

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

Tao Feng<sup>1</sup>, Shufang Wu<sup>1</sup>, Hongwei Zhou<sup>1\*</sup>, and Zhencheng Fang<sup>1\*</sup>

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

\*Correspondence address. Hongwei Zhou, Microbiome Medicine Center, Department of Laboratory Medicine, Zhujiang Hospital, Southern Medical University, 253 Gongye Middle Avenue, Haizhu District, Guangzhou, Guangdong 510280, China. E-mail: [hzhou@smu.edu.cn](mailto:hzhou@smu.edu.cn); Zhencheng Fang, Microbiome Medicine Center, Department of Laboratory Medicine, Zhujiang Hospital, Southern Medical University, 253 Gongye Middle Avenue, Haizhu District, Guangzhou, Guangdong 510280, China. E-mail: [fangzc@smu.edu.cn](mailto:fangzc@smu.edu.cn)

## Abstract

**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.

**Results:** We developed MOBFinder, an 11-class classifier, for categorizing plasmid fragments into 10 MOB types and a nonmobilizable 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 2 diabetes (T2D), MOBFinder offered insights suggesting that the MOBF type plasmid, which is widely present in *Escherichia* and *Klebsiella*, and the MOBQ type plasmid might accelerate antibiotic resistance transmission in patients with T2D.

**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>.

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

## Introduction

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

Plasmid classification is important for investigating multiple properties of plasmids, such as host range, replication patterns, and mobilization mechanisms [4]. Many classification schemes have been developed according to the distinct characteristics of plasmids, including taxonomic classification, replicon typing (Rep), incompatibility typing (Inc), mate-pair formation typing (MPF), and mobilization typing (MOB). In taxonomic classification, plasmids are categorized based on their host bacteria [5]. Rep typing classifies plasmids according to genes controlling their replication, known as replication initiation genes [4, 6]. Inc typing takes advantage of the fact that plasmids with similar replication or partition systems are incompatible within the same cell, categorizing plasmids based on compatibility [6]. MPF typing is based on genes encoding the MPF system, which consists of proteins that

mediate contact and DNA exchange between donor and recipient cells during conjugation [4, 7]. Finally, MOB typing classifies plasmids based on the relaxase gene, which is present in all transmissible plasmids [2, 8, 9]. Plasmids with different relaxase types are categorized as different MOB types, with each possessing a distinct transmission mechanism that determines its taxonomic host range [4, 10]. This variation among different MOB types is critical in researching the spread of virulence traits, the emergence of antibiotic resistance, and the adaptation and evolution of bacteria. Moreover, MOB typing has been found to be effective for identifying novel mobilizable plasmids that were previously unassigned to any Rep or Inc types and for investigating the mobilization characteristics of plasmids with similar mobilization systems [9, 11].

Recently, many experimental and computational schemes have been devised for plasmid typing, as well as to explore the diversity and functionality of plasmids (Table 1). For example, plasmid taxonomic PCR (PlasTax-PCR) [12], PCR-based replicon typing (PBRT) [13], and degenerate primer MOB typing (DPMT) [11] are multiplex PCR methods for identifying plasmids with analogous replication or mobilization systems. PlasTrans, based on deep learning, identifies mobilizable metagenomic plasmid fragments [14]. Web servers such as PlasmidFinder [6], pMLST, and oriTfinder [15] were established based on collected maker gene databases and

Received: February 28, 2024. Revised: May 31, 2024. Accepted: June 24, 2024

© The Author(s) 2024. Published by Oxford University Press on behalf of GigaScience. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

**Table 1:** Experimental and computational schemes developed for plasmid classification

Technology category	Method	Classification scheme	Material	Description
Experimental	DPMT [11]	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 [12]	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 [13]	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 [16]	MOB typing	Plasmid protein sequences	Used the HMMER model to annotate the relaxases and further perform MOB typing
	MOB-suite [17, 18]	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 the database, then classified plasmid assembly clusters with BLAST
	PlasTans [14]	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 [15]	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

alignment-based methods to facilitate Rep, Inc, and MOB typing. COPLA [5], based on average nucleotide identity, performs taxonomic classifications of complete plasmid genomes with an overall accuracy of 41%. For the MOB typing, MOBscan [16] uses the HMMER model to annotate relaxase genes and classify plasmids accordingly. MOB-suite [17, 18] performs plasmid typing for plasmid assemblies. First, it uses Mash distance to cluster plasmid assemblies into clusters; then, it uses marker gene databases to annotate them.

Metagenomic sequencing makes it possible to obtain all plasmid DNA from microbial communities at once, and a number

of computational tools for identifying plasmid fragments from metagenomic data have been developed, such as PlasFlow [19], PlasmidSeeker [20], PlasClass [21], PPR-Meta [22], and PlasForest [23]. As DNA fragments of plasmids and bacteria are intermingled in metagenomic data [24], recognizing the transmission mechanisms and host ranges of plasmids can be challenging. To this end, it is crucial to annotate MOB types of metagenomic plasmid fragments. However, this is difficult when plasmid assembly fragments are incomplete and essential genes for annotation are lacking. Therefore, it is worthwhile to consider alternative methods. Given that plasmids of the same MOB type have similar

transmission mechanisms and host ranges, their genomic signatures (e.g., GC content and codon usage) tend to also be alike, not only relaxase [4, 25]. In this context, neural networks, which have demonstrated strong performance in the classification and identification of biological sequences [26, 27], could be useful. Furthermore, language models [28, 29] derived from such neural networks have also showcased their impressive ability to characterize sequence features [30, 31]. In this methodology, short sequences of nucleotides (referred to as *k*-mers) or amino acids are analogous to “words,” and the longer sequences of DNA or proteins are analogous to “sentences.” Through the application of unsupervised learning on large datasets, each “word” is linked to a feature vector that captures its context, offering a more sophisticated analysis than the traditional *k*-mer frequency method, which simply counts the occurrence of nucleotide sequences without acknowledging their biochemical characteristics. Unlike the conventional method, this language model-based approach assesses sequences based on their contextual importance across different genetic environments, positioning contextually similar sequences close together in a multidimensional space. This technique provides deeper insights into the biochemical complexities of nucleotide sequences, thereby furnishing a more comprehensive understanding of an organism’s functional biology [32]. To characterize the features of plasmids within the same MOB type, we employed language models to perform the MOB annotation. In addition to the relaxase-coding gene, language models exhibit the ability to capture more biological features and associations within comparable mobilization systems, making it possible to perform MOB annotation for metagenomic plasmid assemblies.

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

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

## Materials and Methods

### The workflow of MOBFinder

To annotate the MOB type of plasmid fragments in metagenomics, we designed MOBFinder (Fig. 1). As MOB-suite [17, 18]

did not offer a quantitative likelihood score for the outcomes and some plasmids would be classified into multiple MOB types (Supplementary Fig. S1), we constructed a benchmark dataset using a high-resolution MOB typing strategy for categorizing complete plasmid genomes (Fig. 1B, C). Then, based on a language model and random forest, we designed an algorithm to perform MOB typing for PMFs (Fig. 1D, E).

### MOB typing of complete plasmid genomes

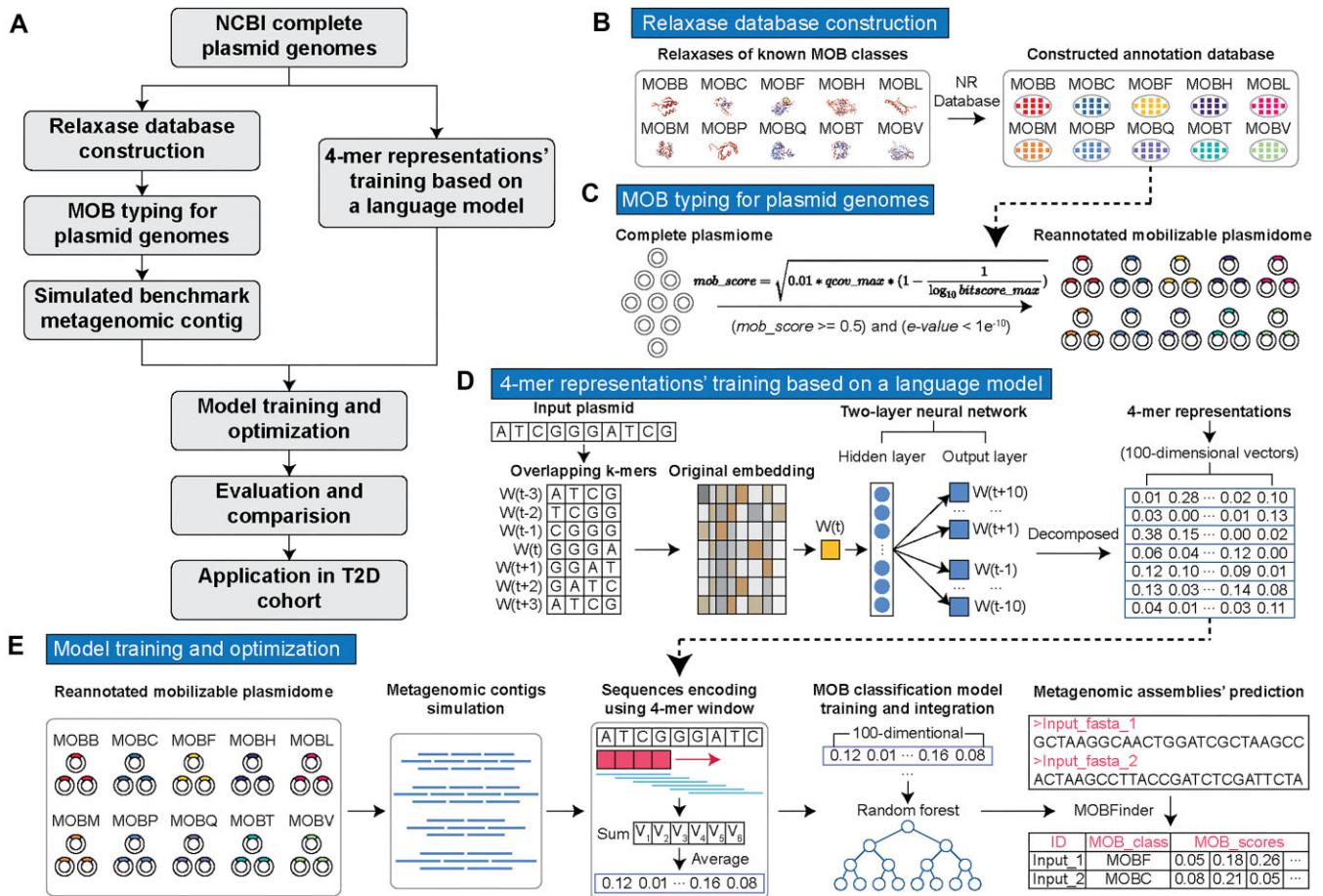
Traditionally, plasmid MOB typing of complete plasmid genomes has been a bioinformatics task based on the analysis of relaxase sequence similarity. The practice of annotating MOB types through BLAST similarity searches using representative sequences of different MOB type relaxases has gradually evolved into the standard method for MOB typing [4, 17, 18]. In this work, we constructed a benchmark dataset of simulated metagenomic contigs based on complete plasmid genomes with known MOB types. Previous studies have included a relatively small number of plasmids in their analyses. To further expand the MOB typing training dataset, we annotated the newly collected plasmid complete genome data for MOB typing according to relaxase information.

Ten validated MOB relaxase protein families were collected, including MOBB, MOBC, MOBF, MOBH, MOBL, MOBM, MOBP, MOBQ, MOBT, and MOBV [2, 7–9, 33, 34] (Fig. 1B). For each MOB category, blastp (RRID:SCR\_001010) [35] was used to search homologous protein sequences against the NCBI nonredundant protein sequence database, with an *e*-value threshold of 1e-10, a *query coverage* threshold of 70%, and an *identity* threshold of 70%. A previous study applied an *e*-value threshold of 1e-5 and minimum requirements for *query coverage* and *identity* set at 50% [4]. However, employing these criteria, we observed that some relaxases were annotated as belonging to multiple MOB types. To eliminate ambiguous annotations and construct a more reliable dataset for the training of MOBFinder, we imposed the stricter criteria mentioned above. After the expansion of protein sequences, local relaxase databases were built using the “makeblastdb” command for MOB typing of plasmid genomes.

Plasmid genomes were retrieved from the NCBI nucleotide database using the keywords “complete” and “plasmid,” and incomplete fragments were removed manually for further analysis. The accession list of these plasmids is provided in Supplementary Table S1. For each plasmid genome, coding sequences were extracted from the genebank file, and blastp [35] was employed to search for the best alignment of local relaxase databases. Here, we defined the *mob\_score* to measure the likelihood of homology:

$$mob\_score = \sqrt{0.01 * qcov\_max * (1 - 1/\log_{10}(bitscore\_max))}$$

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



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

## Word embeddings using a language model

To characterize the features and patterns within each MOB category and use numerical word vectors to represent them, we utilized a skip-gram language model [28, 29] to learn from plasmid genomes. Using a sliding window, the model calculated the likelihood between segmented words and outputted a probability distribution over the context words. The training steps were as follows (Fig. 1D):

- (1) Word generation. Since DNA sequences are composed of different nucleotide characters, we used a  $k$ -mer sliding window to generate overlapping input words. For example, with  $k = 4$ , "ATCGCTGA" would be segmented into "ATCG," "TCGC," "CGCT," "GCTG," and "CTGA." In this step, unique words were generated.
- (2) Word encoding initialization. Each word was initially assigned a random vector.
- (3) Skip-gram model. We employed a standard skip-gram model as described in previous studies [28, 29] to generate word vectors through the dna2vec module [29]. A 2-layer neural network was used to construct the skip-gram model. The initialized vectors were used as input, and the output was

a probability distribution over the input words. Layer 1 was a hidden layer to convert the initialized vectors into a 100-dimensional word vector representation as predefined by Ng [29]. Layer 2 was used to compute and maximize the probability of the correct context words using the negative sampling function, with the size of context words set to 20 (10 words for upstream and downstream, respectively) as preset by Ng [29].

- (4) Model training. For each input plasmid genome, we used an optimization algorithm to minimize the loss function. Then, using the default settings, we used backpropagation to update the neural network parameters (word vectors) for 10 epochs.
- (5) Word vector extraction. After the training process, the word vectors in the hidden layer were extracted to characterize the plasmid fragments.

## Benchmark dataset construction

Because there are no real metagenomic data to serve as a benchmark, using simulated data as a benchmark dataset is a common approach when developing bioinformatics tools [14, 22].

Therefore, in the development of MOBFinder, we artificially generated simulated datasets through the following steps:

- (1) For classified plasmid genomes in each MOB category, we randomly split them at a proportion of 70% and 30% to construct the training and test datasets.
- (2) Training dataset. To predict plasmid fragments with different lengths, we generated contigs of different length ranges: 100–400 bp, 401–800 bp, 801–1,200 bp, and 1,201–1,600 bp. For each MOB class in each length range, we randomly generated 90,000 artificial contigs. Plasmid fragments longer than 1,600 bp were segmented into shorter contigs and predicted using models designed for the corresponding lengths.
- (3) Test dataset. Because some plasmid fragments in real metagenomics datasets were much longer, we generated 4 length groups to assess the performance of MOBFinder: group A with a length range of 801–1,200 bp, group B with a length range of 1,201–1,600 bp, group C with a length range of 3,000–4,000 bp, and group D with a length range of 5,000–10,000 bp. For each MOB class in these 4 groups, 500 fragments were randomly extracted.

## Classification algorithm

To efficiently handle the training dataset and improve the robustness of MOBFinder, we employed random forest to train 4 predictive models using the training dataset. The detailed steps are as follows (Fig. 1E):

- (1) Word representation calculation. For each contig in the training dataset, we used a 4-mer sliding window to generate overlapping words and transformed them into numerical word vectors using trained word embeddings. To characterize the underlying features and patterns of the input contigs, we summed all the word vectors to compute their average as input of random forest.
- (2) Classification model training. To improve the performance of MOBFinder, we trained 4 classification models on different lengths in the training dataset: 100–400 bp, 401–800 bp, 801–1,200 bp, and 1,201–1,600 bp. The number of trees was set to 500 to generate predictive models.
- (3) Model ensemble. The 4 trained models were ensembled into MOBFinder to make more accurate predictions. For fragments shorter than 100 bp, we used a model designed for 100–400 bp to predict the MOB type. For those longer than 1,600 bp, we segmented them into short contigs and made predictions using the corresponding model. For example, a fragment with a length of 4,000 bp would be segmented into 3 contigs: 2 with a length of 1,600 bp and 1 of 800 bp. After predicting fragments with the corresponding models, we aggregated and calculated the weighted average scores for each MOB class, and the MOB type with the highest score was selected as the final prediction result for the input fragment.
- (4) Plasmid bin classification. Metagenomic binning is an essential step in the reconstruction of genomes from individual microorganisms. Thus, we designed MOBFinder to perform MOB typing on both plasmid contigs and plasmid bins. If the input is a plasmid bin, MOBFinder predicts the likelihood of each MOB class for fragments within the bin. For each MOB category, MOBFinder aggregates the scores of each sequence within the bin and calculates the weighted average scores based on the sequence length. The MOB category with the maximum score is selected as the prediction result.

## Performance validation

A test dataset was used to assess the performance of MOBFinder and compare it to MOB-suite and MOBscan. Because MOBscan can only predict MOB type using plasmid protein sequences rather than DNA sequences, we first annotated the proteins in the plasmid fragments of the test set using Prokka (RRID:SCR\_014732) [36] and then used MOBscan to predict the MOB type based on the annotated proteins. We calculated overall *accuracy*, *kappa*, and *run time* by comparing the predicted classes and true classes. We used the online server of MOBscan to perform the MOB annotation, and the calculation of *run time* for MOBscan was confined to the duration spent on preprocessing with Prokka locally. The overall *accuracy* was the proportion of accurate predictions. The *kappa* (a) was calculated to assess the overall consistency between the predictions and true classes, which took into account the possibility of random prediction.  $P_o$  represented observed accuracy [ $P_o = (A_{11} + A_{22} + \dots + A_{nn})/N$ ], where  $A_{11}$ ,  $A_{22}$ , and  $A_{nn}$  represented the values on the diagonal of the confusion matrix and  $n$  represented the number of MOB categories.  $N$  represented the total number of samples.  $P_e$  represented the expected accuracy [ $P_e = (E_{11} + E_{22} + \dots + E_{nn})/N^2$ ], where  $E_{11}$ ,  $E_{22}$ , and  $E_{nn}$  were the expected values in each cell of the confusion matrix;  $n$  was the number of MOB classes; and  $N$  was the total number of samples. The *run time* was recorded using the command “time” in Linux.

$$kappa = (P_o - P_e) / (1 - P_e) \quad (a)$$

$$balanced\ accuracy = (TPR + TNR) / 2 \quad (b)$$

$$harmonic\ mean = 2 * S_n * S_p / (S_n + S_p) \quad (c)$$

$$F1\ score = 2 * precision * recall / (precision + recall) \quad (d)$$

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

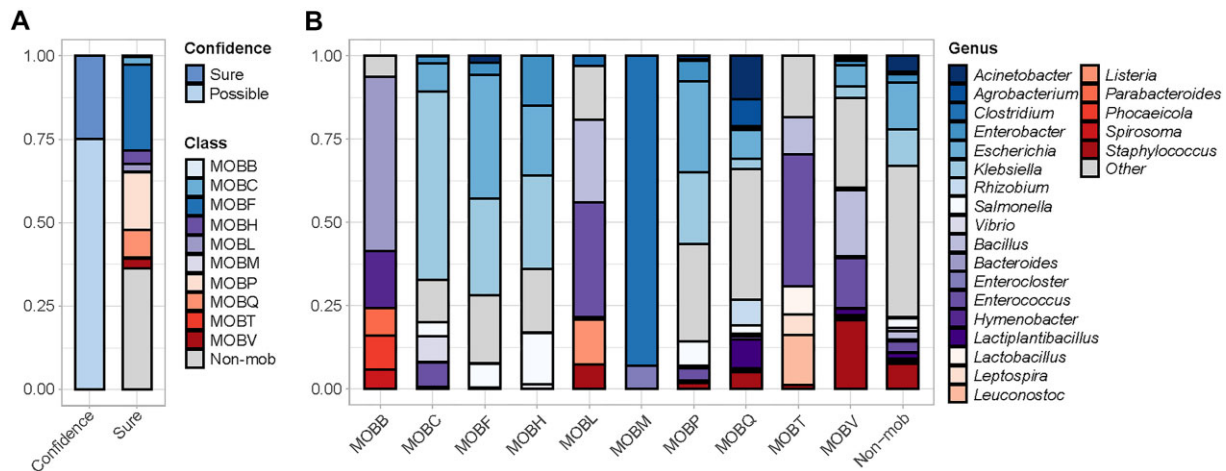
A receiver operating characteristic (ROC) curve was used to visualize the performance of MOBFinder in predicting each MOB category, where the x-axis and y-axis were the false-positive rate (FPR) and true-positive rate (TPR). Plots closer to the left and top indicate higher TPR and lower FPR, which means better performance. For each MOB class, the area under the curve (AUC) value was calculated to quantify the performance of MOBFinder. An AUC value between 0.5 and 1 indicates that the model performs better than random chance, and a higher AUC value indicates better prediction capability.

## Annotation and analysis of T2D metagenomic data

Metagenomic sequencing data (SRA045646 and SRA050230) were retrieved from the NCBI Short Read Archive (SRA) database to investigate whether the plasmids within different MOB classes

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

Class	Number	Average length	GC (%)
MOBB	623	10,921.77	51.27
MOBC	3,218	19,965.28	47.14
MOBF	21,268	103,802.80	52.07
MOBH	4,880	151,108.10	48.37
MOBL	3,446	51,430.63	34.57
MOBM	1,761	2,684.14	27.12
MOBP	15,617	32,237.88	49.70
MOBQ	9,347	89,357.64	56.77
MOBT	1,181	11,643.24	36.92
MOBV	4,405	6,595.43	37.75
Non-MOB	24,649	37,581.85	49.84



**Figure 2:** Benchmark dataset construction using a high-resolution strategy. (A) Proportion of classified plasmid genomes. A confidence level of “sure” means that the classified plasmid genomes had a *mob\_score* of more than 0.5 and an *e-value* of less than  $1e-10$ , while “possible” did not. Plasmid genomes identified as “sure” were used to generate benchmark datasets. Non-MOB, nonmobilizable plasmid. (B) Host range of the classified plasmid genomes at the genus level. Different colors represent different genera, and genera accounting for less than 5% of the total abundance are grouped under the category “other.”

were associated with antibiotic resistance enrichment in T2D patients, as suggested by previous studies [37, 38]. All metagenomic data were preprocessed using the same protocols. PRINSEQ (RRID:SCR\_005454) [39] was used to remove low-quality reads and bowtie2 (RRID:SCR\_016368) [40] was used to remove host reads by aligning them to the human GRCH38 reference genome downloaded from the ENSEMBL database. We excluded metagenomic samples that did not pass quality control. Because the abundance of plasmids in metagenomes was much lower than that of bacteria, we only retained samples with more than 10,000,000 paired-end reads for downstream analysis (Supplementary Table S2).

To improve the efficiency and accuracy of assembly, we used MEGAHIT (RRID:SCR\_018551) [41] to generate metagenomic contigs. PPR-Meta (RRID:SCR\_016915) [22] was utilized to identify and extract plasmid fragments from the assembled fragments while filtering out bacteria and phage sequences. COCACOLA [42] was employed to cluster plasmid fragments into bins based on sequence similarity and composition. This allowed us to investigate the plasmid fragments from same originate and enabled better annotation and analysis of their functions.

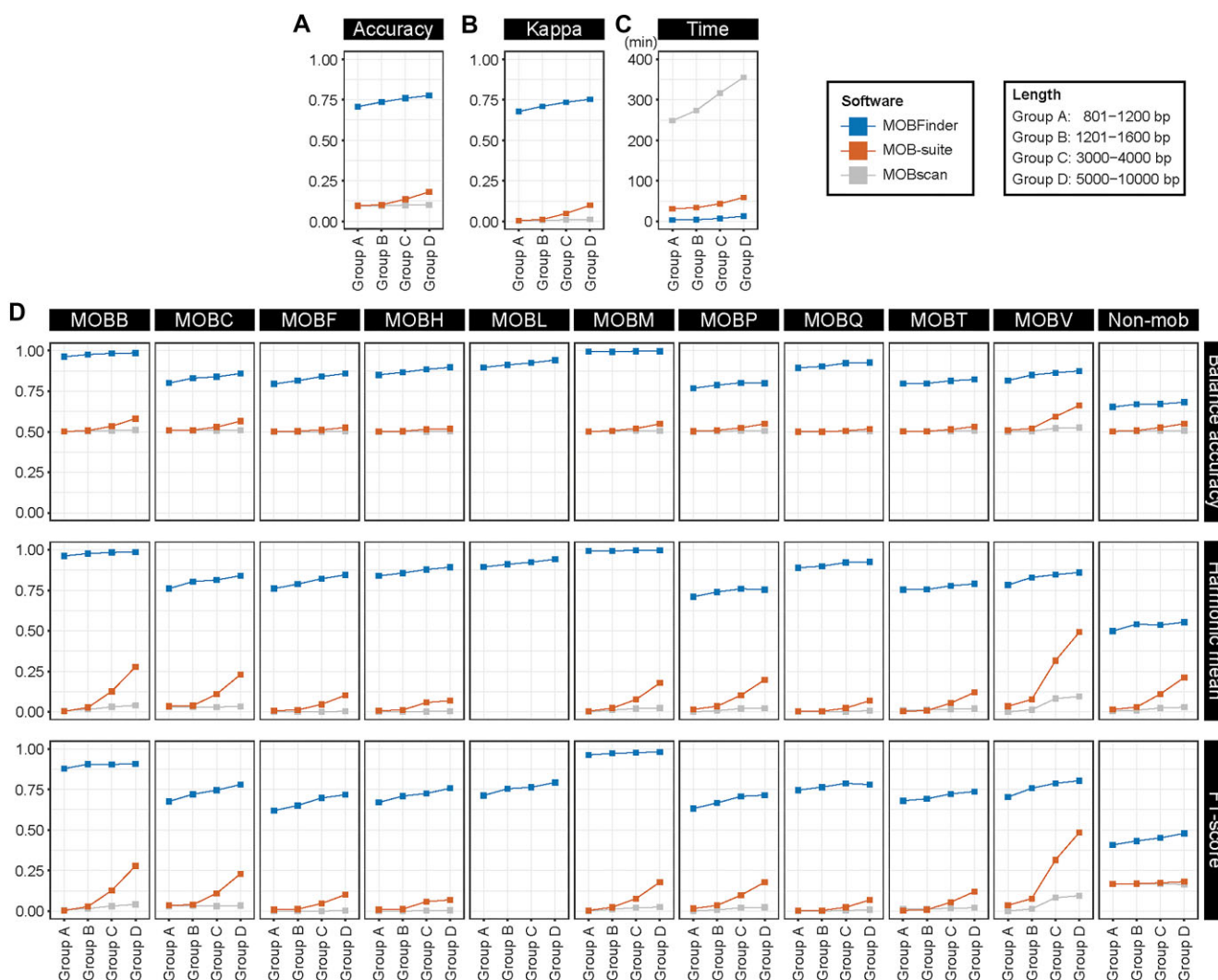
MOBFinder was applied to annotate the MOB types in each plasmid bin. The average fragments per kilobase per million of each plasmid bin was calculated using bowtie2 to represent its

abundance. Next, we analyzed the significance of differences in plasmid bins and various MOB types between healthy and T2D groups using the Wilcoxon rank-sum test. The calculation of *P* values was adjusted for multiple comparisons using the Benjamini–Hochberg method (denoted as *p.adjust*). ABRicate (RRID:SCR\_021093) [43] was utilized to annotate antibiotic resistance genes (*identity* >50% and *qcov* >50%) in each plasmid bin, based on 4 antibiotic resistance gene databases [44–47]. The Tukey honest significant difference test was performed to compare the identified resistance genes among different MOB classes. All statistical analyses were conducted using R.

## Results

### MOB typing of plasmid genomes

To construct the benchmark datasets, we obtained 90,395 complete plasmid genomes and categorized them into 11 MOB categories using blast (Table 2). We removed 22,470 of them potentially classified into more than 1 MOB class, leaving 67,925 classified genomes for the training and optimization of MOBFinder (Fig. 2A). Our analysis results revealed significant differences in the number, average length, and GC content of plasmid genomes among MOB types. Notably, non-MOB types included the genomes



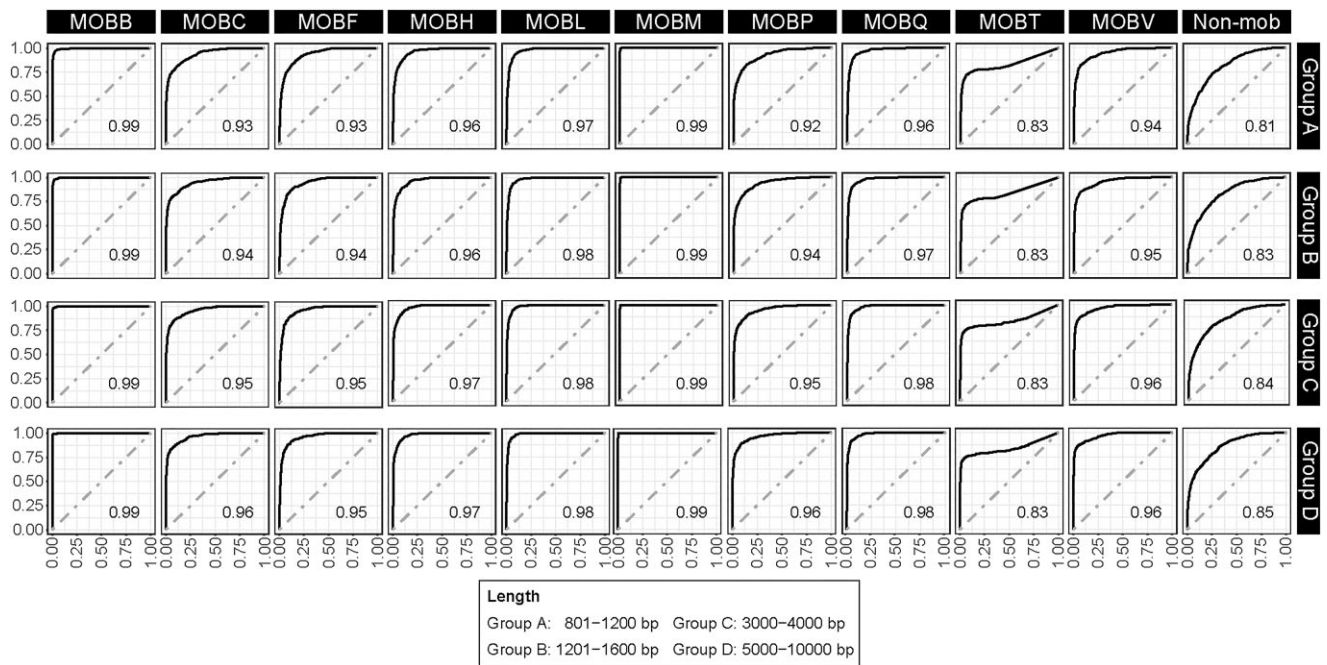
**Figure 3:** Overall performance of MOBFinder and comparison to MOB-suite and MOBScan. Evaluation and comparison in terms of (A) accuracy, (B) kappa, and (C) run time (min). The 4 fragment length groups in the test dataset were group A (801–1,200 bp), group B (1,201–1,600 bp), group C (3,000–4,000 bp), and group D (5,000–10,000bp). (D) For each MOB type, the balanced accuracy, harmonic mean, and F1-score were used to assess the performance of MOBFinder and compared to MOB-suite and MOBScan. Since MOB-suite and MOBScan do not include the prediction of MOBL, only the results of MOBL from MOBFinder are provided. MOBFinder, MOB-suite, and MOBScan are represented by blue lines, orange lines, and gray lines, respectively.

with the most and longest average length, whereas MOBB and MOBM had the fewest plasmid genomes and shortest average length, respectively. In terms of GC content, MOBL had the lowest and MOBQ had the highest amounts. Moreover, plasmids of different MOB types exhibited diverse host ranges at the genus level (Fig. 2B). MOBB was predominantly found in *Bacteroides*, *Hymenobacter*, *Parabacteroides*, *Phocaeicola*, and *Spirosoma*. Particularly, *Phocaeicola* has been detected in the human gut and possessed the gene for porphyrin degradation through horizontal gene transfer [48]. MOBC, MOBF, MOBH, and MOBP were all found in *Escherichia* and *Klebsiella*. *Klebsiella* also is a multidrug-resistant bacterium that has demonstrated resistance to multiple antibiotics [49]. MOBL, MOBT, and MOBV were mainly discovered in *Bacillus* and *Enterococcus*. Almost all MOBM-type plasmid genomes were present in *Clostridium* and *Enterocloster*, and some species in *Clostridium* could cause various diseases [50]. MOBQ demonstrated a broader host range, including *Acinetobacter*, *Agrobacterium*, *Escherichia*, *Rhizobium*, *Lactiplantibacillus*, and *Staphylococcus*. Non-MOB plasmids were detected in the majority of bacteria. These results illustrate the relationship between different MOB types and their host

ranges, as well as demonstrate that MOB typing of plasmid fragments is feasible in the absence of relaxases.

### Overall performance of MOBFinder

We evaluated the overall performance of MOBFinder in terms of accuracy, kappa, and run time and compared the tool to MOBScan and MOB-suite. 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 (Fig. 3A, B). In comparison, the accuracy of MOBFinder ranged from 70% to 77%, a significant improvement of at least 59% over MOB-suite (Fig. 3A). The kappa of MOBFinder ranged between 67% and 75% and was approximately 65% higher than that of MOB-suite (Fig. 3B). Moreover, MOBFinder exhibited a shorter run time in the test dataset, with a more gradual increase trend (Fig. 3C). In general, these results indicate that MOBFinder greatly outperformed the other tools and consistently improved in accuracy and consistency as the sequence length increased.



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

### Evaluation by MOB category

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

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

### Application to T2D metagenomic data

In a previous study, enrichment analysis of fecal samples identified antibiotic resistance pathways in patients with T2D [38]. The precise mechanism of this enrichment, however, remained elusive. We used MOBFinder to analyze real T2D metagenomic data [37]. After preprocessing and assembly, 2,217,064 metagenomic fragments were generated, and plasmid assemblies were identified using PPR-Meta. Subsequently, the plasmid fragments were clustered into 55 bins and annotated using MOBFinder. By employing MOBFinder, we assigned 2 bins to the MOBF class, 8 bins to MOBL, and 17 bins to MOBQ and identified 28 bins as non-MOB (Fig. 5A). Furthermore, we detected 15 bins that exhib-

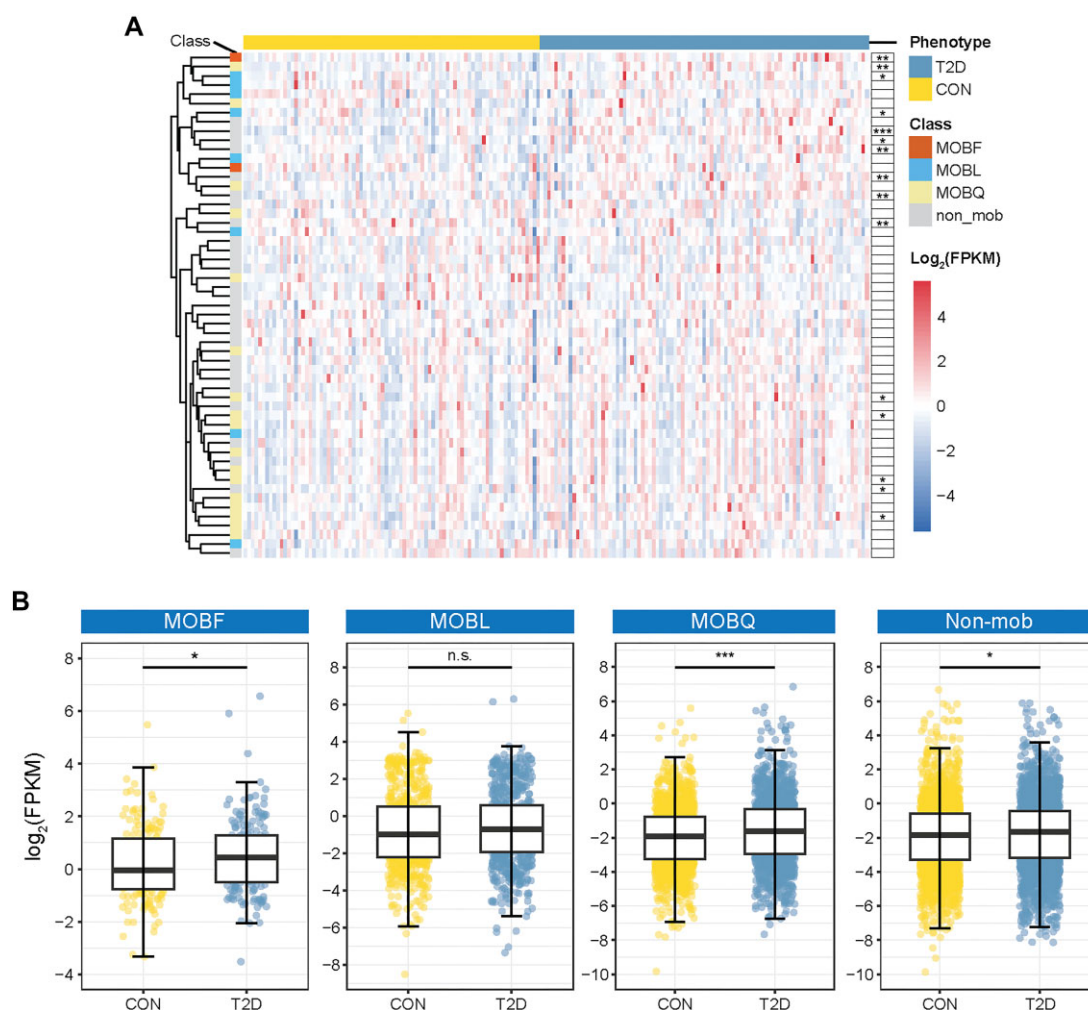
ited significant differences between the T2D group and a control group. Among them, 1 bin was classified as MOBF, 2 as MOBL, 5 as MOBQ, and 7 as non-MOB (Supplementary Fig. S2). Among above MOB types, MOBQ contains the highest number of bins enriched in T2D, while MOBF is widely present in *Escherichia* and *Klebsiella* (Fig. 2B); some strains of *Klebsiella* are resistant to multiple antibiotics, including carbapenems [51], and these 2 MOB types might contribute to antibiotic resistance in T2D patients. Indeed, when we compared the average abundance of each MOB type between the T2D group and the control group (Fig. 5B), the abundances of MOBF and MOBQ were significantly greater in the T2D group.

In addition, these 2 MOB types can be transferred among multiple bacterial species. This suggests that an increase in these 2 MOB types could potentially raise the risk of bacterial infection among individuals with T2D. Subsequently, we used 4 databases [44–47] to detect drug resistance genes in 4 MOB types (Fig. 6). The number of such genes was significantly higher in MOBF than in the other 3 MOB types. This suggests that MOBF plasmids may carry more drug resistance genes than the other MOB types. Furthermore, the increase in MOBF and MOBQ plasmids could result in more bacteria acquiring drug resistance genes, thereby leading to more antibiotic resistance pathways in T2D patients. In summary, our results demonstrate the utility of MOBFinder for annotating plasmid fragments in metagenomes, uncovering the potential mechanisms underlying the antibiotic resistance enrichment in metagenomic analysis.

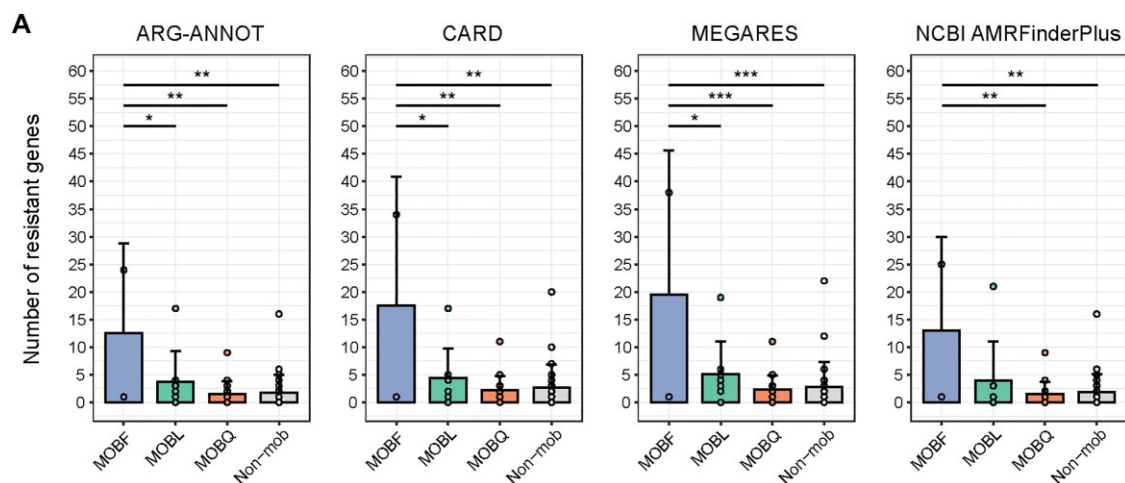
### Use of MOBFinder

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





**Figure 5:** Annotation of T2D-related plasmid bins using MOBFinder. (A) Heatmap of plasmid bins between T2D patients and controls. Each column represents a sample, and each row represents a plasmid bin. (B) Comparison of the abundance of the 4 identified MOB types between T2D patients and controls. The P value was calculated using the Wilcoxon rank-sum test, adjusted using the Benjamini-Hochberg method for multiple comparisons (\**P*-adjust < 0.05, \*\**P*-adjust < 0.01, and \*\*\**P*-adjust < 0.001).



**Figure 6:** Comparison of resistance genes among different MOB types. Four databases were used to identify antibiotic resistance gene within each MOB type, and the P value was calculated using the Tukey honest significant difference test. The 2 groups without significance markings indicate no statistical difference (\**P*-adjust < 0.05, \*\**P*-adjust < 0.01, and \*\*\**P*-adjust < 0.001).

For plasmid metagenomic bins, MOBFinder requires 2 input files: a FASTA file containing the plasmid fragments and a metatable that records the mapping between plasmid fragment IDs and bin IDs. The output results are similar to those of plasmid fragments. The first column is the plasmid bin ID. The second is the predicted MOB class of the plasmid bins. The other columns present the MOB scores of the different MOB types.

## Discussion

We developed MOBFinder based on a language model and the random forest algorithm to classify plasmid fragments and bins from metagenomics data into MOB types. First, using the relaxe-alignment method, plasmid genomes were classified into distinct MOB categories. Analyses revealed substantial differences in parameters such as the number, average length, and GC content of plasmid genomes across MOB types. Additionally, there were noteworthy differences in the host ranges among different MOB classes. These results suggest the potential of utilizing sequence features from different MOB types for PMF MOB typing. To characterize the plasmids within each MOB type, we used the skip-gram model to generate word vectors. Our tool demonstrated superior overall performance compared to other tools. Specifically, for each MOB category, MOBFinder exhibited significant improvements in *balanced accuracy*, *harmonic mean*, and *F1-score*, with values reaching up to 99% for the first 2 measures in the MOBM category.

Traditionally, *k*-mer frequency models and one-hot encoding have commonly been employed to digitize biological sequences, extensively applied across various machine learning algorithms [52]. However, both models simply mark or count the frequency of various characters in sequences, failing to reflect the biological significance underlying each character. These models may also encounter dimensionality issues [52]. For instance, in the *k*-mer model, if *k* is set to 8, the dimensionality of the *k*-mer vector of each DNA sequence becomes  $4^8$ , which is problematic in metagenomics where most fragment lengths do not reach this magnitude. This would result in significant noise in the feature vector and cause overfitting. Similarly, in the one-hot model, for a sequence of length *L* using 4-mers as the base unit, it would require *L* one-hot vectors each with a dimensionality of  $4^4$ . In such instances, if the dataset for training is not sufficiently large, this representation method could also lead to overfitting due to high dimensionality. In contrast, word vector models offer a superior solution to these problems. Such models initially perform a random initialization of vectors for each “word.” Taking the skip-gram algorithm utilized in this study as an example, the dimension of a random vector can be 1-of-*n*, where *n* represents the size of the vocabulary [28]. Following unsupervised pretraining on large datasets, the algorithm maps characters with similar contexts to similar feature spaces. The dimensions of the coordinates (i.e., the word vectors) of these feature spaces will be lower than those of the initial random vectors. Thus, through unsupervised pretraining on large datasets, language models can compress high-dimensional initial vectors into lower-dimensional word vectors (e.g., MOBFinder’s word vectors have a dimensionality of 100), enabling the feature vectors to contain more character information while effectively avoiding dimensionality issues during supervised training.

In a metagenomic sequences classification task, 4-mer is widely used as the basic unit in various bioinformatics tools [53], and thus MOBFinder takes this as a “word.” To assess the impact of training word vectors with different *k*-mer lengths on perfor-

mance, we compared models with *k*-mer lengths of 2, 3, 4, 5, 6, 7, and 8 (Supplementary Fig. S3). We observed lower overall accuracy and *kappa* values for *k* = 2. At *k* = 4, the *balanced accuracy*, *harmonic mean*, *F1-score*, and *AUC* values stabilized across different MOB types. Subsequently, as the *k*-mer length increased, there was no significant improvement in accuracy or other metrics, while the *run time* gradually increased. Therefore, we chose a *k*-mer length of 4 for training word vectors and developing MOBFinder.

Interestingly, in an analysis of T2D metagenomic sequencing data [37], we noted a significant increase in MOB and MOBQ type plasmids in T2D patients. Moreover, we found more drug resistance genes in the MOB class, whose dominant hosts are *Klebsiella* and *Escherichia*, which are associated with the spread of multidrug resistance. Although previous analyses of gut metagenomic data from patients with T2D have reported enrichment of drug resistance pathways [38], our results suggest a potential reason for it: the increased abundance of MOB- and MOBQ-type plasmids in the guts of individuals with T2D may disseminate more antibiotic resistance genes, resulting in such enrichment.

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

## Conclusions

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

## Availability of Source Code and Requirements

- Project name: MOBFinder
- Project homepage: <https://github.com/FengTaoSMU/MOBFinder>
- Operating system(s): Linux
- Programming language: Python, R script
- Other requirements: BLAST, biopython
- License: GPL-3.0
- RRID:SCR\_024451
- biotoolsID: MOBFinder

## Additional Files

**Supplementary Table S1.** Accessions list of classified plasmid genomes.

**Supplementary Table S2.** List of metagenomic samples used in our analysis.

**Supplementary Fig. S1.** MOB typing using MOB-suite. Single-class, plasmid genomes classified into 1 MOB type; multiclass, plasmid genomes classified into more than 1 MOB category; non-MOB, nonmobilizable plasmids.

**Supplementary Fig. S2.** Abundance of each significantly different plasmid bin from various MOB types between patients with type 2 diabetes and controls.

**Supplementary Fig. S3.** Comparison results for the development of MOBFinder using word vectors trained with different *k*-mer lengths. (A–C) Overall accuracy, *kappa*, and run time of the MOB classification model trained with word vectors trained using different lengths of *k*-mers. (D) Balanced accuracy, harmonic mean, F1-score, and AUC of word vectors trained with different *k*-mer lengths across different MOB types.

## Abbreviations

AUC: the area under the curve; BLAST: Basic Local Alignment Search Tool; DPMT: degenerate primer MOB typing; Inc: incompatibility typing; MOB: mobilization typing; MPF: mate-pair formation typing; NCBI: National Center for Biotechnology Information; non-MOB: nonmobilizable; PBRT: PCR-based replicon typing; PlasTax-PCR: plasmid taxonomic PCR; Rep: replicon typing; ROC: receiver operating characteristic; Sn: sensitivity; Sp: specificity; SRA: Short Read Archive; T2D: type 2 diabetes; TNR: true-negative rate; TPR: true-positive rate.

## Author Contributions

T.F., Z.C.F., and H.W.Z. proposed and designed this work. T.F. and Z.C.F. developed and optimized the software. T.F., Z.C.F., S.F.W., and H.W.Z. wrote and revised the manuscript.

## Funding

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

## Data Availability

Snapshots of our code and other data further supporting this work are openly available in the GigaScience repository, GigaDB [54].

## Competing Interests

The authors declare that they have no competing interests.

## References

- Helinski DR. A brief history of plasmids. *EcoSal Plus* 2022;10(1):eESP00282021. <https://doi.org/10.1128/ecosalplus.esp-0028-2021>.
- Garcillán-Barcia MP, Francia MV, de la Cruz F. The diversity of conjugative relaxases and its application in plasmid classification. *FEMS Microbiol Rev* 2009;33(3):657–87. <https://doi.org/10.1111/j.1574-6976.2009.00168.x>.
- Rodríguez-Beltrán J, DelaFuente J, León-Sampedro R, et al. Beyond horizontal gene transfer: the role of plasmids in bacterial evolution. *Nat Rev Micro* 2021;19(6):347–59. <https://doi.org/10.1038/s41579-020-00497-1>.
- Shintani M, Sanchez ZK, Kimbara K. Genomics of microbial plasmids: classification and identification based on replication and transfer systems and host taxonomy. *Front Microbiol* 2015;6:242. <https://doi.org/10.3389/fmicb.2015.00242>.
- Redondo-Salvo S, Bartomeus-Peñalver R, Vielva L, et al. COPLA, a taxonomic classifier of plasmids. *BMC Bioinf* 2021;22(1):390. <https://doi.org/10.1186/s12859-021-04299-x>.
- Carattoli A, Hasman H. PlasmidFinder and In Silico pMLST: identification and typing of plasmid replicons in whole-genome sequencing (WGS). *Methods Mol Biol* 2020;2075:285–94. [https://doi.org/10.1007/978-1-4939-9877-7\\_20](https://doi.org/10.1007/978-1-4939-9877-7_20).
- Smillie C, Garcillán-Barcia MP, Francia MV, et al. Mobility of plasmids. *Microbiol Mol Biol Rev* 2010;74(3):434–52. <https://doi.org/10.1128/MMBR.00020-10>.
- Francia MV, Varsaki A, Garcillán-Barcia MP, et al. A classification scheme for mobilization regions of bacterial plasmids. *FEMS Microbiol Rev* 2004;28(1):79–100. <https://doi.org/10.1016/j.femsre.2003.09.001>.
- Garcillán-Barcia MP, Alvarado A, de la Cruz F. Identification of bacterial plasmids based on mobility and plasmid population biology. *FEMS Microbiol Rev* 2011;35(5):936–56. <https://doi.org/10.1111/j.1574-6976.2011.00291.x>.
- Bradley P, den Bakker HC, Rocha EPC, et al. Ultrafast search of all deposited bacterial and viral genomic data. *Nat Biotechnol* 2019;37(2):152–59. <https://doi.org/10.1038/s41587-018-0010-1>.
- Alvarado A, Garcillán-Barcia MP, de la Cruz F. A degenerate primer MOB typing (DPMT) method to classify gamma-proteobacterial plasmids in clinical and environmental settings. *PLoS One* 2012;7(7):e40438. <https://doi.org/10.1371/journal.pone.0040438>.
- Cuartas R, Coque TM, de la Cruz F, et al. PLASmid TAXonomic PCR (PlasTax-PCR), a multiplex relaxase MOB typing to assort plasmids into taxonomic units. *Methods Mol Biol* 2022;2392:127–42. [https://doi.org/10.1007/978-1-0716-1799-1\\_10](https://doi.org/10.1007/978-1-0716-1799-1_10).
- Carattoli A, Bertini A, Villa L, et al. Identification of plasmids by PCR-based replicon typing. *J Microbiol Methods* 2005;63(3):219–28. <https://doi.org/10.1016/j.mimet.2005.03.018>.
- Fang Z, Zhou H. Identification of the conjugative and mobilizable plasmid fragments in the plasmidome using sequence signatures. *Microb Genom* 2020;6(11):mgen000459. <https://doi.org/10.1099/mgen.0.000459>.
- Li X, Xie Y, Liu M, et al. oriTfinder: a web-based tool for the identification of origin of transfers in DNA sequences of bacterial mobile genetic elements. *Nucleic Acids Res* 2018;46(W1):W229–34. <https://doi.org/10.1093/nar/gky352>.
- Garcillán-Barcia MP, Redondo-Salvo S, Vielva L, et al. MOB-scan: automated annotation of MOB relaxases. *Methods Mol Biol* 2020;2075:295–308. [https://doi.org/10.1007/978-1-4939-9877-7\\_21](https://doi.org/10.1007/978-1-4939-9877-7_21).
- Robertson J, Nash JHE. MOB-suite: software tools for clustering, reconstruction and typing of plasmids from draft assemblies. *Microb Genom* 2018;4(8):e000206. <https://doi.org/10.1099/mgen.0.000206>.
- Robertson J, Bessonov K, Schonfeld J, et al. Universal whole-sequence-based plasmid typing and its utility to prediction of host range and epidemiological surveillance. *Microb Genom* 2020;6(10):mgen000435. <https://doi.org/10.1099/mgen.0.000435>.

19. Krawczyk PS, Lipinski L, Dziembowski A. PlasFlow: predicting plasmid sequences in metagenomic data using genome signatures. *Nucleic Acids Res* 2018;46(6):e35. <https://doi.org/10.1093/nar/gkx1321>.
20. Roosaare M, Puustusmaa M, Möls M, et al. PlasmidSeeker: identification of known plasmids from bacterial whole genome sequencing reads. *PeerJ* 2018;6:e4588. <https://doi.org/10.7717/peerj.4588>.
21. Pellow D, Mizrahi I, Shamir R. PlasClass improves plasmid sequence classification. *PLoS Comput Biol* 2020;16:e1007781. <https://doi.org/10.1371/journal.pcbi.1007781>.
22. Fang Z, Tan J, Wu S, et al. PPR-Meta: a tool for identifying phages and plasmids from metagenomic fragments using deep learning. *Gigascience* 2019;8:1–14. <https://doi.org/10.1093/gigascience/giz066>.
23. Pradier L, Tissot T, Fiston-Lavier AS, et al. PlasForest: a homology-based random forest classifier for plasmid detection in genomic datasets. *BMC Bioinf* 2021;22(1):349. <https://doi.org/10.1186/s12859-021-04270-w>.
24. Sobecky PA, Hazen TH. Horizontal gene transfer and mobile genetic elements in marine systems. *Methods Mol Biol* 2009;532:435–53. [https://doi.org/10.1007/978-1-60327-853-9\\_25](https://doi.org/10.1007/978-1-60327-853-9_25).
25. Suzuki H, Yano H, Brown CJ, et al. Predicting plasmid promiscuity based on genomic signature. *J Bacteriol* 2010;192(22):6045–55. <https://doi.org/10.1128/JB.00277-10>.
26. Wu S, Fang Z, Tan J, et al. DeePhage: distinguishing virulent and temperate phage-derived sequences in metavirome data with a deep learning approach. *Gigascience* 2021;10(9):giab056. <https://doi.org/10.1093/gigascience/giab056>.
27. Fang Z, Feng T, Zhou H, et al. DeePVP: identification and classification of phage virion proteins using deep learning. *Gigascience* 2022;11:giac076. <https://doi.org/10.1093/gigascience/giac076>.
28. Mikolov T, Chen K, Corrado G, et al. Efficient estimation of word representations in vector space. *arXiv preprint* 2013. <https://doi.org/10.1371/journal.pone.0141287>. Accessed 16 November 2021.
29. Ng P. dna2vec: consistent vector representations of variable-length k-mers. *arXiv* 2017. <https://doi.org/10.48550/arXiv.1701.06279>. Accessed 18 November 2021.
30. Tsukiyama S, Hasan MM, Fujii S, et al. LSTM-PHV: prediction of human-virus protein-protein interactions by LSTM with word2vec. *Brief Bioinform* 2021;22(6):bbab228. <https://doi.org/10.1093/bib/bbab228>.
31. Sharma R, Shrivastava S, Kumar Singh S, et al. Deep-ABPpred: identifying antibacterial peptides in protein sequences using bidirectional LSTM with word2vec. *Brief Bioinform* 2021;22(5):bbab065. <https://doi.org/10.1093/bib/bbab065>.
32. Asgari E, Mofrad MR. Continuous distributed representation of biological sequences for deep proteomics and genomics. *PLoS One* 2015;10(11):e0141287. <https://doi.org/10.1371/journal.pone.0141287>.
33. Wisniewski JA, Traore DA, Bannam TL, et al. TcpM: a novel relaxase that mediates transfer of large conjugative plasmids from *Clostridium perfringens*. *Mol Microbiol* 2016;99(5):884–96. <https://doi.org/10.1111/mmi.13270>.
34. Ramachandran G, Miguel-Arribas A, Abia D, et al. Discovery of a new family of relaxases in Firmicutes bacteria. *PLoS Genet* 2017;13(2):e1006586. <https://doi.org/10.1371/journal.pgen.1006586>.
35. Camacho C, Coulouris G, Avagyan V, et al. BLAST+: architecture and applications. *BMC Bioinf* 2009;10:421. <https://doi.org/10.1186/1471-2105-10-421>.
36. Seemann T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* 2014;30(14):2068–69. <https://doi.org/10.1093/bioinformatics/btu153>.
37. Qin J, Li Y, Cai Z, et al. A metagenome-wide association study of gut microbiota in type 2 diabetes. *Nature* 2012;490(7418):55–60. <https://doi.org/10.1038/nature11450>.
38. Wu H, Tremaroli V, Schmidt C, et al. The gut microbiota in prediabetes and diabetes: a population-based cross-sectional study. *Cell Metab* 2020;32(3):379–90. e3. <https://doi.org/10.1016/j.cmet.2020.06.011>.
39. Schmieder R, Edwards R. Quality control and preprocessing of metagenomic datasets. *Bioinformatics* 2011;27(6):863–64. <https://doi.org/10.1093/bioinformatics/btr026>.
40. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Methods* 2012;9(4):357–59. <https://doi.org/10.1038/nmeth.1923>.
41. Li D, Liu CM, Luo R, et al. MEGAHIT: an ultra-fast single-node solution for large and complex metagenomics assembly via succinct de Bruijn graph. *Bioinformatics* 2015;31(10):1674–76. <https://doi.org/10.1093/bioinformatics/btv033>.
42. Lu YY, Chen T, Fuhrman JA, et al. COCACOLA: binning metagenomic contigs using sequence COMposition, read Coverage, CO-alignment and paired-end read LinkAge. *Bioinformatics* 2017;33(6):791–98. <https://doi.org/10.1093/bioinformatics/btw290>.
43. Seemann T. Abricate. GitHub 2020. <https://github.com/tseeman/abricate>. Accessed 20 June 2022.
44. Gupta SK, Padmanabhan BR, Diene SM, et al. ARG-ANNOT, a new bioinformatic tool to discover antibiotic resistance genes in bacterial genomes. *Antimicrob Agents Chemother* 2014;58(1):212–20. <https://doi.org/10.1128/AAC.01310-13>.
45. Jia B, Raphenya AR, Alcock B, et al. 2017: expansion and model-centric curation of the comprehensive antibiotic resistance database. *Nucleic Acids Res* 2017;45(D1):D566–73. <https://doi.org/10.1093/nar/gkw1004>.
46. Doster E, Lakin SM, Dean CJ, et al. MEGARes 2.0: a database for classification of antimicrobial drug, biocide and metal resistance determinants in metagenomic sequence data. *Nucleic Acids Res* 2020;48(D1):D561–69. <https://doi.org/10.1093/nar/gkz1010>.
47. Feldgarden M, Brover V, Haft DH, et al. Validating the AMRFinder tool and resistance gene database by using antimicrobial resistance genotype-phenotype correlations in a collection of isolates. *Antimicrob Agents Chemother* 2019;63(11):e00483–19. <https://doi.org/10.1128/AAC.00483-19>.
48. Hehemann JH, Correc G, Barbeyron T, et al. Transfer of carbohydrate-active enzymes from marine bacteria to Japanese gut microbiota. *Nature* 2010;464(7290):908–12. <https://doi.org/10.1038/nature08937>.
49. Fu S, Wang R, Xu Z, et al. Metagenomic sequencing combined with flow cytometry facilitated a novel microbial risk assessment framework for bacterial pathogens in municipal wastewater without cultivation. *iMeta* 2023;2:e77. <https://doi.org/10.1002/imt2.77>.
50. Dieterle MG, Rao K, Young VB. Novel therapies and preventative strategies for primary and recurrent *Clostridium difficile* infections. *Ann N Y Acad Sci* 2019;1435(1):110–38. <https://doi.org/10.1111/nyas.13958>.
51. Yang X, Dong N, Chan EW, et al. Carbapenem resistance-encoding and virulence-encoding conjugative plasmids in *Klebsiella pneumoniae*. *Trends Microbiol* 2021;29(1):65–83. <https://doi.org/10.1016/j.tim.2020.04.012>.

52. Jaillard M, Palmieri M, van Belkum A, et al. Interpreting k-mer-based signatures for antibiotic resistance prediction. *GigaScience* 2020;9(10):giaa110. <https://doi.org/10.1093/gigascience/giaa110>.
53. Sedlar K, Kupkova K, Provaznik I. Bioinformatics strategies for taxonomy independent binning and visualization of sequences in shotgun metagenomics. *Comput Struct Biotechnol J* 2016;15:48–55. <https://doi.org/10.1016/j.csbj.2016.11.005>.
54. Feng T, Wu S, Zhou H, et al. Supporting data for “MOBFinder: A Tool for Mobilization Typing of Plasmid Metagenomic Fragments Based on a Language Model.” *GigaScience Database*. 2024. <https://doi.org/10.5524/102559>.