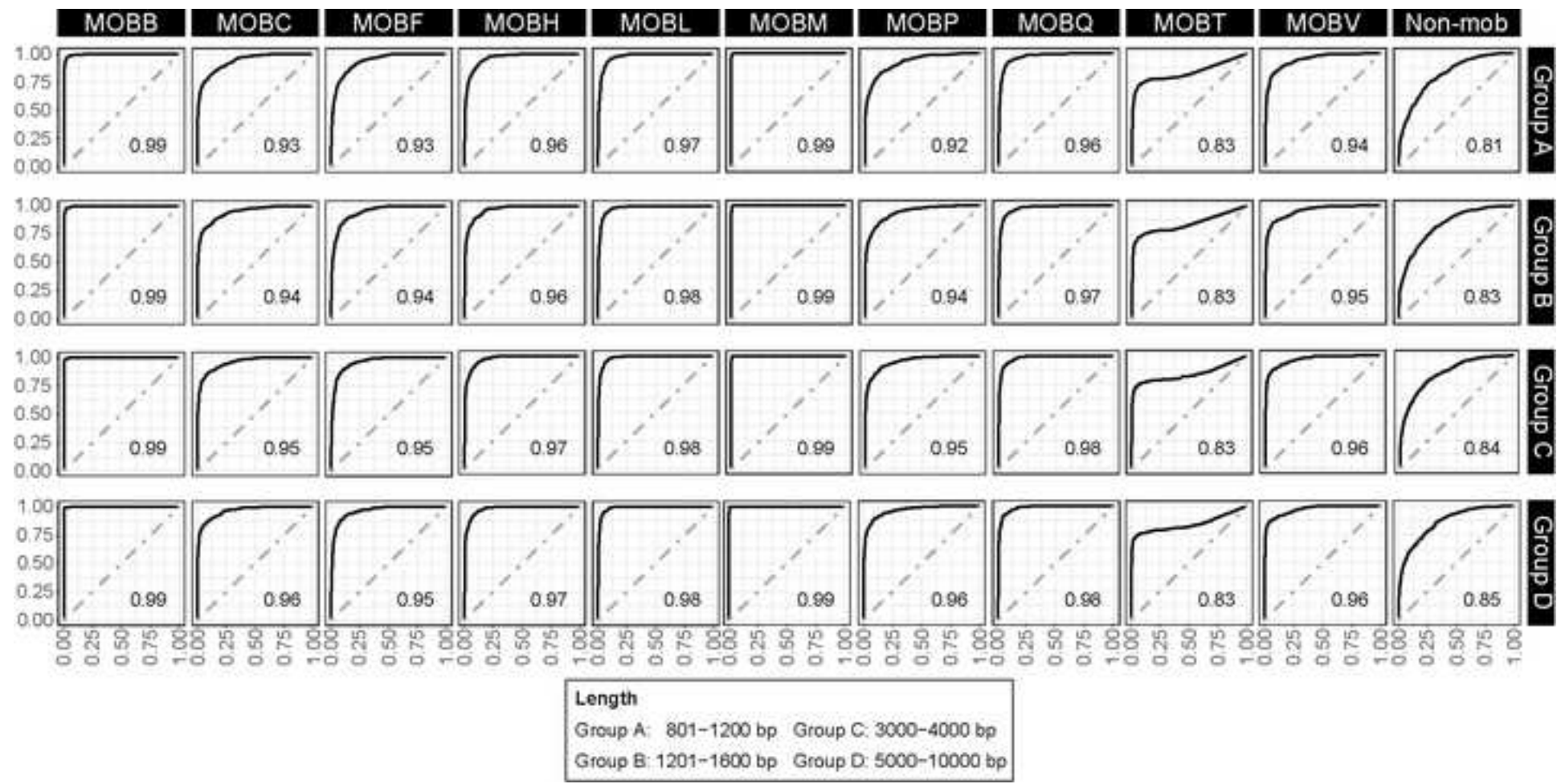
Table 2. Number, average length, and GC content of plasmid genomes for each MOB type.

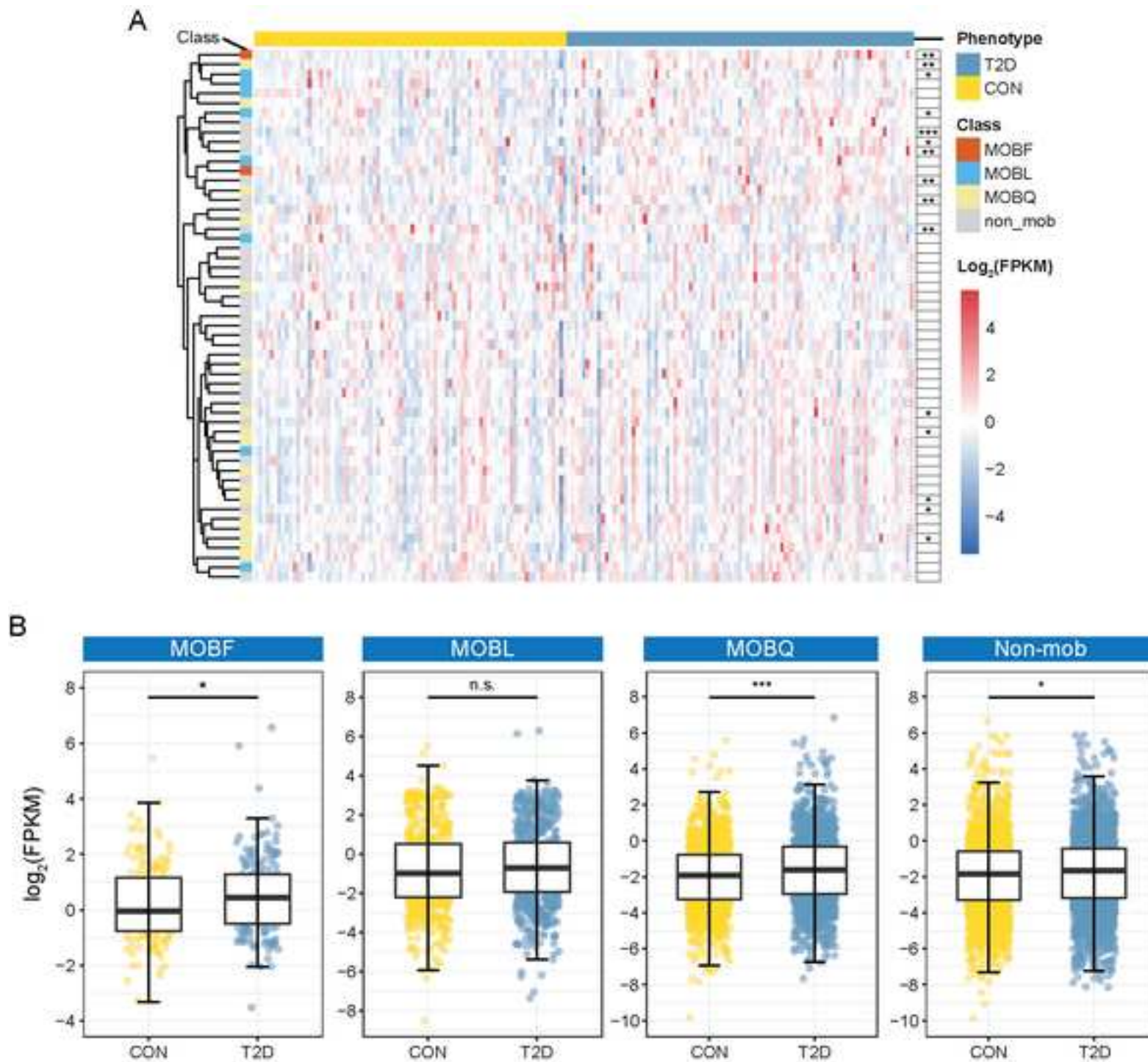
Class	Number	Average length	GC (%)
MOBB	623	10921.77	51.27
MOBC	3218	19965.28	47.14
MOBF	21268	103802.80	52.07
MOBH	4880	151108.10	48.37
MOBL	3446	51430.63	34.57
MOBM	1761	2684.14	27.12
MOBP	15617	32237.88	49.70
MOBQ	9347	89357.64	56.77
MOBT	1181	11643.24	36.92
MOBV	4405	6595.43	37.75
Non-MOB	24649	37581.85	49.84



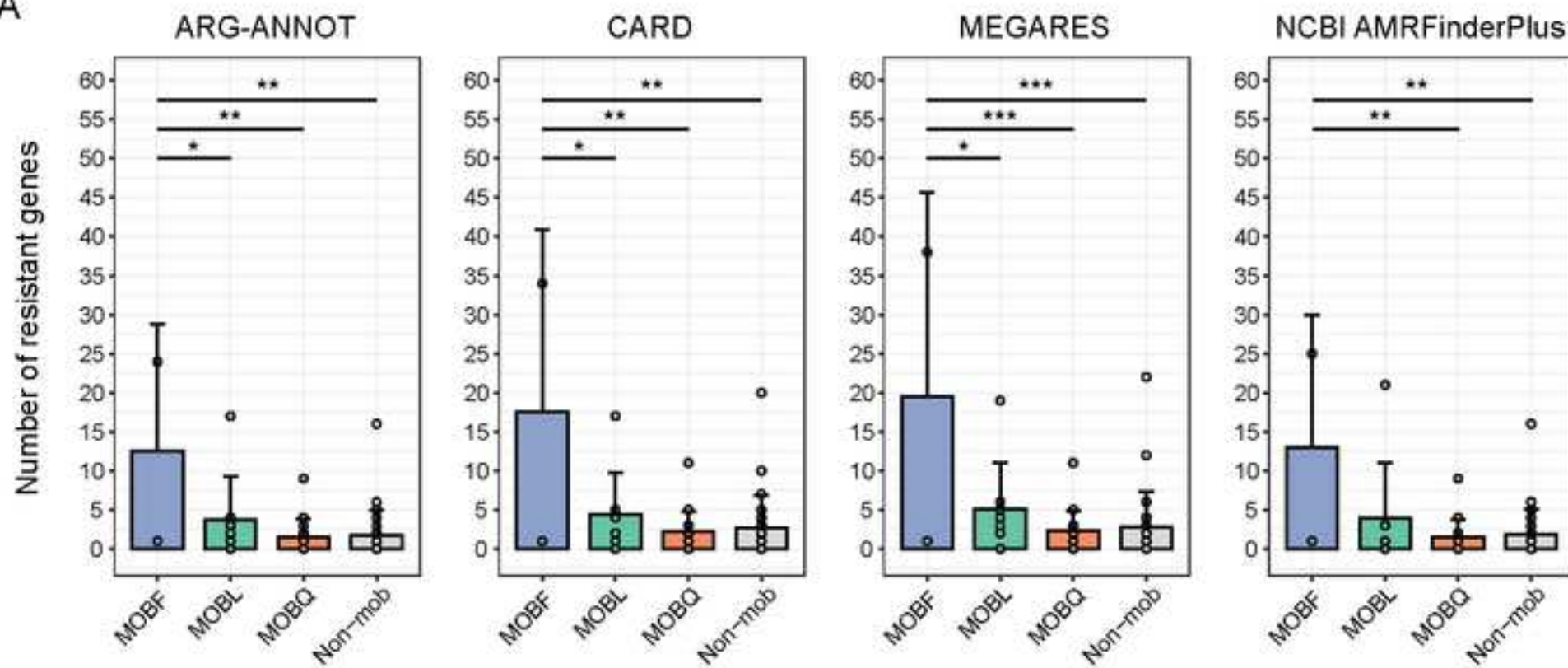








A



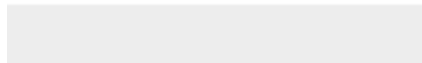


Click here to access/download
Supplementary Material
Suplmentaty Figures.docx





Click here to access/download
Supplementary Material
Supplementary Table 1.xlsx





Click here to access/download
Supplementary Material
Supplementary Table 2.xlsx



Cover Letter

Dear Editor,

Thank you very much for your previous E-mail on May 16, 2024, regarding our manuscript, “MOBFinder: a tool for MOB typing for plasmid metagenomic fragments based on language model” (Manuscript Number: GIGA-D-24-00070). We are very pleased to know that our manuscript is potentially acceptable for publication in the journal, subject to the further revisions suggested by the reviewers. We are very grateful for your substantial and helpful advice with respect to our manuscript, and we are pleased to receive the reviewers’ overall positive comments about our work. We thank the reviewers for their substantial and valuable comments, including their careful reading and checking of the manuscript, which greatly helped us improve the paper.

Our revisions and responses to the editor’s and two reviewers’ comments (italic text) are provided below.

To Editor:

1. One of the reviewers suggested you to improve the language, GigaScience is providing copy editing service, you can contact Qi Chen (chenqi@genomics.cn) if you need.

For this revision, we used the copy editing service of the *GigaScience* journal to refine the language. We thoroughly reviewed the entire text again to ensure that there are no serious errors in spelling, grammar, or meaning in each sentence. Specifically, based on the suggestions from the language editing, we have changed our title to “MOBFinder: **A** tool for **mobilization** typing of plasmid metagenomic fragments based on **a** language model”.

2. In addition, please register any new software application in the bio.tools and SciCrunch.org databases to receive RRID (Research Resource Identification Initiative ID) and biotoolsID identifiers, and include these in your manuscript. Computational workflows should be registered in workflowhub.eu and the DOIs cited in the relevant places in the manuscript. These will facilitate tracking, reproducibility and re-use of your tool.

According to the journal requirement, we have registered the tool in the bio.tools and SciCrunch.org databases. In Lines 585-586 of the “Availability of Source Code and Requirements” section of the revised manuscript, we have added the following statement:

RRID: SCR_024451.
biotoolsID: MOBFinder.

Also, all the related scripts and data have also submitted to the GigaDB server.

To Reviewer #1:

Specific Comments:

1. *the unpaired Wilcoxon signed-rank two-sided test.*

-> *should be corrected to either*

"Wilcoxon rank-sum test" or "Mann-Whitney U test"

https://en.wikipedia.org/wiki/Mann%E2%80%93Whitney_U_test

"Wilcoxon rank-sum test" redirects here. For Wilcoxon signed-rank test, see Wilcoxon signed-rank test.

https://en.wikipedia.org/wiki/Wilcoxon_signed-rank_test

Not to be confused with Wilcoxon rank-sum test.

We apologize for the confusion between the Wilcoxon rank-sum test and the Wilcoxon signed-rank test. We have corrected this mistake in the revised manuscript. In Line 368 and 472 of the revised manuscript, the statistical method has been corrected to “**the Wilcoxon rank-sum test**”.

2. *Since MOBscan can only predict the MOB type with plasmid proteins, we annotated the plasmids in the test set with Prokka, then manually submitted them to the MOBscan website for MOB type annotation.*

Given that MOBScan operates as an online tool and cannot be executed locally, the calculation of MOBScan's run time was confined to the duration spent on preprocessing with Prokka locally." (Please refer to Line 313-319 in the revised manuscript).

-> *Actually, it can be executed locally using the scripts included in <https://github.com/santirdnd/COPLA/>. It may not be necessary to run MOBscan locally (it may be okay that they manually submitted them to the MOBscan website), but I'll*

inform you regardless.

We are very grateful to Reviewer 1 for reminding us that MOBscan can be run locally. In the revised manuscript, we have removed the statement “MOBscan can only predict the MOB type with plasmid proteins” and revised the corresponding description to “**We used the online server of MOBscan to perform the MOB annotation, and...**” (Please refer to Line 312-313 in the revised manuscript).

3. Line 418-421

In the comparison, it was observed that MOBscan did not perform well, achieving low accuracy and kappa values across sequences of varying lengths, while MOB-suite exhibited marginally better performance than MOBscan when handling sequences of greater length (Figure 3A, 3B). (Please refer to Line 418-421 in the revised manuscript).

-> Do the authors' results contradict the following general expectation? MOB-typer utilizes BLAST, whereas MOBscan utilizes hmmscan, and therefore, MOBscan is expected to retrieve more distantly related proteins than MOB-typer.

We would like to thank Reviewer 1 for the discussion regarding BLAST and HMMscan. Firstly, we acknowledge that for more distantly related proteins, sequence searching based on HMM exhibits higher sensitivity than BLAST. However, we believe that our results do not contradict this theory. There are two reasons that might explain why, in this manuscript, the performance of tools based on HMM appears slightly inferior to those based on BLAST.

(1) The number of reference sequences can impact the performance of the tools. In MOB-suite, a large number of reference sequences are used for BLAST sequence alignment, whereas in MOBscan, the number of relaxase sequences used to profile HMM files for some MOB types is not very large. For instance, for the MOBF type, MOBscan utilizes 146 relaxase sequences for configuring HMM files, while MOB-suite employs 396 sequences to construct the BLAST database. The difference in the number of reference sequences could potentially lead to MOBscan's performance being slightly inferior to that of MOB-suite.

(2) The aim of this study is not to design new methods for identifying novel relaxases. In our test data, the relaxases all come from sequenced plasmids, so there is some homology with the relaxases in the database. When the query sequence and the database sequence have high homology, the performance of BLAST may not necessarily be

worse than HMM. In fact, existing studies have shown that methods based on BLAST can sometimes outperform those based on HMM when the homology is high (Ref: PMID: 25140992).

4. *MOB-suit and MOBscan are represented by blue lines, orange lines and gray lines respectively.*

-> *should be*

"MOB-suite"

We thank Reviewer 1 for the careful checking of the manuscript. In Line 427-428 of the revised manuscript, we have revised "MOB-suit" to "**MOB-suite**".

5. *I suggest receiving English language editing before publishing the paper.*

"For the MOB typing, MOBscan [18] uses the HMMER model to annotated the relaxases and further perform MOB typing."

-> *should be*

"For the MOB typing, MOBscan [18] uses the HMMER model to annotate the relaxases and further perform MOB typing."

We are sorry for the grammatical error. In Line 109 of the revised manuscript, we have revised the sentence as "For the MOB typing, MOBscan [18] uses the HMMER model to **annotate** relaxase genes and classify plasmids accordingly".

In addition, we have used the copy editing service of the *GigaScience* journal to refine the language through the whole manuscript.

To Reviewer #2:

General Comments:

I would like to commend you on the revisions made to your manuscript following the initial round of reviews. It is evident that considerable effort has been put into addressing the concerns and suggestions raised during the first review. The changes and additions you have implemented have significantly enhanced the clarity, depth, and scholarly value of your paper. The manuscript has been improved substantially and all the initial concerns have been addressed satisfactorily. I support the publication of this manuscript in GigaScience.

Here, we would like to express our sincere gratitude to Reviewer 2 for the positive

comments on our work, describing our revised manuscript: “*The manuscript has been improved substantially and all the initial concerns have been addressed satisfactorily.*” We are very thankful for Reviewer 2's suggestions during the first revision process, which greatly enhanced the clarity, depth, and academic value of our paper.

In hoping that the above revision has clarified all the points by two reviewers and given a point-by-point response to all the concerns, we hereby submit our revised manuscript to the journal. We thank you for your kind consideration.

Note: The initial version of the article has been made public on bioRxiv (<https://doi.org/10.1101/2023.12.06.570414>), and beyond that, it has not been published in any other traditional journals.

Sincerely yours,

Zhencheng Fang, Ph.D.

Microbiome Medicine Center, Department of Laboratory Medicine, Zhujiang Hospital,
Southern Medical University, Guangzhou, 510280, China

Email: fangzc@smu.edu.cn