


# 1 plASgraph - using graph neural networks to detect 2 plasmid contigs from an assembly graph

3 **Janik Sielemann** 


4 Computational Biology, Faculty of Biology, Center for Biotechnology (CeBiTec) & Graduate School  
5 DILS, Bielefeld Institute for Bioinformatics Infrastructure (BIBI), Bielefeld University, 33615  
6 Bielefeld, Germany  
7 Department of Mathematics, Simon Fraser University, Burnaby, Canada  
8 [janik.sielemann@uni-bielefeld.de](mailto:janik.sielemann@uni-bielefeld.de)

9 **Katharina Sielemann** 

10 Genetics and Genomics of Plants, Faculty of Biology, Center for Biotechnology (CeBiTec) &  
11 Graduate School DILS, Bielefeld Institute for Bioinformatics Infrastructure (BIBI), Bielefeld  
12 University, 33615 Bielefeld, Germany  
13 Department of Mathematics, Simon Fraser University, Burnaby, Canada  
14 [ksielemann@cebitec.uni-bielefeld.de](mailto:ksielemann@cebitec.uni-bielefeld.de)

15 **Broňa Brejová** 

16 Department of Computer Science, Faculty of Mathematics, Physics and Informatics, Comenius  
17 University in Bratislava, Slovakia  
18 Department of Mathematics, Simon Fraser University, Burnaby, Canada  
19 [brejova@dcs.fmph.uniba.sk](mailto:brejova@dcs.fmph.uniba.sk)

20 **Tomáš Vinar** 

21 Department of Applied Informatics, Faculty of Mathematics, Physics and Informatics, Comenius  
22 University in Bratislava, Slovakia  
23 Department of Mathematics, Simon Fraser University, Burnaby, Canada  
24 [tomas.vinar@fmph.uniba.sk](mailto:tomas.vinar@fmph.uniba.sk)

25 **Cedric Chauve** 

26 Department of Mathematics, Simon Fraser University, Burnaby, Canada  
27 [cedric.chauve@sfu.ca](mailto:cedric.chauve@sfu.ca)

## 28 — Abstract —

---

29 Identification of plasmids from sequencing data is an important and challenging problem related  
30 to antimicrobial resistance spread and other One-Health issues. In our work, we provide a new  
31 architecture for identifying plasmid contigs in fragmented genome assemblies built from short-read  
32 data. Unlike previous machine-learning approaches for this problem, which classify individual contigs  
33 separately, we employ graph neural networks (GNNs) to include information from the assembly  
34 graph. Propagation of information from nearby nodes in the graph allows accurate classification of  
35 even short contigs that are difficult to classify based on sequence features or database searches alone.

36 Our new species-agnostic software tool plASgraph outperforms recently developed PlasForest,  
37 which uses database searches to supplement sequence-based features. Since our tool does not rely  
38 on existing plasmid databases, it is more suitable for classification of contigs in novel species and  
39 discovery of previously unknown plasmid sequences. Our tool can also be trained on a specific  
40 species, and in that scenario it outperforms mlplasmids trained on the same species.

41 On one hand, our work provides a new, accurate, and easy to use tool for plasmid classification;  
42 on the other hand, it serves as a motivation for more widespread use of GNNs in bioinformatics, such  
43 as in pangenome sequence analysis, where sequence graphs serve as a fundamental data structure.

44 **Availability:** <https://github.com/cchauve/plASgraph>

45 **2012 ACM Subject Classification** Applied computing - Life and medical sciences - Computational  
46 biology - Computational genomics

47 **Keywords and phrases** Contig classification, graph neural network, machine learning, plasmids

48 **Funding** *Janik Sielemann*: Bielefeld University, Graduate School DILS (Digital Infrastructure for

## 2 pIASgraph

49 the Life Sciences); EU Horizon 2020 grant No. 872539 (PANGAIA)  
50 *Katharina Sielemann*: Bielefeld University, Graduate School DILS (Digital Infrastructure for the  
51 Life Sciences); EU Horizon 2020 grant No. 872539 (PANGAIA)  
52 *Broňa Brejová*: VEGA grant 1/0463/20; EU Horizon 2020 grant No. 872539 (PANGAIA)  
53 *Tomáš Vinař*: VEGA grant 1/0538/22; EU Horizon 2020 grant No. 872539 (PANGAIA)

54 **Acknowledgements** This research was enabled in part by support provided by Compute Canada  
55 ([www.computecanada.ca](http://www.computecanada.ca)). We thank Aniket Mane for introducing the idea of using GNN with  
56 assembly graphs. Most of the work was conducted during a visit of JS, KS, BB and TV to Simon  
57 Fraser University enabled by the PANGAIA EU project.

## 58 1 Introduction

59 Plasmids are mobile genetic elements that are involved in horizontal gene transfers and  
60 have been shown to be a major vector for the spread of antimicrobial resistance (AMR)  
61 genes [8, 20]. Plasmids are extra-chromosomal DNA molecules, often circular and significantly  
62 shorter than bacterial chromosomes, and can occur in multiple copies in a bacterial cell.  
63 Whereas some bacteria do not contain any plasmid, it is common to observe several plasmids  
64 co-existing within a bacterial cell, often with different copy numbers. Due to their high  
65 mobility and impact in AMR spread, the detection of plasmids from sequencing data is an  
66 important question in One-Health epidemiologic surveillance approaches, see e.g. [10].

67 Given sequencing data, either from a bacterial isolate or from a metagenome, the detection  
68 of plasmids can be approached at various levels of detail. The most elementary task, *contig*  
69 *classification*, aims at detecting which assembled contigs are likely to originate from a plasmid.  
70 *Plasmid binning* aims at grouping contigs into groups likely to originate from the same  
71 plasmid. Last, *plasmid assembly* aims at reconstructing full plasmid sequences. While  
72 obtaining full plasmids provides the most accurate information, the ability to extract plasmid  
73 contigs from assembled sequencing data (the contig classification problem) already provides  
74 very useful information, allowing e.g. to identify genes that might be susceptible to transfer  
75 to other bacteria. Moreover, the prediction of plasmid contigs can be used as an input for  
76 plasmid binning or assembly. For example, the plasmid binning method *gplas* [4] relies on a  
77 preliminary contig classification obtained with *mlplasmids* [5] and the metagenome plasmid  
78 assembly method *SCAPP* [22] relies on classifying contigs using *PlasClass* [21].

79 While the analysis of plasmids from sequencing data has been a very active research  
80 area, the problems mentioned above are still challenging, especially when sequencing data  
81 are provided in the form of Illumina short reads [6]. In the present paper, we propose a  
82 novel method for the contig classification problem, specifically designed to analyse short-read  
83 contigs from a single bacterial isolate.

84 There exists a large corpus of algorithms for the contig classification problem, most of  
85 them developed recently. These methods rely mainly on machine-learning approaches. The  
86 earliest method for contig classification was *cBar* [30], which introduced the use of the *k*-mer  
87 profile of a contig as the main feature in a machine-learning classification model; in *cBar*,  
88 the model was trained on a large dataset of closed bacterial genome assemblies. The general  
89 principle of using *k*-mer properties as classification features has also been used in several  
90 recent machine-learning classifiers, namely *PlasFlow* [16], *mlplasmids* [5], and *PlasClass* [21].  
91 *PPR-Meta* [11] is a deep-learning method that relies on one-hot encoded contig sequences.  
92 *PlasForest* [23] and *Deeplasmid* [3] are two recent methods based on machine-learning models  
93 that use different features for a given contig, such as its GC content (generally plasmids have  
94 a GC content different from chromosomes) and the presence of plasmid-specific sequences,

95 detected through the mapping against a reference plasmid database. RFPlasmid [27] combines  
96 both kinds of features, the  $k$ -mer profile and plasmid-specific sequences. Among the methods  
97 introduced above, both mlplasmids and RFPlasmid are species-specific methods, i.e. require  
98 a model to be trained per bacterial species; in contrast, PlasFlow, PlasClass, PlasForest and  
99 Deeplasmid are tools that do not target a specific species.

100 The recent method 3CAC [24] introduced the idea that the classification of a contig can be  
101 improved from the knowledge of the classification of the neighbouring contigs in the assembly  
102 graph. Most current assembly programs [7, 28] output an assembly graph containing final  
103 contigs as nodes and possible connections between them supported by sequencing data as  
104 edges. Individual molecules, such as chromosomes or plasmids, ideally correspond to walks  
105 in this graph, but some edges may be missing, disconnecting the walk. Conversely, the walks  
106 for individual molecules often form complicated tangled structures joined at shared and  
107 repeated sequences. Nonetheless, adjacent nodes often share the same molecule of origin and  
108 thus the same class. 3CAC applies simple heuristics to improve machine learning predictions  
109 for individual contigs based on their adjacency in the graph. Our aim is to integrate the  
110 information from the assembly graph directly to the underlying machine learning model.

111 Here, we introduce a novel machine-learning method, plASgraph, for the problem of  
112 classifying short-read contigs as plasmidic or chromosomal. Our method is based on combining  
113 features of existing methods with a novel approach incorporating a graph neural network  
114 (GNN) [12]. More precisely, plASgraph associates to each contig of a bacterial genome  
115 assembly a set of features that have been shown to differentiate plasmids and chromosomes:  
116 read coverage, used as a proxy of copy number, GC content and contig length, together  
117 with two novel features, the node degree in the assembly graph and the distance between  
118 the contig  $k$ -mer profile and the whole assembly  $k$ -mer profile. The rationale to integrate  
119 the  $k$ -mer profile by comparing it to the assembly-wide profile is to allow our model to be  
120 species-agnostic, i.e. not learning a species-specific  $k$ -mer profile, as is done in species-specific  
121 models such as mlplasmids and RFPlasmid. Moreover, plASgraph is a *de novo* tool that does  
122 not require the comparison of the input contigs with a database of known plasmids. Based  
123 on these features, plASgraph trains a GNN model whose core is a set of graph convolutional  
124 layers aimed at propagating the information from neighbouring contigs in the assembly  
125 graph. To the best of our knowledge, plASgraph is the first method that applies GNNs to  
126 contig classification in an assembly graph, building on the idea (introduced in 3CAC) that  
127 information from neighbouring contigs can improve accuracy. Outside of classification, GNNs  
128 were also used recently on assembly graphs for metagenomic contigs binning [17].

129 The output of plASgraph is a pair of scores for each graph node, a plasmid score and a  
130 chromosomal score, used to determine if a given contig is likely to originate from a plasmid  
131 or a chromosome or both. Unlike other methods, the two scores associated to a contig allow  
132 to detect *ambiguous* contigs that have shared sequences of both plasmidic and chromosomal  
133 origin. In order to train plASgraph, we rely on the availability of hybrid sequencing data  
134 composed of short and long reads; our method thus does not depend on the availability of a  
135 reference plasmid database.

136 We evaluated the performance of plASgraph in two contexts: species-specific and species-  
137 agnostic. In the species-specific context, we trained plASgraph on data from a single bacterial  
138 species and compared the model accuracy to mlplasmids [5] when applied to isolates from the  
139 same species and from other species. In the species-agnostic context, we trained plASgraph  
140 on data combined from several species and compared it to PlasForest [23], which in addition  
141 to the sequence-based features also uses the information from database searches. In both  
142 scenarios, plASgraph outperforms the competing tool.

## 4 pIASgraph

### 143 **2** Methods

#### 144 **2.1** Overview

145 The input to our problem is an assembly graph of a bacterial isolate in which nodes correspond  
146 to contigs and edges correspond to adjacencies supported by sequencing data. This graph is  
147 typically created from short reads and as a result can contain a large number of contigs of  
148 various sizes. For example, in our *E. faecium* training set (average genome size 2.84 Mbp),  
149 the number of contigs ranged from tens to hundreds with an N50 value between 34 kbp and  
150 253 kbp. Our goal is to classify individual contigs as originating from a plasmid or from a  
151 chromosome. However, some contigs in fact correspond to sequences that occur as parts of  
152 both plasmids and chromosomes within the same sample (for example, mobile elements or  
153 low-complexity sequences); we call such contigs *ambiguous*. Due to the presence of ambiguous  
154 contigs, we treat the problem as two separate classification tasks, generating independent  
155 scores for chromosome and plasmid labels.

156 Ideally, the input assembly graph would consist of a few connected components, each  
157 corresponding either to a single chromosome or a plasmid. In a fully resolved assembly,  
158 each chromosome and plasmid would actually correspond to an isolated vertex. However,  
159 in graphs created from short-read data, walks corresponding to individual molecules are  
160 typically interconnected by spurious edges or through ambiguous contigs, creating complex  
161 structures (see Figure 5 for an illustration). The main novelty of our approach is to use  
162 the graph neighbourhood of a node as a source of information, under the assumption that  
163 other contigs from the same molecule (chromosome or plasmid) are likely to belong to this  
164 neighbourhood.

165 Given a comprehensive database of closed genomes for a given bacterial species, longer  
166 contigs can likely be classified simply based sequence homology; in our work, we concentrate  
167 on classifying contigs in a *de novo* framework that does not rely on existing reference genomes.  
168 The *de novo* contigs classification problem is of interest for e.g. the analysis of samples from  
169 poorly sampled bacterial species.

#### 170 **2.2** Input features

171 As an input to the classification task, each contig is characterized by five input features:

- 172 1. the *degree* of the corresponding node in the assembly graph;
- 173 2. the *relative contig length*, defined by the contig length divided by the length of the longest  
174 contig;
- 175 3. the *relative GC content*, defined by subtracting the average GC content (expressed as a  
176 percentage) of the whole assembly from the contig GC content;
- 177 4. the *relative coverage*, defined as the contig read depth divided by the median read depth  
178 over the whole assembly (in our experiments, we use the read depth value provided by  
179 the Unicycler assembler [28]);
- 180 5. the *relative pentamer distance*, defined as the Euclidean distance between the pentamer  
181 profile of the contig and the pentamer profile of the whole assembly; we define the  
182 pentamer profile of a contig or a set of contigs as the count vector for all pentamers  
183 (reverse complements were aggregated), shifted and scaled so that the smallest vector  
184 entry is 0 and the largest vector entry is 1.

185 The motivation to rely on relative features instead of absolute features is to enable the  
186 model to generalize across species, and thus to not be dependent on species-specific values.  
187 Regarding the relative pentamer distance, one can expect that large chromosomal contigs

188 will have a value closer to zero, while shorter plasmid contigs will exhibit a large value for  
189 this feature. Moreover, by abstracting the pentamer content of a sample by the relative  
190 pentamer distance, we expect that our model will be less susceptible to learning to classify  
191 chromosome sequences by simply recognizing the pentamer frequencies characteristic for a  
192 particular species or a clade.

### 193 2.3 Model architecture

194 We employed a deep neural network to solve our classification task. The key part of our  
195 architecture is the use of a *graph convolutional network* (GNN) [15] to account for the  
196 assembly graph structure. The propagation of information between individual nodes is  
197 accomplished by *graph convolutional layers* (GCLs). Briefly, the input to a GCL contains  
198 a vector of  $k$  features for each of the  $n$  nodes of the graph and the adjacency matrix of  
199 the graph. The layer first combines the feature vectors corresponding to the node and its  
200 neighbours, with weight of nodes depending on their degree. The feature vector of each node  
201 is then transformed by a fully-connected layer with  $\ell$  output features followed by a non-linear  
202 activation. More precisely, if we organize the  $n$  feature vectors into the  $n \times k$  matrix  $X$ , the  
203 graph convolutional layer can be expressed as

$$204 \quad Z = \sigma(\tilde{D}^{-1/2} \tilde{A} \tilde{D}^{-1/2} X \Theta), \quad (1)$$

205 where  $\tilde{A}$  is the graph adjacency matrix with ones along the diagonal,  $\tilde{D}$  is a diagonal matrix  
206 where  $\tilde{D}_{ii} = \sum_j \tilde{A}_{ij}$ ,  $\Theta$  is a  $k \times \ell$  matrix of trainable weights,  $\sigma$  is a non-linear activation  
207 function, and  $Z$  is the  $n \times \ell$  matrix of output feature vectors. A single GCL integrates  
208 information from an immediate neighbourhood of a node; by employing  $d$  GCLs one integrates  
209 the information from the distance of at most  $d$  for each node.

210 Figure 1 shows the architecture we have designed for our task. The five input features for  
211 each node are first transformed by a fully connected layer to a vector of length 32 per node.  
212 This is followed by six GCLs using the same weight matrix  $\Theta$ . The last two fully connected  
213 layers operate on each node separately, the first producing a vector of length 32, and the  
214 second producing two output scores, loosely interpretable as probabilities of the node being  
215 part of a chromosome and plasmid, respectively. Since these two outputs correspond to two  
216 separate classification tasks, we do not require these two scores to sum to one. Each layer is  
217 followed by the ReLU activation, except for the last layer, which uses the sigmoid activation.  
218 All layers excluding the last one are followed by 10% dropout to prevent overfitting.

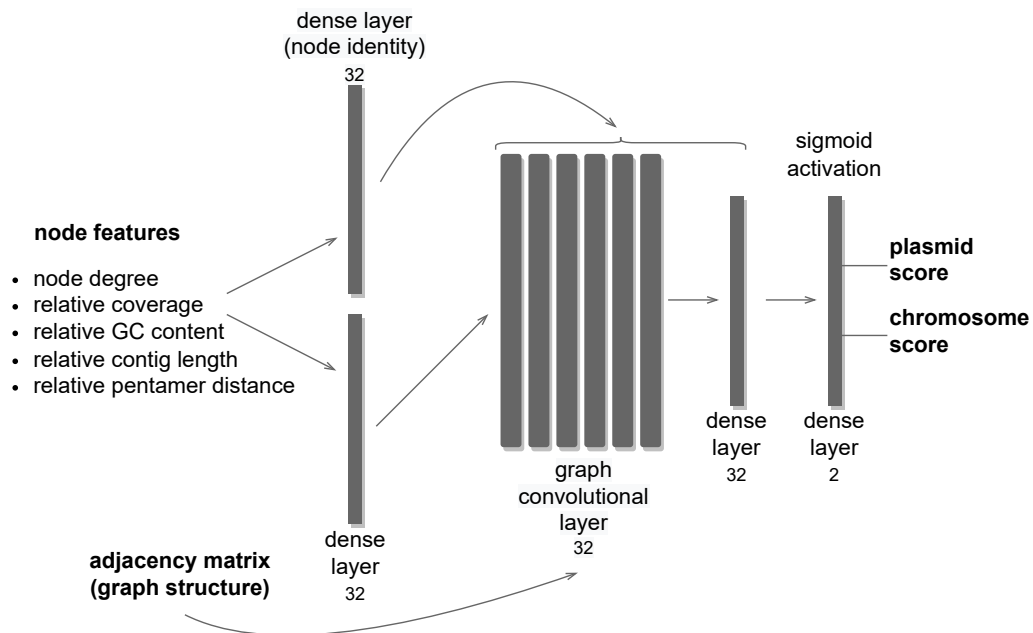
219 GCLs combine features of each node with features of the neighbours and over time, the  
220 influence of the original features is greatly diminished. In our task, the original features can  
221 be highly informative, especially for nodes corresponding to longer contigs; therefore we want  
222 to maintain node identity throughout the computation. To accomplish this, each GCL and  
223 the penultimate dense layer receive as an input an additional vector of length 32 for each  
224 node, representing a separate encoding of the five input features. Thus, each of these layers  
225 starts with the input vectors of length 64 and reduces them to vectors of length 32.

226 The network is trained using Adam optimizer [14] with binary cross entropy loss function  
227 and a constant learning rate of 0.005. The model is implemented using Keras [9] and  
228 TensorFlow v2.8.0 [1], with GCLs from Spektral v1.0.8 [12].

### 229 2.4 Data preparation

230 For training and testing, we used data from eight bacterial species, listed in Tab. 1. Individual  
231 bacterial isolates were sequenced both by short-read (Illumina) and long-read (Oxford

## 6 pIASgraph



■ **Figure 1 Model architecture of pIASgraph.** The model uses as inputs the assembly graph structure and five input features per node (contig). The core of the architecture is composed of six graph convolutional layers. The model generates two outputs per node, which facilitate the classification of plasmids and chromosomes as two separate classification tasks.

■ **Table 1 Datasets used in this study.** The second column lists the number of samples used for training or testing. Either the SRA project ID or the DOI are listed for data access.

Species	# samples	Data accession	Reference
<i>Bacillus megaterium</i>	2	PRJNA658106	[26]
<i>Citrobacter freundii</i>	88	PRJNA605147	[25]
<i>Escherichia coli</i>	166	PRJNA605147 PRJNA761884	[25] [2]
<i>Enterococcus faecium</i>	53	PRJEB28495 10.6084/m9.figshare.7046804 10.6084/m9.figshare.7047686	[5]
<i>Klebsiella oxytoca</i>	22	PRJNA605147	[25]
<i>Klebsiella pneumoniae</i>	45	PRJNA605147 PRJNA761884	[25] [2]
<i>Staphylococcus pseudintermedius</i>	3	PRJNA521119	[19]
<i>Vibrio parahaemolyticus</i>	3	PRJEB31923	[29]



232 Nanopore or Pacific Biosciences) technologies [2, 5, 19, 25, 26, 29]. We have followed the  
233 general methodology introduced by mlplasmids [5], that does not rely on databases of known  
234 plasmids and chromosomes (Fig. S1). By combining both short and long reads, we create a  
235 *hybrid assembly* with UniCycler [28], which typically has a small number of contigs, mostly  
236 corresponding to complete circular chromosomes or plasmids. A *short-read assembly* is then  
237 constructed from short reads only. The hybrid assembly is used to derive ground truth  
238 classification of short-read contigs, as explained in the next paragraph. The short-read  
239 assembly graph is then used as an input for our model, both in training and testing scenarios.

240 In hybrid assemblies, the ground truth labels are determined based on the contig length.  
241 In particular, all contigs longer than a species-specific threshold (Fig. S2) are labeled as  
242 'chromosome', while shorter circular contigs are labeled as 'plasmid'. The remaining short  
243 linear contigs can possibly be a part of an unfinished plasmid or chromosome, and consequently  
244 they remain unlabeled. The ground truth labels for short-read assemblies are determined  
245 by mapping the contigs to the corresponding hybrid assembly contigs, from which they  
246 inherit the labels. The key difference between our pipeline and mlplasmids is that if a contig  
247 matches equally well to both chromosomal and plasmidic contigs of the hybrid assembly, it is  
248 labeled as 'ambiguous'. Such contigs are considered as positive examples for both plasmid and  
249 chromosome classification tasks. We have observed that without considering such ambiguous  
250 matches, the assembly graphs of the short-read assemblies often contained paths with nodes  
251 labeled by alternating classes, which is a clearly inconsistent labeling. The use of ambiguous  
252 labels allows us to avoid such artefacts. Contigs that matched an unlabeled contig of the  
253 hybrid assembly were left unlabeled, and samples that contained more than 5% of unlabeled  
254 contigs were discarded from further analysis. The distribution of short read contig labels is  
255 shown in Fig. S3. In general, most contigs are labeled as chromosome and less than 1.3% of  
256 all contigs are left unlabeled.

257 Both hybrid and short-read assemblies were created by Unicycler v0.5.0 [28]. Short-read  
258 contigs were mapped to the hybrid assembly by minimap2 v2.24 [18] with -c option for  
259 accurate alignment.

260 For training and prediction, all contigs shorter than 100 bp were removed from the  
261 short-read assembly graphs and their neighbours were connected by direct edges as part of  
262 the feature extraction process. Thus, plASgraph is not predicting the class of contigs shorter  
263 than 100 bp. Note that such short contigs often do not have a reliable ground-truth label  
264 (often corresponding to SNPs and the like).

### 265 **3 Experimental evaluation**

266 In this section, we evaluate the performance of our plASgraph model and compare it to two  
267 recent tools mlplasmids [5] and PlasForest [23] (see Table S1 for overview of experiments).  
268 The scripts used for training and evaluation, as well as detailed results are available at  
269 [https://github.com/cchauve/plASgraph\\_WABI\\_2022](https://github.com/cchauve/plASgraph_WABI_2022).

270 Similarly to plASgraph, mlplasmids is a *de novo* tool, using only sequence-derived features  
271 to classify contigs. However, it requires training on each species separately and was designed  
272 for contigs longer than 1 kbp. In contrast, PlasForest is a species-agnostic tool, designed to  
273 work also on short contigs (< 1 kbp), but it is dependent on the comparison of the input  
274 sequences to sequence databases. Neither of these tools use the assembly graph information.

275 Since plASgraph was designed to explicitly handle ambiguous contigs by including separate  
276 plasmid and chromosomal classification tasks, we evaluate its prediction accuracy for each  
277 of these tasks separately. A contig is predicted as a chromosome if the chromosome score

## 8 pIASgraph

278 output of the neural network is at least 0.5; and similarly it is predicted as plasmid if the  
279 plasmid score is at least 0.5. The true and predicted contig labels then induced the counts  
280 of true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN)  
281 for each classification task. Each contig was counted as one unit, regardless of its length,  
282 and unlabeled contigs were not counted in the evaluation. Ambiguous contigs (those with  
283 score  $> 0.5$  in both classification tasks) are considered as being labeled both plasmid and  
284 chromosome in our accuracy evaluation. As the main accuracy measure, we use the F1 score,  
285 which is the harmonic mean of precision ( $TP/(TP + FP)$ ) and recall ( $TP/(TP + FN)$ ).

### 286 3.1 PIASgraph architecture leads to accurate species-specific models

287 Although our main goal is to produce a single model that can be used across species, we  
288 have also trained three species-specific models: *E. faecium* (46 training isolates), *E. coli* (66  
289 training isolates), and *K. pneumoniae* (35 training isolates). The isolates for training were  
290 chosen randomly from all available samples except for the *E. faecium*, where we used the  
291 same training set as mlpasmids. Out of the training samples, 20% were used as a validation  
292 set for the training procedure. The accuracy was evaluated and compared to mlpasmids  
293 on held-out testing data sets (*E. faecium* 7 isolates, *E. coli* 100 isolates, *K. pneumoniae* 10  
294 isolates). The development of the model architecture was solely performed on the *E. faecium*  
295 data set.

296 Figure 2 shows that the F1-score of pIASgraph is higher than mlpasmids on most samples  
297 from *E. faecium* and *K. pneumoniae*. Although mlpasmids outperforms pIASgraph on many  
298 *E. coli* samples, the overall score for plasmid classification is higher for pIASgraph, whereas  
299 the performance for the chromosome classification task is comparable (Tab. 2).

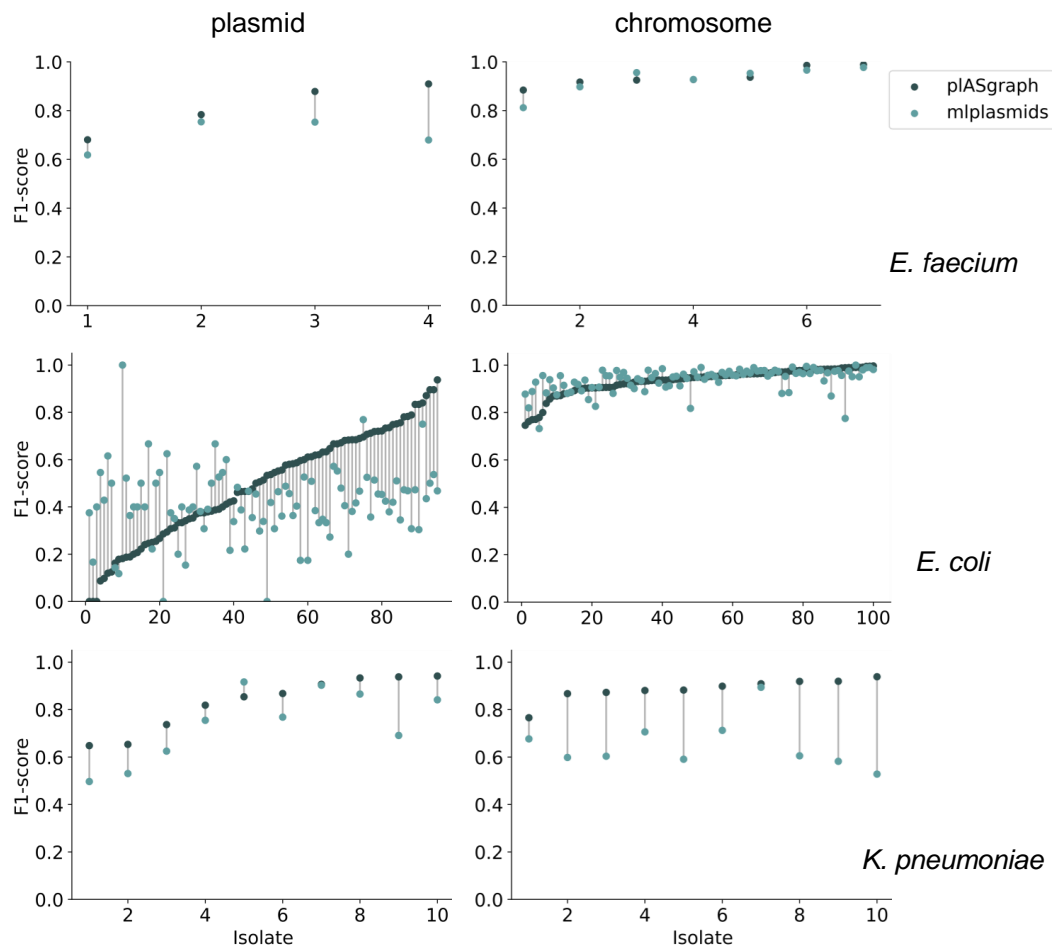
300 Our results also show that the main advantage gained by pIASgraph in the species-specific  
301 setting comes from classification of short contigs (100-1000 bp), since considering only contigs  
302 above 1 kbp, mlpasmids achieves better accuracy (Tab. 2). This observation can be explained  
303 by the  $k$ -mer frequency feature used by mlpasmids, which provides more detailed information  
304 for longer contigs. For contigs shorter than 1 kbp, the  $k$ -mer frequency feature may only  
305 contain counts for a few distinct  $k$ -mers. This can lead to a relatively high number of  
306 zero values in the feature vector which does not allow the model to classify the respective  
307 contig accurately. In contrast, pIASgraph uses only a single  $k$ -mer related feature, which is  
308 supplemented by information from graph neighbourhood; this combination likely helps to  
309 classify short contigs more accurately. Figure 3 demonstrates that indeed, the addition of  
310 GCLs significantly increases the prediction accuracy compared to a simpler model with the  
311 same input features.

### 312 3.2 Species-specific pIASgraph models generalize well to other species

313 Table 2 also shows that in cross-species application of species-specific models, pIASgraph has  
314 a distinct advantage compared to mlpasmids. For example, the pIASgraph model trained on  
315 *E. coli* achieves plasmid F1-score 0.76 on *K. pneumoniae*, which is only slightly lower than  
316 F1-score 0.83 achieved by the model trained on the same species. In the case of mlpasmids,  
317 we see a significant decrease of the F1-score, from 0.74 on the same species to 0.44 across  
318 species. Similar trends persist when we restrict evaluation to contigs above 1 kbp. Note  
319 the *E. faecium* models are less accurate on *E. coli* and *K. pneumoniae* testing sets due  
320 to the large phylogenetic distance, but pIASgraph still performs significantly better than  
321 mlpasmids. The distribution of per sample F1-scores is shown in Fig. S4.

322 We attribute this better cross-species generalization of pIASgraph to our use of input



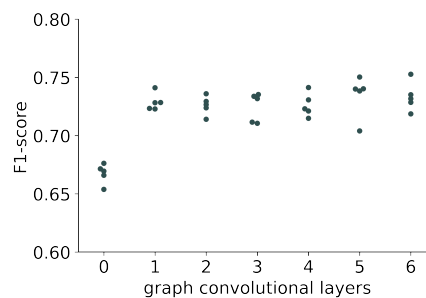


■ **Figure 2 Evaluation of the species-specific pIASgraph models in comparison to mplasmids.** Each row compares the F1-score of pIASgraph to mplasmids on individual testing isolates for both plasmid (left) and chromosome (right) classification tasks. The isolates without plasmids are not shown in graphs on the left, as F1-score is then undefined. Rows from top to bottom correspond to *E. faecium*, *E. coli*, and *K. pneumoniae* data sets respectively. The models are always trained and tested on the same species.

## 10 pIASgraph

■ **Table 2 Species-specific and cross-species evaluation.** Table shows the average F1-scores across isolates for all analyses of the species-specific models trained (rows) and tested (columns) on a specific combination of species. The cases where the model was trained and tested on the same species are highlighted by teal color. The numbers in parentheses show the results on contigs with length > 1 kbp. Plasmid and chromosome classification tasks are evaluated separately. Evaluation on *E. faecium* is not shown, since we have developed the model architecture on this dataset.

Model	Plasmid		Chromosome	
	<i>E. coli</i>	<i>K. pneumoniae</i>	<i>E. coli</i>	<i>K. pneumoniae</i>
<i>E. faecium</i> pIASgraph	0.299 (0.386)	0.469 (0.548)	0.939 (0.933)	0.788 (0.785)
<i>E. faecium</i> mlplasmids	0.139 (0.188)	0.296 (0.397)	0.860 (0.737)	0.665 (0.535)
<i>E. coli</i> pIASgraph	0.493 (0.566)	0.761 (0.824)	0.935 (0.928)	0.856 (0.909)
<i>E. coli</i> mlplasmids	0.415 (0.663)	0.436 (0.514)	0.939 (0.954)	0.719 (0.663)
<i>K. pneumoniae</i> pIASgraph	0.582 (0.657)	0.830 (0.851)	0.949 (0.944)	0.885 (0.928)
<i>K. pneumoniae</i> mlplasmids	0.305 (0.474)	0.739 (0.917)	0.430 (0.810)	0.650 (0.954)



■ **Figure 3 Impact of the number of layers on the model performance.** Different numbers of graph convolutional layers (GCLs) were tested for the *E. faecium* model. Each architecture was trained five times using random seeds and the F1-score was calculated on the 20% split validation set. The largest improvement in performance is visible with introduction of the first GCL; a slight additional increase in accuracy is observed for five and six GCLs.

323 features that are relative to the sample-wide statistics, and thus are less species dependent and  
324 are able to use the overall context provided by the whole assembly. In contrast, mlplasmids  
325 directly uses  $k$ -mer frequencies, which are highly specific to individual species.

### 326 3.3 Species-agnostic plASgraph model

327 One of the goals of the plASgraph model was to create a tool that could be applied to  
328 newly identified species for which no information is available in sequence databases and no  
329 training sets assembled with long reads are readily available. To this end, we have trained a  
330 species-agnostic model on a mixed training set from *E. faecium*, *E. coli*, and *K. pneumoniae*,  
331 using 20 isolates from each as a part of the training set. Again, 20% of the training set was  
332 withheld for validation during training. We evaluated the performance of the species-agnostic  
333 model on five species, different from those included in the training set. All available isolates  
334 for the five species (*B. megaterium* (2), *C. freundii* (88), *K. oxytoca* (22), *S. pseudintermedius*  
335 (3) and *V. parahaemolyticus* (3)) were used for testing. Since mlplasmids is not suitable for  
336 this application, we compared the plASgraph model with the recent PlasForest model [23]  
337 which has been designed for cross-species use.

338 Figure 4 shows that plASgraph has better average plasmid accuracy than PlasForest  
339 on four out of five species (see also Fig. S5 for the distribution of F1-scores for individual  
340 isolates). For the *S. pseudintermedius* data set, which contained only a single plasmid, both  
341 tools failed to predict that plasmid correctly. Both tools have high accuracy in predicting  
342 chromosomal contigs, with plASgraph being more accurate on *C. freundii*, *K. oxytoca*, and  
343 *B. megaterium*, but performing slightly worse for the remaining two species. The performance  
344 for ambiguous contigs is shown in Figure S6. Figure S7 shows scatter plot of all *C. freundii*  
345 contigs based on their chromosome and plasmid score and their true label. It shows that  
346 overall plASgraph classifies a large majority contigs accurately.

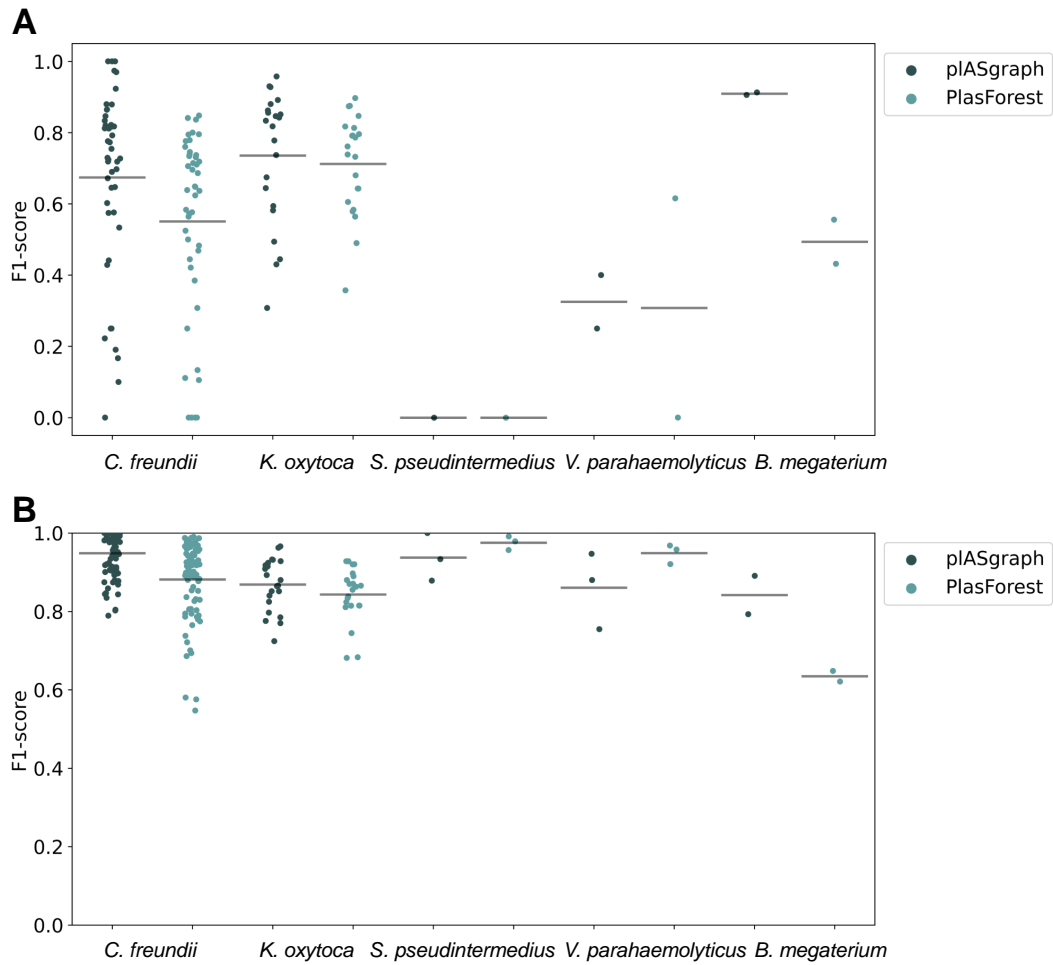
347 PlasForest uses features derived from querying input sequences against a reference  
348 database, which is not used in our predictions. The higher average F1-score of plASgraph  
349 therefore suggests that the contribution of information present in the assembly graph  
350 combined with relative features of the contigs exceeds the contribution from homology search.  
351 Moreover, independence from sequence databases makes our tool more suitable for application  
352 to completely novel species.

353 PlASgraph not only provides a score for plasmidic and chromosomal contigs but also  
354 outputs a visualization of an assembly graph labeled according to the predictions. Figure 5  
355 shows parts of the assembly graph for *C. freundii* isolate SAMN15148288 with nodes colored  
356 according to the ground truth and both plASgraph and PlasForest predictions. The ground  
357 truth supports our initial reasoning to incorporate the information provided in the assembly  
358 graph, as linked contigs are more likely to belong to the same class. While both tools  
359 make some incorrect predictions, visualization clearly shows several isolated chromosome  
360 predictions among plasmid contigs and vice versa in the PlasForest prediction, whereas  
361 plASgraph has only one such isolated false positive. In general, plASgraph predictions are  
362 more consistent with the assembly graph topology.

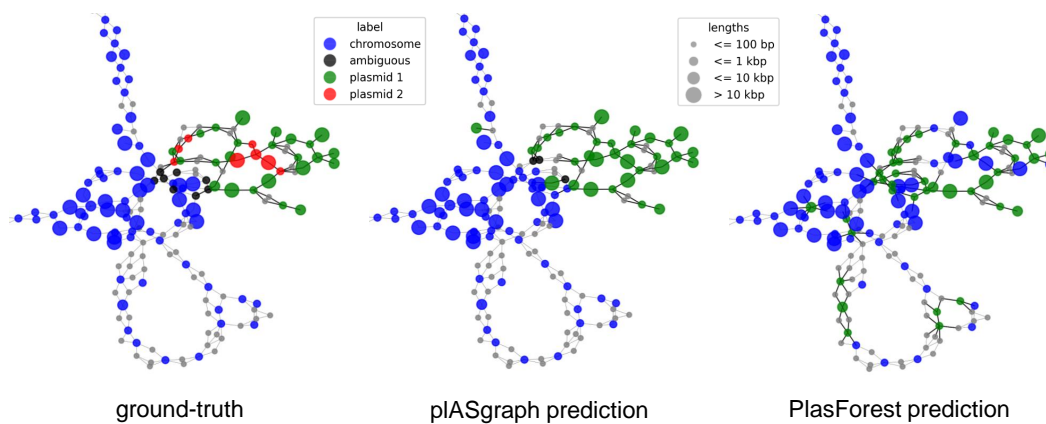
## 363 4 Discussion

364 PlASgraph is a GNN that can be used to identify plasmidic, chromosomal, and ambiguous  
365 contigs directly from a bacterial assembly graph. Our tool is easy to use, only requires a  
366 short-read assembly graph file as an input, and outperforms other state-of-the-art methods.

## 12 pIASgraph



■ **Figure 4 Evaluation of the species-agnostic pIASgraph model and comparison to PlasForest.** A: plasmid classification, B: chromosome classification. Horizontal lines represent the average F1-score of all isolates used for testing.



■ **Figure 5 Contig classification in the context of the assembly graph of *C. freundii* isolate SAMN15148288.** Chromosomal contigs are colored in blue and ambiguous contigs are colored in black. Left: The ground-truth, including two different plasmids (green and red). Middle: plASgraph predictions. Right: PlasForest predictions. Note that the classification tasks do not include binning of contig plasmids, thus all predicted plasmid contigs are color in green. The assembly graph extends to the upper left as a loop of chromosomal contigs alternating with unlabelled SNPs, which is not shown.

367        plASgraph is not dependent on specific species and can therefore be used also for newly  
368 sequenced bacteria for which no closed genome sequence is available yet. Our species-  
369 agnostic plASgraph model outperforms recently published, database-dependent PlasForest  
370 [23] approach when compared across different bacterial species. *De novo* classification  
371 (database independence) allows more accurate identification of previously unknown plasmids  
372 and chromosomes which can be critical for diverse One-Health epidemiologic surveillance.

373        However, when desired, plASgraph can also be trained for a particular species. Our species-  
374 specific models are more accurate when compared to mlplasmids [5], although mlplasmids  
375 performs better than plASgraph on longer contigs above 1 kbp. We hypothesize that  
376 mlplasmids is able to learn to recognize chromosome contigs of a particular species through  
377 their pentamer distributions at the cost of cross-species generalizability. Furthermore, this  
378 approach is unreliable for classifying contigs of lengths in the range of 100-1000 bp. Accurate  
379 classification of shorter contigs by plASgraph may enable identification of more complete  
380 plasmids from incomplete assemblies and has a potential to facilitate novel plasmid discovery.

381        Another novel feature of plASgraph is the separation of plasmid and chromosome classi-  
382 fication tasks, recognizing that some contigs are ambiguous, being parts of both types of  
383 molecules. These ambiguous contigs are an interesting subject for further study by themselves;  
384 our preliminary analysis of ambiguous contigs in our datasets suggests that the majority of  
385 them are related to transposons and phages. These mobile elements can integrate into both  
386 plasmids and chromosomes within the cell.

387        The simplicity of the architecture of the plASgraph model makes it amenable to extensions.  
388 For example, the use of additional information about plasmids, such as the presence of  
389 plasmid-specific genes in a contig, could allow further increase in classification accuracy as  
390 this additional information would propagate to nearby nodes thanks to the GNN architecture.  
391 In addition, it will be interesting to investigate how plASgraph could be adapted for accurate  
392 plasmid identification in metagenomic datasets, like wastewater samples, which play an  
393 increasingly crucial role in monitoring antibiotic resistance [13].

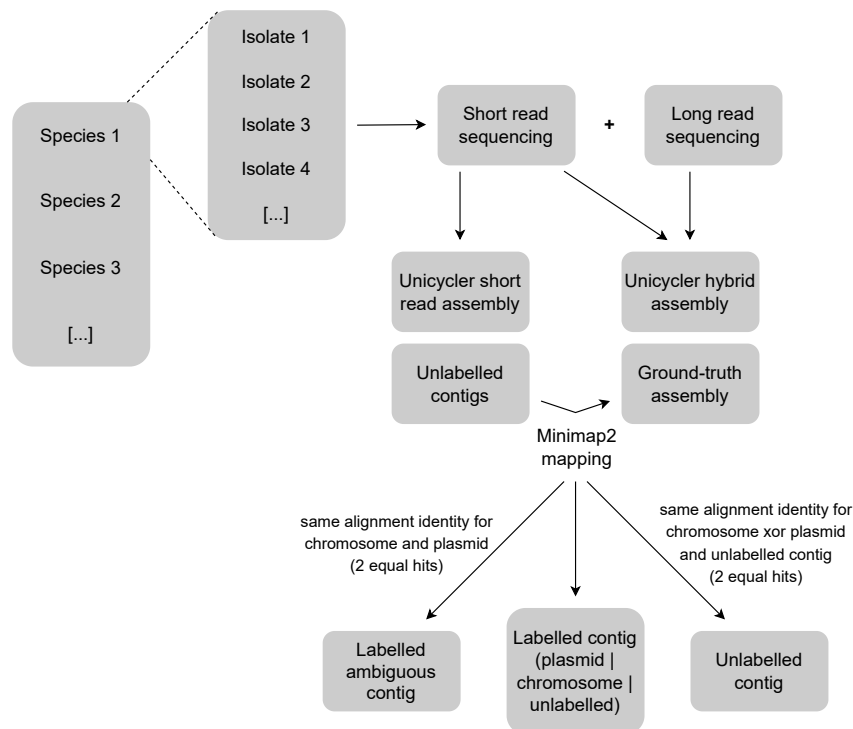
394 ——— **References** ———

- 395 **1** Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro,  
396 Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfel-  
397 low, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz  
398 Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore,  
399 Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever,  
400 Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol  
401 Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng.  
402 TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available  
403 from tensorflow.org. URL: <https://www.tensorflow.org/>, doi:10.5281/zenodo.4724125.
- 404 **2** Mislav Acman, Ruobing Wang, Lucy van Dorp, Liam P Shaw, Qi Wang, Nina Luhmann,  
405 Yuyao Yin, Shijun Sun, Hongbin Chen, Hui Wang, et al. Role of mobile genetic elements in  
406 the global dissemination of the carbapenem resistance gene blaNDM. *Nature Communications*,  
407 13(1):1–13, 2022. doi:10.1038/s41467-022-28819-2.
- 408 **3** William B Andreopoulos, Alexander M Geller, Miriam Lucke, Jan Balewski, Alicia Clum,  
409 Natalia N Ivanova, and Asaf Levy. Deepplasmid: deep learning accurately separates plasmids  
410 from bacterial chromosomes. *Nucleic Acids Research*, 50(3):e17–e17, 2021. doi:10.1093/nar/  
411 gkab1115.
- 412 **4** Sergio Arredondo-Alonso, Martin Bootsma, Yaír Hein, Malbert R C Rogers, Jukka Corander,  
413 Rob J L Willems, and Anita C Schürch. gplas: a comprehensive tool for plasmid analysis using  
414 short-read graphs. *Bioinformatics*, 36(12):3874–3876, 2020. doi:10.1093/bioinformatics/  
415 btaa233.
- 416 **5** Sergio Arredondo-Alonso, Malbert RC Rogers, Johanna C Braat, Tess D Verschuuren, Janetta  
417 Top, Jukka Corander, Rob JL Willems, and Anita C Schürch. mlplasmids: A user-friendly tool  
418 to predict plasmid-and chromosome-derived sequences for single species. *Microbial Genomics*,  
419 4(11), 2018. doi:10.1099/mgen.0.000224.
- 420 **6** Sergio Arredondo-Alonso, Rob J. Willems, Willem van Schaik, and Anita C. Schürch. On  
421 the (im)possibility of reconstructing plasmids from whole-genome short-read sequencing data.  
422 *Microbial Genomics*, 3(10):e000128, 2017. doi:10.1099/mgen.0.000128.
- 423 **7** Anton Bankevich, Sergey Nurk, Dmitry Antipov, Alexey A Gurevich, Mikhail Dvorkin,  
424 Alexander S Kulikov, Valery M Lesin, Sergey I Nikolenko, Son Pham, Andrey D Prjibelski,  
425 et al. SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing.  
426 *Journal of computational biology*, 19(5):455–477, 2012. doi:10.1089/cmb.2012.0021.
- 427 **8** Alessandra Carattoli. Plasmids and the spread of resistance. *International Journal of Medical*  
428 *Microbiology*, 303(6):298–304, 2013. doi:10.1016/j.ijmm.2013.02.001.
- 429 **9** François Chollet et al. Keras. <https://keras.io>, 2015.
- 430 **10** G. W. Cox, E. J. Parmley, B. P. Avery, R. J. Irwin, R. J. Reid-Smith, A. E. Deckert, R. L.  
431 Finley, D. Daignault, D. C. Alexander, V. Allen, S. El Bailey, S. Bekal, L. Chui, G. J. German,  
432 D. Haldane, L. Hoang, J. Minion, G. Zahariadis, M. R. Mulvey, and A. Bharat. A One-Health  
433 Genomic Investigation of Gentamicin Resistance in *Salmonella* from Human and Chicken  
434 Sources in Canada, 2014 to 2017. *Antimicrobial Agents and Chemotherapy*, 65(12):e0096621,  
435 2021. doi:10.1128/AAC.00966-21.
- 436 **11** Zhencheng Fang, Jie Tan, Shufang Wu, Mo Li, Congmin Xu, Zhongjie Xie, and Huaiqiu Zhu.  
437 PPR-Meta: a tool for identifying phages and plasmids from metagenomic fragments using  
438 deep learning. *GigaScience*, 8(6), 2019. doi:10.1093/gigascience/giz066.
- 439 **12** Daniele Grattarola and Cesare Alippi. Graph Neural Networks in TensorFlow and Keras with  
440 Spektral [application notes]. *IEEE Computational Intelligence Magazine*, 16(1):99–106, 2021.  
441 doi:10.1109/MCI.2020.3039072.
- 442 **13** Jianhua Guo, Jie Li, Hui Chen, Philip L Bond, and Zhiguo Yuan. Metagenomic analysis  
443 reveals wastewater treatment plants as hotspots of antibiotic resistance genes and mobile  
444 genetic elements. *Water research*, 123:468–478, 2017. doi:10.1016/j.watres.2017.07.002.

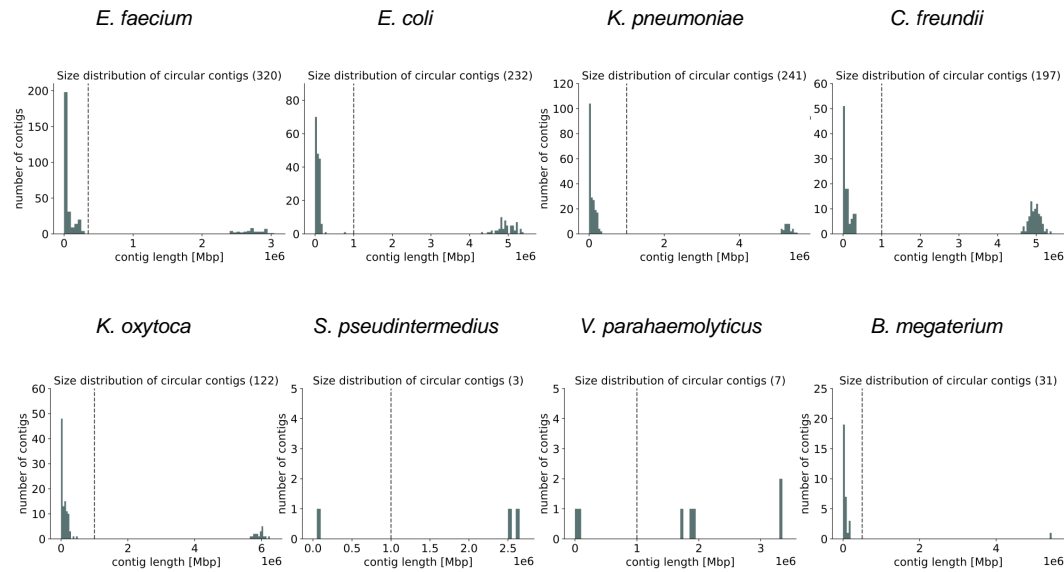


- 445 14 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv*  
446 *preprint arXiv:1412.6980*, 2014. doi:10.48550/arXiv.1412.6980.
- 447 15 Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional  
448 networks. *arXiv preprint arXiv:1609.02907*, 2016. doi:10.48550/arXiv.1609.02907.
- 449 16 Pawel S Krawczyk, Leszek Lipinski, and Andrzej Dziembowski. PlasFlow: predicting plasmid  
450 sequences in metagenomic data using genome signatures. *Nucleic Acids Research*, 46(6):e35–e35,  
451 2018. doi:10.1093/nar/gkx1321.
- 452 17 Andre Lamurias, Mantas Sereika, Mads Albertsen, Katja Hose, and Thomas Dyhre Nielsen.  
453 Metagenomic binning with assembly graph embeddings. *bioRxiv*, 2022. doi:10.1101/2022.  
454 02.25.481923.
- 455 18 Heng Li. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*, 34(18):3094–  
456 3100, 2018. doi:10.1093/bioinformatics/bty191.
- 457 19 Sara V Little, Laura K Bryan, Andrew E Hillhouse, Kranti Konganti, and Sara D Lawhon.  
458 Whole-genome sequences of *Staphylococcus pseudintermedius* isolates from canine and human  
459 bacteremia infections. *Microbiology Resource Announcements*, 8(28):e00735–19, 2019. doi:  
460 10.1128/MRA.00735-19.
- 461 20 Alex Orlek, Nicole Stoesser, Muna F. Anjum, Michel Doumith, Matthew J. Ellington, Tim  
462 Peto, Derrick Crook, Neil Woodford, A. Sarah Walker, Hang Phan, and Anna E. Sheppard.  
463 Plasmid classification in an era of whole-genome sequencing: Application in studies of antibiotic  
464 resistance epidemiology. *Frontiers in Microbiology*, 8:182, 2017. doi:10.3389/fmicb.2017.  
465 00182.
- 466 21 David Pellow, Itzik Mizrahi, and Ron Shamir. PlasClass improves plasmid sequence classifica-  
467 tion. *PLoS Computational Biology*, 16(4), 2020. doi:10.1371/journal.pcbi.1007781.
- 468 22 David Pellow, Alvah Zorea, Maraike Probst, Ori Furman, Arik Segal, Itzhak Mizrahi, and Ron  
469 Shamir. SCAPP: an algorithm for improved plasmid assembly in metagenomes. *Microbiome*,  
470 9:144, 2021. doi:10.1186/s40168-021-01068-z.
- 471 23 Léa Pradier, Taztio Tissot, Anna-Sophie Fiston-Lavier, and Stéphanie Bedhomme. PlasForest:  
472 a homology-based random forest classifier for plasmid detection in genomic datasets. *BMC*  
473 *Bioinformatics*, 22(1):1–17, 2021. doi:10.1186/s12859-021-04270-w.
- 474 24 Lianrong Pu and Ron Shamir. 3CAC: improving the classification of phages and plasmids  
475 in metagenomic assemblies using assembly graphs. *bioRxiv*, 2022. doi:10.1101/2021.11.05.  
476 467408.
- 477 25 Liam P Shaw, Kevin K Chau, James Kavanagh, Manal AbuOun, Emma Stubberfield, H Soon  
478 Gweon, Leanne Barker, Gillian Rodger, Mike J Bowes, Alasdair TM Hubbard, et al. Niche  
479 and local geography shape the pangenome of wastewater-and livestock-associated *Enterobac-*  
480 *teriaceae*. *Science Advances*, 7(15):eabe3868, 2021. doi:10.1126/sciadv.abe3868.
- 481 26 Philip S Shwed, J Crosthwait, K Weedmark, E Hoover, and F Dussault. Complete genome  
482 sequences of *Priestia megaterium* type and clinical strains feature complex plasmid arrays.  
483 *Microbiology Resource Announcements*, 10(27):e00403–21, 2021. doi:10.1128/MRA.00403-21.
- 484 27 Linda van der Graaf-van Bloois, Jaap A. Wagenaar, and Aldert L. Zomer. RFPlasmid:  
485 predicting plasmid sequences from short-read assembly data using machine learning. *Microbial*  
486 *Genomics*, 7(11), 2021. doi:10.1099/mgen.0.000683.
- 487 28 Ryan R Wick, Louise M Judd, Claire L Gorrie, and Kathryn E Holt. Unicycler: resolving  
488 bacterial genome assemblies from short and long sequencing reads. *PLoS Computational*  
489 *Biology*, 13(6):e1005595, 2017. doi:10.1371/journal.pcbi.1005595.
- 490 29 Lucy Witherall, Sariqa Wagley, Clive Butler, Charles R Tyler, and Ben Temperton. Genome  
491 sequences of four *Vibrio parahaemolyticus* strains isolated from the English Channel and  
492 the River Thames. *Microbiology Resource Announcements*, 8(24):e00392–19, 2019. doi:  
493 10.1128/MRA.00392-19.
- 494 30 Fengfeng Zhou and Ying Xu. cBar: a computer program to distinguish plasmid-derived from  
495 chromosome-derived sequence fragments in metagenomics data. *Bioinformatics*, 26(16):2051–  
496 2052, 2010. doi:10.1093/bioinformatics/btq299.

497 **5** Supplementary figures and tables



**Figure S1 Labeling workflow.** A short-read and hybrid assemblies are generated for each isolate. Unlabeled contigs of the short-read assembly are then mapped against the ground-truth hybrid assembly using minimap2 [18]. In case of a unique best alignment, the short-read contig is labeled according to the matching hybrid contig. If two equally good hits are identified to a chromosome and a plasmid hybrid assembly contig, the short-reads contig is labeled as 'ambiguous'. If two equally good hits are identified to a chromosome or a plasmid contig and to an unlabeled contig, the contig is 'unlabeled'.



**Figure S2 Circular contig size distribution in hybrid assemblies.** Above each plot, the total number of circular contigs is shown in parentheses. A vertical line marks the species-specific threshold in each plot; circular contigs shorter than the threshold are considered plasmidic, whereas longer contigs are considered chromosomal.

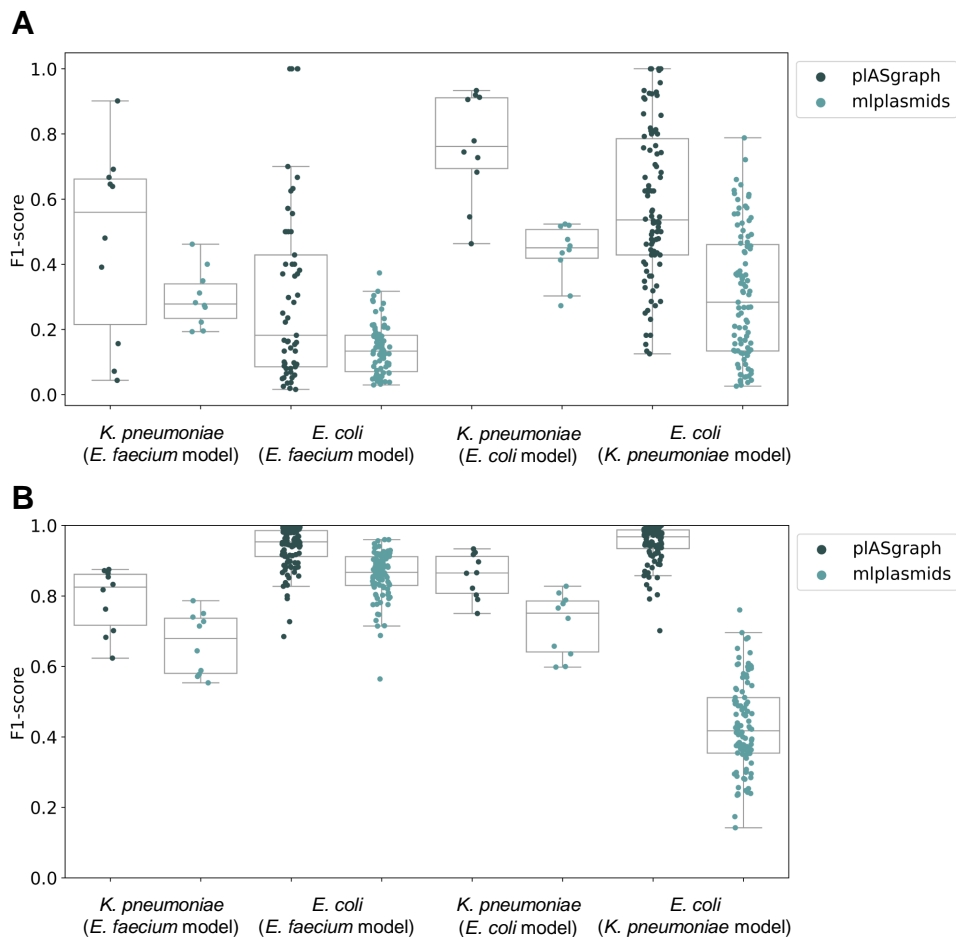
<i>E. faecium</i>		<i>E. coli</i>		<i>K. pneumoniae</i>		<i>C. freundii</i>	
no label	chromosome	no label	chromosome	no label	chromosome	no label	chromosome
0.59% (65)	62.44% (6885)	0.4% (167)	80.36% (33461)	0.26% (19)	53.08% (3838)	0.16% (13)	78.92% (6341)
414,386 bp	138,199,710 bp	1,595,167 bp	950,916,864 bp	282,939 bp	234,420,434 bp	180,679 bp	432,113,168 bp
plasmid	ambiguous	plasmid	ambiguous	plasmid	ambiguous	plasmid	ambiguous
26.23% (2892)	10.74% (1184)	14.63% (6094)	4.61% (1918)	37.57% (2717)	9.09% (657)	14.74% (1184)	6.19% (497)
11,280,374 bp	695,401 bp	37,172,761 bp	1,106,557 bp	13,824,751 bp	378,997 bp	9,836,037 bp	301,554 bp
<i>K. oxytoca</i>		<i>S. pseudintermedius</i>		<i>V. parahaemolyticus</i>		<i>B. megaterium</i>	
no label	chromosome	no label	chromosome	no label	chromosome	no label	chromosome
0.05% (2)	57.92% (2487)	1.2% (2)	94.58% (157)	1.01% (3)	93.6% (278)	1.29% (4)	33.23% (103)
80,418 bp	128,494,420 bp	435 bp	7,523,737 bp	1,771 bp	15,246,620 bp	126,852 bp	10,059,564 bp
plasmid	ambiguous	plasmid	ambiguous	plasmid	ambiguous	plasmid	ambiguous
35.44% (1522)	6.59% (283)	1.81% (3)	2.41% (4)	1.35% (4)	4.04% (12)	51.29% (159)	14.19% (44)
7,960,839 bp	226,571 bp	40,596 bp	5,005 bp	109,047 bp	5,702 bp	1,120,186 bp	29,909 bp

**Figure S3 Distribution of short-read contig labels.** Darker colors represent a higher percentage of the respective label in each dataset. Each small square shows the percentage, the absolute number, as well as the cumulative length of the short-read contigs with the respective label.

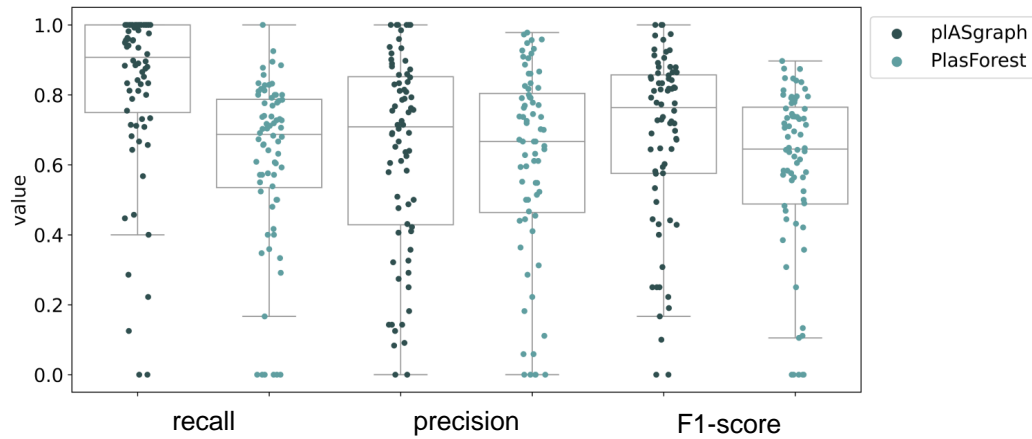
## 18 pIASgraph

■ **Table S1 Overview of the performed comparative analyses.**

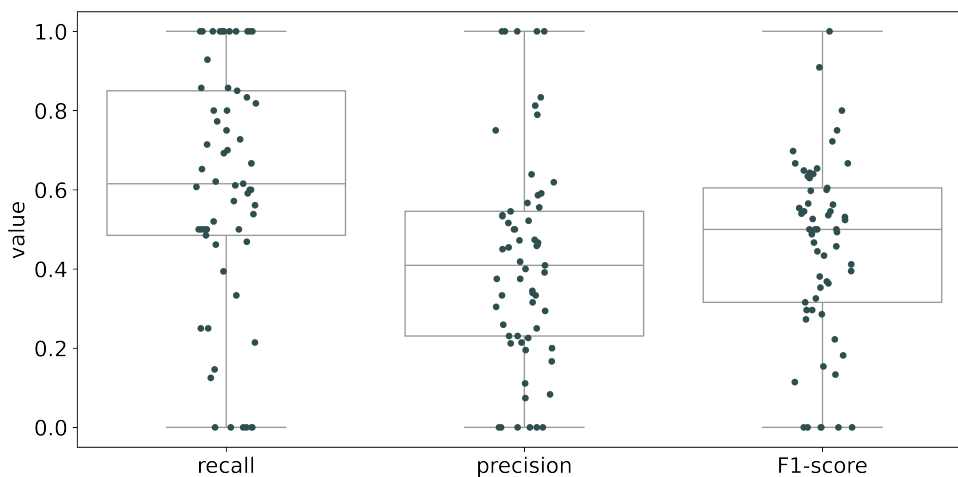
Evaluated on	Trained on	Benchmarked against	Experiment type
<i>E. coli</i>	<i>E. coli</i>		species-specific
<i>E. faecium</i>	<i>E. faecium</i>	mplasmids	
<i>K. pneumoniae</i>	<i>K. pneumoniae</i>		
<i>E. coli</i>	<i>E. faecium</i>	mplasmids	cross-species
<i>E. coli</i>	<i>K. pneumoniae</i>		
<i>K. pneumoniae</i>	<i>E. faecium</i>		
<i>K. pneumoniae</i>	<i>E. coli</i>		
<i>B. megaterium</i>	<i>E. faecium</i>	PlasForest	generalized
<i>C. freundii</i>	+		
<i>K. oxytoca</i>	<i>E. coli</i>		
<i>S. pseudintermedius</i>	+		
<i>V. parahaemolyticus</i>	<i>K. pneumoniae</i>		



■ **Figure S4 Cross-species evaluation of the species-specific pIASgraph models in comparison to mplasmids. A: Plasmid classification. B: Chromosome classification.**

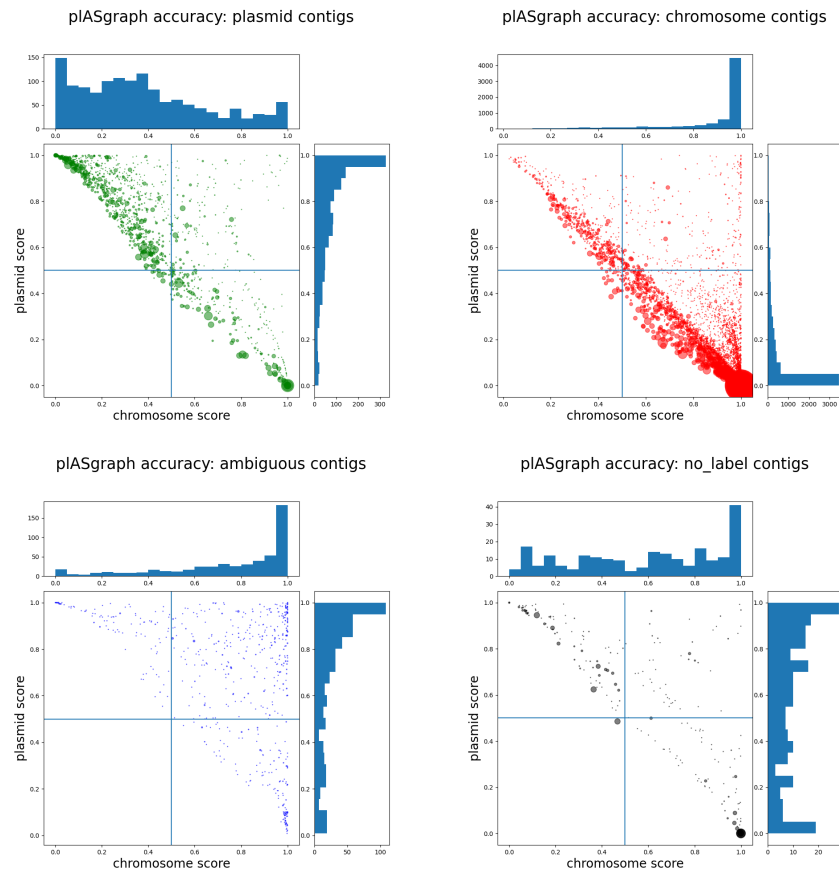


■ **Figure S5 Accuracy of plasmid classification for species-agnostic pIASgraph model in comparison to PlasForest.** Each data point represents one isolate; the plot combines isolates from all species used in Fig. 4.



■ **Figure S6 Accuracy of ambiguous contig classification for species-agnostic pIASgraph model.** Each data point represents one isolate; the plot combines isolates from all species used in Fig. 4. The isolates with no ambiguous contigs are not shown.

20 pIASgraph



■ **Figure S7** Results on the 96 *C. freundii* samples (9118 contigs). Each dot represents a contig, its radius being proportional to the contig length. Contigs are split in four panels according to their ground truth label. Each contig is shown with coordinates being its chromosome score (x-axis) and its plasmid score (y-axis).