1 **Accurate plasmid reconstruction from metagenomics data using assembly-alignment**
2 **graphs and contrastive learning**

3

4 **AUTHORS**

5 Pau Piera Líndez[1], Lasse Schnell Danielsen[1], Iva Kovačić [2], Marc Pielies Avellí[1], Joseph Nesme [2],

6 Lars Juhl Jensen[3], Jakob Nybo Nissen[1*], Søren Johannes Sørensen[2,*], Simon Rasmussen[1,*]

7

8 [1] Novo Nordisk Foundation Center for Basic Metabolic Research, Faculty of Health and Medical

9 Sciences, University of Copenhagen, Denmark

10 [2] Section of Microbiology, Department of Biology, University of Copenhagen,

11 Universitetsparken 15, DK-2100, Copenhagen, Denmark

12 [3] Novo Nordisk Foundation Center for Protein Research, Faculty of Health and Medical

13 Sciences, University of Copenhagen, Copenhagen, Denmark

14

15 * To whom correspondence should be addressed. Email: jakob.nissen@sund.ku.dk,

16 sjs@bio.ku.dk, srasmuss@sund.ku.dk

17

## ABSTRACT

Plasmids are extrachromosomal DNA molecules that enable horizontal gene transfer in bacteria, often conferring advantages such as antibiotic resistance. Despite their significance, plasmids are underrepresented in genomic databases due to challenges in assembling them, caused by mosaicism and micro-diversity. Current plasmid assemblers rely on detecting circular paths in single-sample assembly graphs, but face limitations due to graph fragmentation and entanglement, and low coverage. We introduce PlasMAAG (Plasmid and organism Metagenomic binning using Assembly Alignment Graphs), a framework to recover plasmids and organisms from metagenomic samples that leverages an approach that we call "assembly-alignment graphs" alongside common binning features. On synthetic benchmark datasets, PlasMAAG reconstructed 50–121% more near-complete plasmids than competing methods and improved the Matthews Correlation Coefficient of geNomad contig classification by 28–106%. On hospital sewage samples, PlasMAAG outperformed all other methods, reconstructing 33% more plasmid sequences. PlasMAAG enables the study of organism-plasmid associations and intra-plasmid diversity across samples, offering state-of-the-art plasmid reconstruction with reduced computational costs.

## INTRODUCTION

Plasmids are extrachromosomal DNA molecules within a host cell that are physically separated from chromosomal DNA and can replicate independently (1–3). Plasmids differ in genome length, copy number, replication mechanism, and conjugation mode. The part of the plasmid that encodes the core replication machinery is typically contiguous and is called the 'backbone'. The replication and maintenance of plasmids incur a metabolic burden for the host, and to avoid purifying selection, the plasmid must carry additional 'cargo' genes that increase the fitness of either the host or the plasmid (4). This may be genes that confer antibiotic resistance (5–7).

Approximately 50% of bacteria carry one or more plasmids (8). Nonetheless, in databases, sequences from plasmids remain underrepresented compared to those from cellular genomes. For instance, RefSeq contains 82,471 bacterial genomes, but only 7,892 plasmids (9). Characterization of environmental plasmids in *in vitro* conditions is inherently limited by the so called "cultivation bottleneck" (10, 11), where laboratory conditions modify microbial

48   diversity, offering a poor representation of the original composition. Therefore, despite the

49   great number of plasmid-related studies, most studies have investigated plasmid virulence

50   and properties in isolated strains (8, 12–14). This mismatch between estimated plasmid

51   prevalence in bacteria and plasmid representation in the databases emphasizes our

52   incomplete understanding of the plasmid genetic structure, diversity, and function (2).

53   Metagenomic offers culture free alternative techniques; however, the genetic complexity of

54   environmental samples complicates the process (10, 11, 15). Besides the challenges of

55   assembling bacterial chromosomes, assembly of plasmids bring additional challenges: a)

56   plasmids undergo frequent recombination, creating groups of plasmids that share a 'backbone'

57   but diverge on their 'cargo' sequence (12); b) plasmids at high copy number have higher

58   mutation rates than chromosomes, which increases micro-diversity and makes them difficult

59   to assemble with de Bruijn-graph based assemblers (16) and c) plasmids are enriched for

60   repeated sequences associated with transposable elements (2). A consequence of this is that

61   plasmid sequences will be fragmented across assemblies and entangled by sharing the same

62   'backbone' and repeated genetic elements (12, 17).

63   To overcome these challenges, dedicated metagenomic plasmid assemblers such as Recycler,

64   metaplasmidSPAdes, and SCAPP have been developed (18–20). These methods rely on the

65   assembly graph, a data structure used by metagenome assemblers, that represents overlaps

66   between sequencing reads. Assembly graphs represent contiguous sequences (contigs) as

67   nodes and overlaps between these as edges (21). By leveraging assembly graphs, plasmid

68   assemblers can identify connected sequences and resolve complex genomic regions (22).

69   Recycler re-interprets the metagenomic assembly graph, leveraging paired-reads information,

70   and attempting to extract subgraph cycles with uniform coverage from the graph, in a process

71   named graph 'peeling' also used by SCAPP (18). MetaplasmidSPAdes iteratively extracts cyclic

72   subgraphs from the assembly graph with uniform coverage, filtering the subgraphs with

73   plasmidVerify (19), a tool that classifies sequences into plasmidic and chromosomal based on

74   gene content using a profile-HMM (19). Finally, SCAPP tries to find plasmid cyclic paths in

75   assembly graphs based on paired read mappings, presence of plasmid-specific genes,

76   sequence length, coverage, and plasmid sequence score annotation based on PlasClass (18,

77   23). Common to the methods is that they operate on single-sample assembly graphs, use the

78   circularity of plasmids, and contig coverage. However, the methods have fundamental

79    limitations. First, low coverage causes the "fragmentation problem", where some plasmids

80    appear disconnected in the graph, making them impossible to identify by graph peeling (14,

81    24). Second, the high recombination rate of plasmids causes entangled assembly graph

82    components where circularity is hard to identify (25). Finally, SCAPP and metaplasmidSPAdes

83    leverage plasmid gene signatures to guide the plasmid candidate's path extraction from the

84    assembly graph.

85    Binning is a computational strategy used to reconstruct genomes by grouping contigs based

86    on their genome of origin, providing an alternative to assembly graph-based methods.

87    Modern binners typically integrate several sequence features, including k-mer composition

88    (26–31), abundance patterns across samples (26–31), assembly graph connectivity (32), and

89    taxonomic markers or annotations (28, 30, 33). Most of these features can be computed on a

90    per-sequence basis and are therefore not vulnerable to the fragmentation problem suffered

91    by assembly graphs with low coverage. Furthermore, it has been shown how binning

92    information can be used to refine contig classification, using binning features to guide

93    classification rather than contig classification to reconstruct the original sequences (34, 35).

94    We have previously developed the binning tool VAMB, which combines several of these

95    features using a variational autoencoder into a latent space which is clustered to form bins (26,

96    36).

97    In this paper, we introduce assembly-alignment graphs (AAGs), which combine the intra-

98    sample sequence overlaps recorded by assembly graphs, with cross-sample overlaps detected

99    by ordinary sequence alignment. Using a contrastive learning approach, we were able to

100    effectively integrate the AAG with ordinary binning features in a new binning framework called

101    PlasMAAG that can reconstruct both plasmids and cellular genomes. We evaluated PlasMAAG

102    on simulated data, where it reconstructed 9-70% more near-complete (≥0.95 precision, ≥0.9

103    recall) (NC) plasmids and cellular genomes than SemiBin2, ComeBin, MetaBAT2, MetaDecoder,

104    and VAMB. Regarding only plasmids, PlasMAAG reconstructed at least 50-121% more NC

105    plasmids than any other binner, and 14-40% more NC plasmids than the unfiltered set of cycles

106    from SCAPP cycles in 4/5 benchmark datasets. When using a confident threshold PlasMAAG

107    reconstructed 21-212% more confident NC plasmids than SCAPP confident. PlasMAAG

108    achieved excellent performance on hospital sewage samples, reconstructing at least 33% more

109    plasmid sequences than any other tool, as evaluated with a robust paired long-read short-

110   read validation setup. Using PlasMAAG's ability to reconstruct plasmids and hosts, we studied

111   host-plasmid associations in hospital sewage samples and intra-plasmid diversity across

112   samples. To our knowledge, PlasMAAG is the only method that enables 'multi-sample'

113   characterization of both plasmids and cellular genomes, achieving state-of-the-art plasmid

114   reconstruction with reduced computational resource requirements than current plasmid

115   binners.

116   **RESULTS**

117   **PlasMAAG: Combining assembly graphs, alignment graphs, TNFs, and co-abundances**

118   **for binning**

119   PlasMAAG is a new deep learning binning algorithm designed to reconstruct cellular and

120   plasmid genomes (**Figure 1**). Compared to our previous developed binning algorithm VAMB,

121   PlasMAAG introduces three novelties. First, we combine multi-sample assembly graphs with

122   contig alignment graphs into a single graph called 'assembly-alignment graph'. The assembly-

123   alignment graph is then projected to an embedding space with fastnode2vec (37), from which

124   communities of contigs can be extracted. Second, we enhanced the training of the variational

125   autoencoder (VAE) by adding contrastive learning, based on information from the assembly-

126   alignment graph. Finally, we leverage binning to ensemble geNomad (38) contig annotation

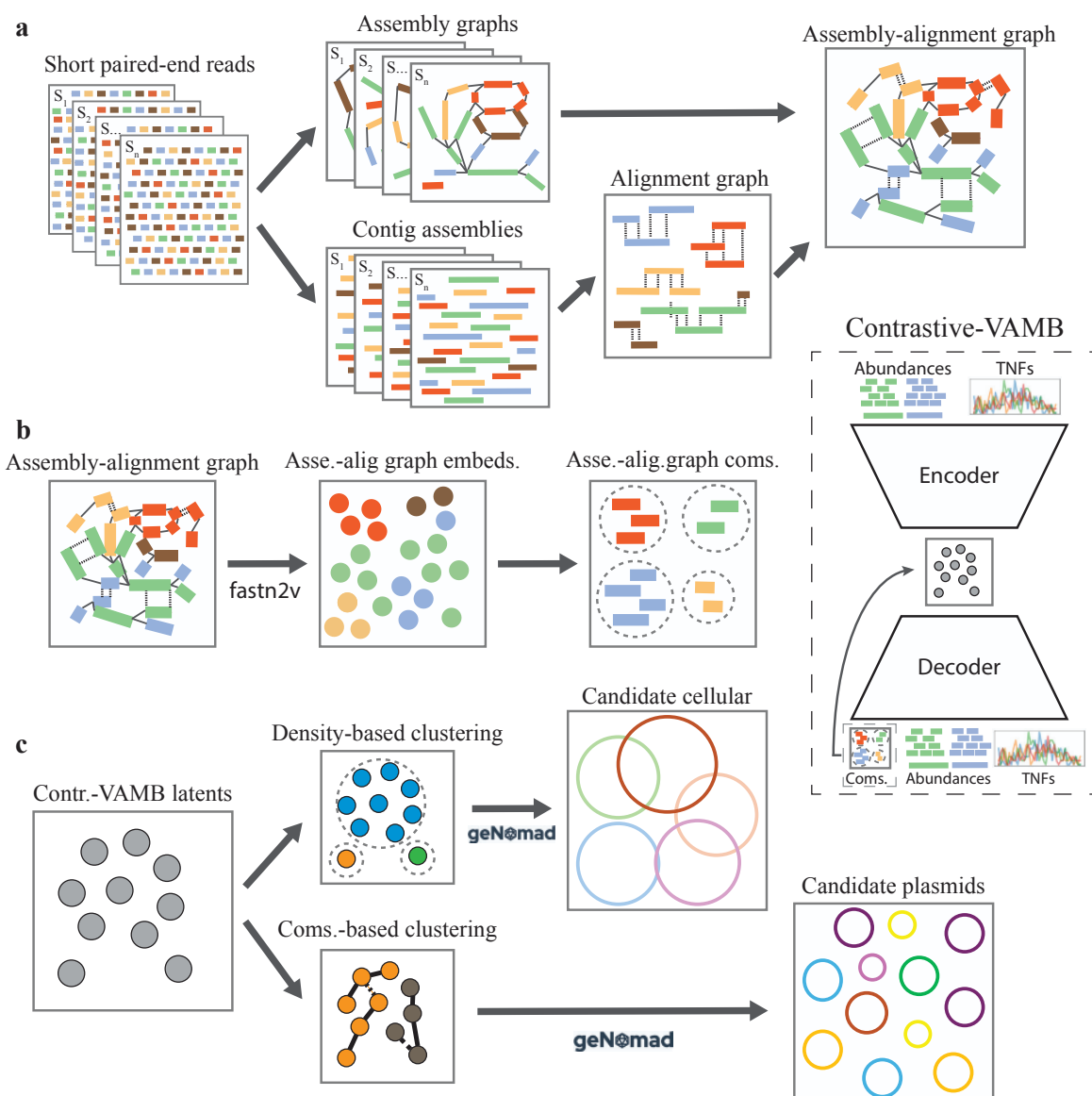127   scores across bins to classify the bins into plasmid or cellular genomes.

**Figure 1. PlasMAAG workflow overview.** PlasMAAG leverages assembly graphs, alignment graphs, k-mer signal, and contig co-abundances for binning, with a final step where bins are classified as cellular or plasmid bins based upon refined geNomad predictions. **a**. Per-sample assembly graphs are merged with the between-sample alignment graph, generating the assembly-alignment graph. **b**. Fastnode2vec is used to generate contig embeddings from the assembly-alignment graph, from where contig communities are extracted. Communities are expanded, merged, and purified using a variational autoencoder with contrastive loss that push communities towards be preserved in the embedding. **c**. Plasmid and cellular candidate bins are extracted from the VAE embedding based on their geNomad scores, using distinct plasmid and cellular clustering strategies.

## PlasMAAG reconstructed 21-212% more plasmid bins compared to SCAPP confident

To develop and test PlasMAAG we re-assembled the simulated CAMI2 short-read human microbiome toy datasets. Re-assembly of CAMI2 was required because the original CAMI2 plasmids were not simulated as independent entities from their hosts cellular genomes, and because assembly graphs were not available (see Methods). We found that PlasMAAG

141    reconstructed 5-64% more NC bins over all 5 benchmark datasets compared to VAMB, the

142    second best performing binner on the benchmark data **(Figure 2.A).** The improvement in

143    binning performance was driven by increased reconstruction of plasmids, where PlasMAAG

144    reconstructed at least 50-121% more NC plasmids candidate bins compared to SemiBin2,

145    ComeBin, MetaBAT2, MetaDecoder, and VAMB across all benchmark datasets. When

146    comparing to SCAPP, PlasMAAG reconstructed 14-40% more NC plasmids than SCAPP cycles

147    over 4/5 benchmark datasets **(Figure 2.B)**. When evaluating confident plasmids bins

148    generated by PlasMAAG (above 0.95 geNomad plasmid threshold), PlasMAAG reconstructed

149    21-212% more NC plasmid bins than SCAPP confident (**Figure 2.C**). Furthermore, PlasMAAG

150    spanned a larger variation of plasmids, since the unique set of confident PlasMAAG plasmids

151    across the benchmark datasets included 172 NC plasmid bins and 223 medium-quality (≥0.9

152    precision, ≥0.5 recall) (MQ) plasmid bins not reconstructed by SCAPP confident. In contrast,

153    SCAPP confident reconstructed 64 NC plasmid bins and 68 MQ plasmid bins not reconstructed

154    by PlasMAAG. The intersecting set of plasmids reconstructed by both methods was 164 NC,

155    and 185 MQ plasmid bins, respectively (**Figure 2.D-E**)**.** Considering cellular binning, PlasMAAG

156    was also competitive, reconstructing 0.7-9% less NC cellular bins than VAMB, the best cellular

157    binner on benchmark datasets **(Figure 2.F)**. The set of PlasMAAG confident plasmids offered

158    a better balance than SCAPP confident between the true positive and true negative plasmids

159    present in the benchmark datasets, with a 14-43% improvement in F1 (**Figure 2.G-H**,

160    **Supplementary Figure 1**, **Supplementary Note 1**). By averaging geNomad scores across

161    PlasMAAG's clusters, we can detect plasmids more accurately than applying geNomad on

162    individual contigs, yielding an improvement over the plasmid/non-plasmid contig

163    classification Area Under Precision-Recall Curve (AUPRC) and Matthews correlation coefficient

164    (MCC) of between 28-69% and 42-131%, respectively (**Figure 2.I-J**, **Supplementary Note 2**,

165    **Supplementary Figure 2**, **Supplementary Table 1**).

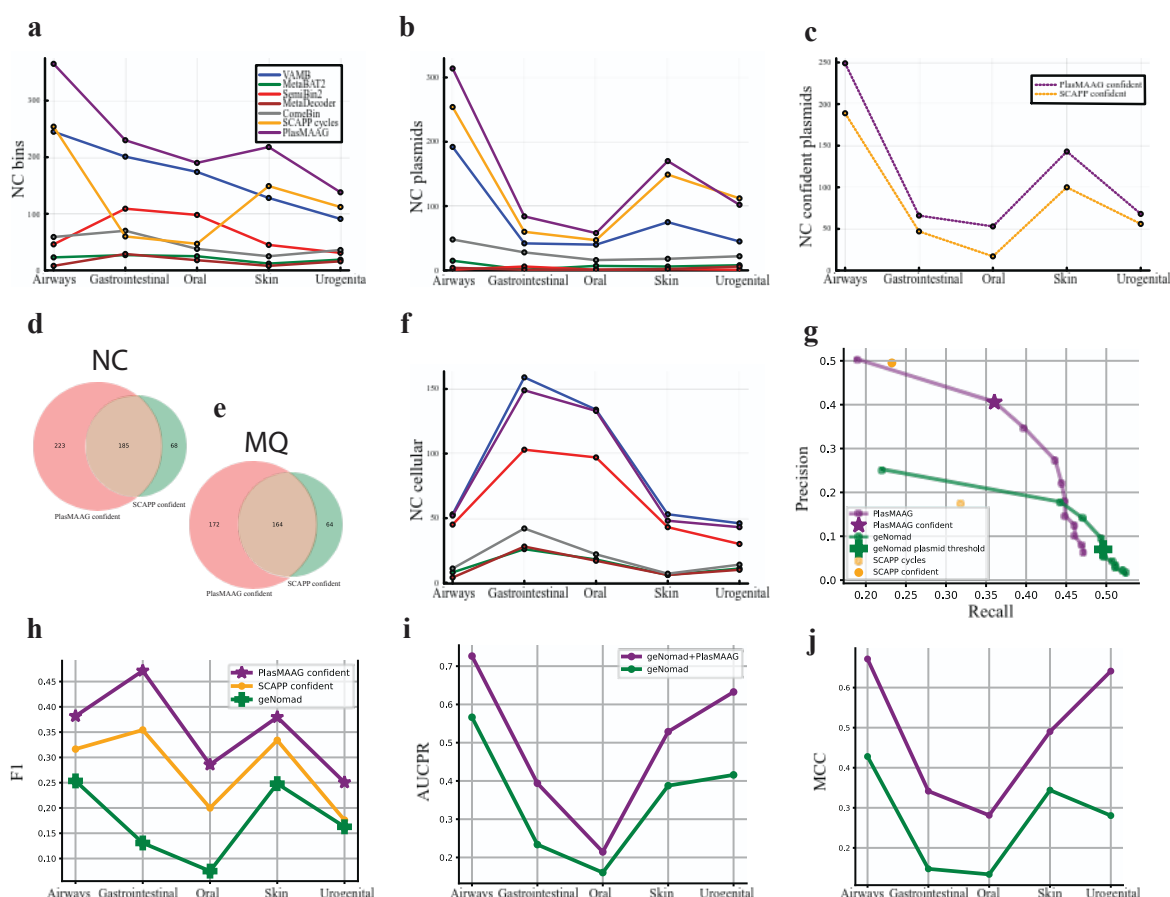**Figure 2. PlasMAAG binning and classification performance across the benchmark datasets. a.** NC bins (cellular + plasmids) reconstructed from the five benchmark datasets for VAMB (blue), MetaBAT2 (green), SemiBin2 (red), MetaDecoder (brown), ComeBin (grey), SCAPP cycles (yellow), and PlasMAAG (purple). **b.** NC plasmid bins reconstructed by all methods. **c.** NC plasmids reconstructed by SCAPP confident (yellow dotted) and PlasMAAG confident (purple dotted). **d.** Set of NC complete unique plasmids reconstructed only by PlasMAAG confident (red), only by SCAPP confident (green), and by both methods (light brown) across all datasets. **e.** Same than **d** but for MQ plasmid bins. **f.** NC cellular bins reconstructed by all methods except SCAPP confident. **g.** Plasmid sample precision-recall (see Methods) from the Airways dataset for PlasMAAG across geNomad thresholds (purple), PlasMAAG confident (purple star), geNomad across thresholds (green), geNomad at the default plasmid threshold (green cross), SCAPP cycles (light yellow), and SCAPP confident (dark yellow). **h.** Sample F1 across the five benchmark datasets for geNomad at the default plasmid threshold (green), SCAPP confident (yellow), and PlasMAAG confident (purple). **i.** Area Under Precision-Recall Curve (AUPRC) for the classification of plasmids by geNomad (green) and when aggregating the geNomad scores per PlasMAAG community-based clusters (purple). **j.** Matthew correlation coefficient (MCC) for the classification of plasmids by geNomad (green) and when aggregating the geNomad scores per PlasMAAG community-based clusters (purple).

## Assembly graphs have a strong signal for binning

In assembly graphs, edges represent sequence overlaps between contigs. Therefore, it has long been known that they are informative for binning (32, 39). To quantify how informative edges were, we weighted them by *normalized linkage* (see Methods), based on the number of overlapping k-mers, and the length of the contigs. Normalized linkage showed a positive correlation with edge accuracy at genome (species) level, with Spearman correlation

186  coefficients 0.49-0.93 (0.86-0.98) across all benchmark datasets (**Figure 3.A**, **Supplementary**

187  **Figure 3**). Additionally, normalized linkage was evaluated for correlation with edge accuracy

188  (i.e. how often two contigs linked by an edge belong to the same genome) by calculating the

189  Area Under Precision-Recall Curve (AUPRC). The resulting AUPRC ranged from 0.66 to 0.74 at

190  the genome level and 0.81 to 0.90 at the species level across the benchmark datasets (**Figure**

191  **3.B**, **Supplementary Figure 3**). We concluded that the assembly graph contains useful signals

192  for binning.

193  **Alignment graphs contain taxonomic information across samples**

194  PlasMAAG uses the multi-split binning workflow due to its superior accuracy (26, 36), where

195  samples are assembled individually. Therefore, assembly graphs only inform about overlaps

196  between intra-sample contigs. To also include between-sample contig overlap information, we

197  aligned contigs across samples with strict criteria to accept a hit (see Methods). The alignments

198  were highly precise with an accuracy at genome (species) level of 57-95% (95-99%) (**Figure**

199  **3.C**, **Supplementary Figure 4**). By adding alignments between pairs of contigs as edges to

200  the assembly graph, we created an alignment-assembly graph (AAG), where we weighed each

201  edge by either alignment metrics (for alignment edges) and normalized linkage (for assembly

202  graph edges, see Methods). Alignment edge weight between two contigs correlated with

203  taxonomic relatedness of the contig's genomes, showing an 82-98 (98-100) Area Under

204  Precision-Recall Curve (AUPRC) across the benchmark datasets at genome (species) taxonomic

205  level (**Figure 3.D**, **Supplementary Figure 5**). Furthermore, there was a positive correlation

206  between the averaged alignment graph edges and the average accuracy, with a Spearman

207  correlation coefficient of 0.71-0.95 across all benchmark datasets (**Figure 3.E**, **Supplementary**

208  **Figure 5**).

209  **Assembly-alignment graphs integrate alignments and assembly graphs**

210  The complementarity between cross-sample alignments and the intra-sample assembly graph

211  connections in the AAG enabled us to integrate these in a unified graph, resulting in a

212  combined graph that we named 'assembly-alignment graph'. We evaluated the edges in the

213  assembly-alignment graphs across the benchmark datasets to assess whether higher edge

214  weights correspond to contigs that are taxonomically close, such as those from the same

215  genome. The edge weights in the assembly-alignment graph reflect taxonomic relationship

216   between sequences, consistent with the original assembly and alignment graphs, achieving a

217   AUPRC of 0.69-0.90 (0.93-0.97) across the benchmark datasets at genome (species) taxonomic

218   level (**Figure 3.F**, **Supplementary Figure 6**). Consistently with the AUPRC findings, we found

219   a positive correlation between the averaged edge weights and the average edge accuracy at

220   genome taxonomic level, with 0.20-0.97 Spearman correlation coefficients across benchmark

221   datasets (**Figure 3.G**, **Supplementary Figure 6**). The assembly-alignment graph integrates

222   assembly graphs and alignment information across samples into a unified object, where edge

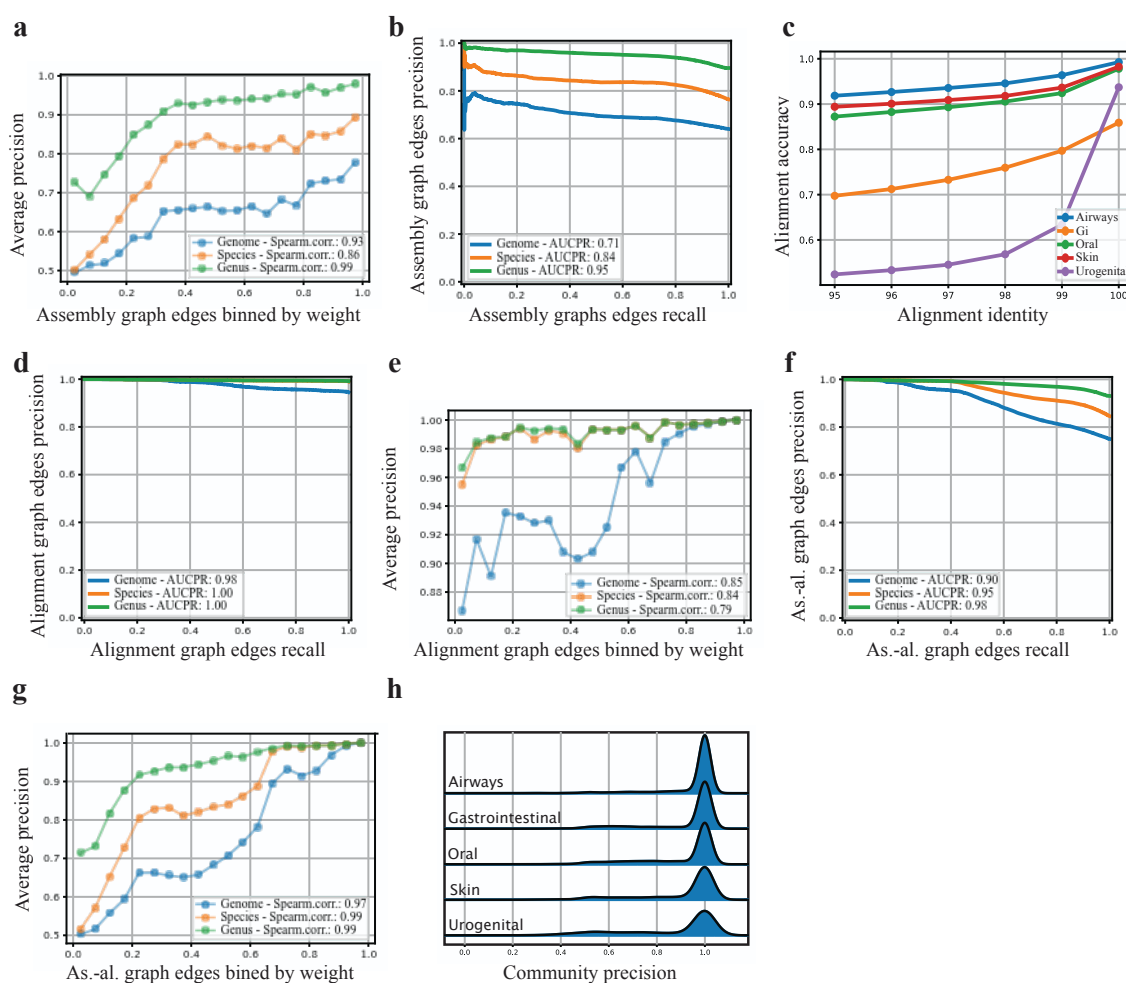223   weights reflect taxonomic relationships.



224
225   **Figure 3. Assembly graph, alignment graph, and assembly-alignment graph-based features for binning. a**. Average precision of the
226   assembly graph edges from the Airways benchmark dataset, sorted by edge weight and grouped into 5% bins, is shown for genome (blue),
227   species (orange), and genus (green) taxonomic levels. **b**. Precision-recall curve of the assembly graph edge weights from the Airways
228   benchmark dataset at genome (blue), species (orange), and genus (green) taxonomic levels. **c**. Alignment accuracy when increasing
229   minimum identity thresholds across benchmark datasets. Results are shown only for restrictive alignments (see Methods) between contigs
230   from different samples. **d**. Precision-recall curve of the alignment graph edge weights from the Airways benchmark dataset at genome
231   (blue), species (orange), and genus (green) taxonomic levels. **e**. Average precision of the alignment graph edges from the Airways
232   benchmark dataset, sorted by weight and grouped into 5% bins is shown for genome (blue), species (orange), and genus (green) taxonomic

233   levels. **f**. Precision-recall curve of the assembly-alignment graph edge weights from the Airways benchmark dataset at genome (blue),

234   species (orange), and genus (green) taxonomic levels. **g**. Average precision of assembly-alignment graph edges from the Airways

235   benchmark dataset, sorted by weight and grouped into 5% bins is shown for genome (blue), species (orange), and genus (green) taxonomic

236   levels. **h.** Precision distribution of communities extracted using FastNode2Vec from assembly-alignment graphs across the five benchmark

237   datasets at genome taxonomic level.

## Extracting high precision, low completeness communities from the assembly-alignment graph

240   The majority of contigs are too short to contain a stable signal for binning, but the AAG

241   cohesion depends on the nodes representing short contigs. Therefore, we condensed the AAG

242   into a set of node communities using fastnode2vec (see Methods). We found that the

243   extracted communities from this graph embedding had high purity, with an average precision

244   at genome (species) level of 86-95% (95-97%) across the benchmark datasets, and where 63-

245   84% (85-91%) of communities had a precision at genome (species) level (**Figure 3.H**,

246   **Supplementary Table 2**)**.** However, we observed that communities were composed of rather

247   few contigs, with 85-91% of the communities were composed of 10 or less contigs across the

248   benchmark datasets. Furthermore, microbial genomes were fragmented in, on average, 12.2-

249   32.8 communities, and plasmids somewhat were less fragmented, split between 1.6-2.5

250   communities on average (**Supplementary Figure 7**, **Supplementary Table 2**). We also

251   noticed that only 31-47% of contigs in the datasets belonged to any community

252   (**Supplementary Table 2**). In conclusion, the communities extracted from the AAG using

253   fastnode2vec were precise, but incomplete and fragmented.

## Contrastive variational autoencoders improve binning through aggregating, merging and splitting communities

256   To address the fragmentation of AAG communities, we leveraged traditional binning features

257   such as contig k-mer composition and abundances (40). In the VAMB framework, these contig

258   features are embedded using a variational autoencoder (VAE), and these embeddings are then

259   used to cluster contigs together. PlasMAAG follows the same approach but also considers

260   community structure during the embedding and clustering process. To encourage contigs of

261   the same community to be close in the embedding, we added an extra term to the loss

262   function of the VAE which penalized high embedding distance between contigs of the same

263   community. We call this term 'contrastive loss'. We then applied a clustering strategy on the

264   contrastive VAMB embeddings, consisting of three key steps: (1) Merging – Communities close

265    in the embedding were merged to reduce genome fragmentation, and increase genome recall.

266    (2) Splitting – Communities with contigs placed far apart in the embedding were split up to

267    increase precision. (3) Expansion – Unassigned contigs located close to a community in latent

268    space were added to the community to improve recall. We refer to these three steps as

269    'community-based' clustering (see Methods, **Supplementary Figure 8**). This community-

270    based clustering resulted in a 46–102% increase in genome recall across benchmark datasets

271    compared to the raw communities, confirming the effectiveness of the community merging

272    step. The splitting step improved precision by 0.03–1% (**Supplementary Figure 9**, **Table 3**),

273    indicating minor but positive impact without compromising recall. On the other hand,

274    community expansion had limited effect, with only a 1–3% increase in community size

275    (**Supplementary Table 3**), suggesting that step 3 had a smaller impact. Since recall increased

276    and precision slightly improved, F1 scores also increased, along with the number of

277    reconstructed near-complete (NC) bins.

### Contrastive loss had a positive impact on binning

279    To better understand the importance of the contrastive loss on the latent representations, we

280    evaluated how it impacted community-based clustering and clustering from the original VAMB,

281    which we call 'density-based' clustering. Community-based clustering with contrastive loss

282    achieved 28–63% higher average F1 scores compared to clustering without the contrastive

283    loss, reconstructing 7–45% more NC bins across the benchmark datasets (**Supplementary**

284    **Figures 10–11**). Contrastive loss also improved density-based clustering, causing a 57–162%

285    increase in F1 scores across all benchmark datasets (**Supplementary Figure 10**), but did not

286    uniformly increase the number of NC bins. NC bin recovery was increased by 1–6% in 4 out of

287    5 datasets but led to 16% fewer NC bins in one dataset due to a small decrease in precision

288    (**Supplementary Figure 11**). Overall, the contrastive loss boosted recall and led to significantly

289    higher F1 scores in both clustering approaches, whereas its effect on precision and final NC

290    bin counts varied depending on the dataset and clustering strategy, highlighting the trade-

291    offs introduced by enforcing graph-based community structures in the latent space.

### Differential embeddings of plasmids and organisms requires tailored clustering leveraging geNomad

294 As previously mentioned, we found the cellular genomes to be fragmented across more
295 communities than plasmids, probably due to the larger size of cellular genomes. Furthermore,
296 we observed distinct patterns in k-mer composition, contig co-abundance, and PlasMAAG
297 latent representations between plasmids and cellular genomes (**Supplementary Note 3**,
298 **Supplementary Figures 12-14**, **Supplementary Tables 3-4**). This suggested that community-
299 based clustering might be more suitable for plasmids, and density-based clustering for cellular
300 genomes, which we indeed verified with our benchmark datasets (**Supplementary Note 3**,
301 **Supplementary Figures 12-14**, **Supplementary Tables 3-4**). Therefore, to identify potential
302 plasmid communities in the AAG, we used geNomad to assign plasmid scores to each
303 community (38). We found that averaging geNomad scores across communities led to more
304 accurate plasmid identification compared to scoring individual contigs (**Supplementary Note**
305 **2**, **Supplementary Figure 2**, **Supplementary Table 1**). This allowed us to extract communities
306 as putative plasmid bins for community-based clustering and clustered the remaining contigs
307 using density-based clustering. Additionally, we found that this was sensitive to the geNomad
308 threshold used for the classification, particularly in the case of organisms (**Supplementary**
309 **Note 4**, **Supplementary Figures 15-16**). For instance, when setting a geNomad plasmid
310 threshold of 0.7, we observed a decrease on the NC cellular genomes (plasmids) of 6-39% (3-
311 18%) (**Supplementary Figure 17**). This indicated that the selection and dereplication process,
312 based on geNomad-identified plasmid clusters, led to a trade-off in cellular genomes recovery.
313 We conclude that integrating geNomad sequence predictions with PlasMAAG's diverse
314 clustering strategies enhanced binning performance, enabling the robust reconstruction of
315 both cellular genomes and plasmids.

**Evaluating PlasMAAG plasmid binning using hospital sewage samples long-read data,**
317 **and short-read plasmidomics data**

318 Validating PlasMAAG binning performance on real data is not straightforward as current tools
319 do not provide quality estimates for plasmids and might show inherent biases when exploring
320 understudied environments such as wastewater. We instead applied a binning validation
321 strategy based on sequencing both short- and long-read metagenomics from the same set of
322 samples (**Fig. 4.A**). We considered a long-read contig to be composed of a set of short-read
323 contigs if they aligned with 97% identity and a long-read contig coverage of 90% (see
324 Methods). By tallying the number of such sets of short-read contigs binned together, we got

325    a measure of recall of short-read contig binning. We observed that PlasMAAG community-

326    based bins reconstructed 21% more long-read contigs than VAMB, the second-best

327    performing binner (**Fig. 4.B**). Superior PlasMAAG binning performance was consistent even

328    when accounting for incompleteness of the long-read assembled contigs (**Supplementary**

329    **Note 5**, **Supplementary Figure 18**). To identify the subset of long-read contigs that

330    originated from plasmids, we sequenced samples after a plasmid enrichment to obtain paired

331    metagenomics and 'plasmidomics' samples as done previously (41) (see Methods). Long-read

332    contigs were defined as plasmid contigs if they were either (1) at least 50% covered by

333    plasmidomics reads or (2) circular and below 500 kb. We identified short-read contigs as

334    originating from plasmid if they aligned well to any long-read contig identified as plasmid (see

335    Methods). Using this criteria, PlasMAAG community-based reconstructed 138 NC plasmids,

336    which was 33% more plasmid long-read fragments than the second best binner VAMB, and

337    431% more NC plasmids than SCAPP cycles (**Figure 4.B**). These results were consistent with

338    the performance validated using unfiltered long-read contigs, demonstrating PlasMAAG's

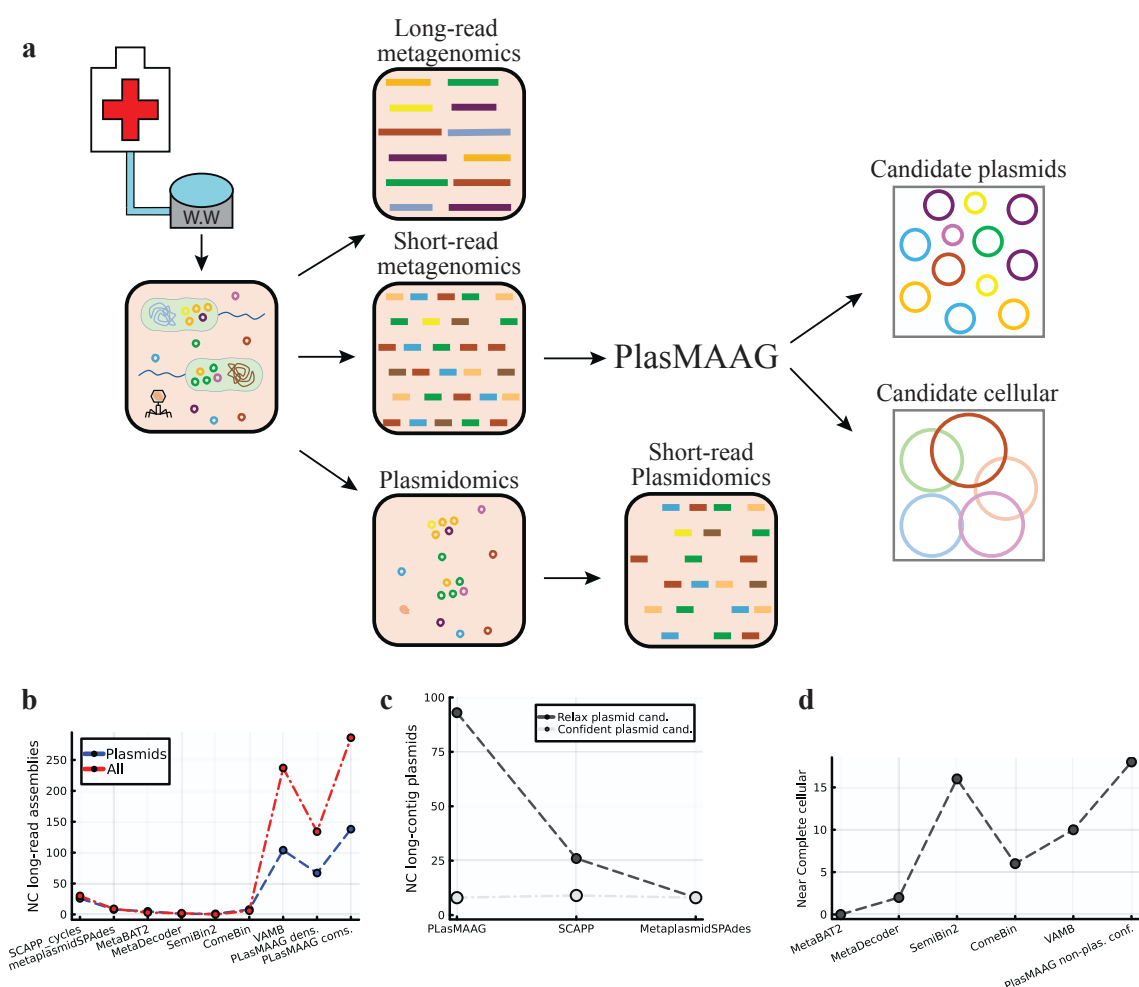339    robust binning capacity across diverse biological entities.

**Figure 4. PlasMAAG on real samples from hospital sewage. a**. Overview of the strategy used to validate PlasMAAG on the five hospital sewage samples. For each sample, long-read metagenomics, short-read metagenomics, and short-read plasmidomics datasets were generated (see Methods). PlasMAAG was applied to the short-read metagenomics data to produce candidate plasmid and cellular bins. These bins were validated against a reference assembly composed by long-read contigs to assess overall binning performance, and against a second reference assembly constructed from long-read contigs with plasmid evidence, identified either by circularity or plasmidomics read coverage. **b**. Binning performance of all methods across the five sewage samples, evaluated using all long-read contigs (red) and long-read contigs with plasmid evidence (blue). PlasMAAG dens.: bins produced using VAMB's density-based clustering algorithm on PlasMAAG's latents. PlasMAAG coms.: bins generated using the community-based clustering algorithm. **c**. Binning performance of PlasMAAG, SCAPP, and MetaPlasmidSPAdes under relaxed (light gray) and strict (dark gray) plasmid filtering criteria. **d**. NC cellular bins according to CheckM2 estimates, produced by all organism binners for the five hospital sewage samples. PlasMAAG non-plas. conf.: PlasMAAG density-based bins after extracting candidate plasmid contigs by aggregating geNomad plasmid contig scores per PlasMAAG community-based clusters (see Methods).

## Identifying plasmids in PlasMAAG bins using aggregated geNomad scores
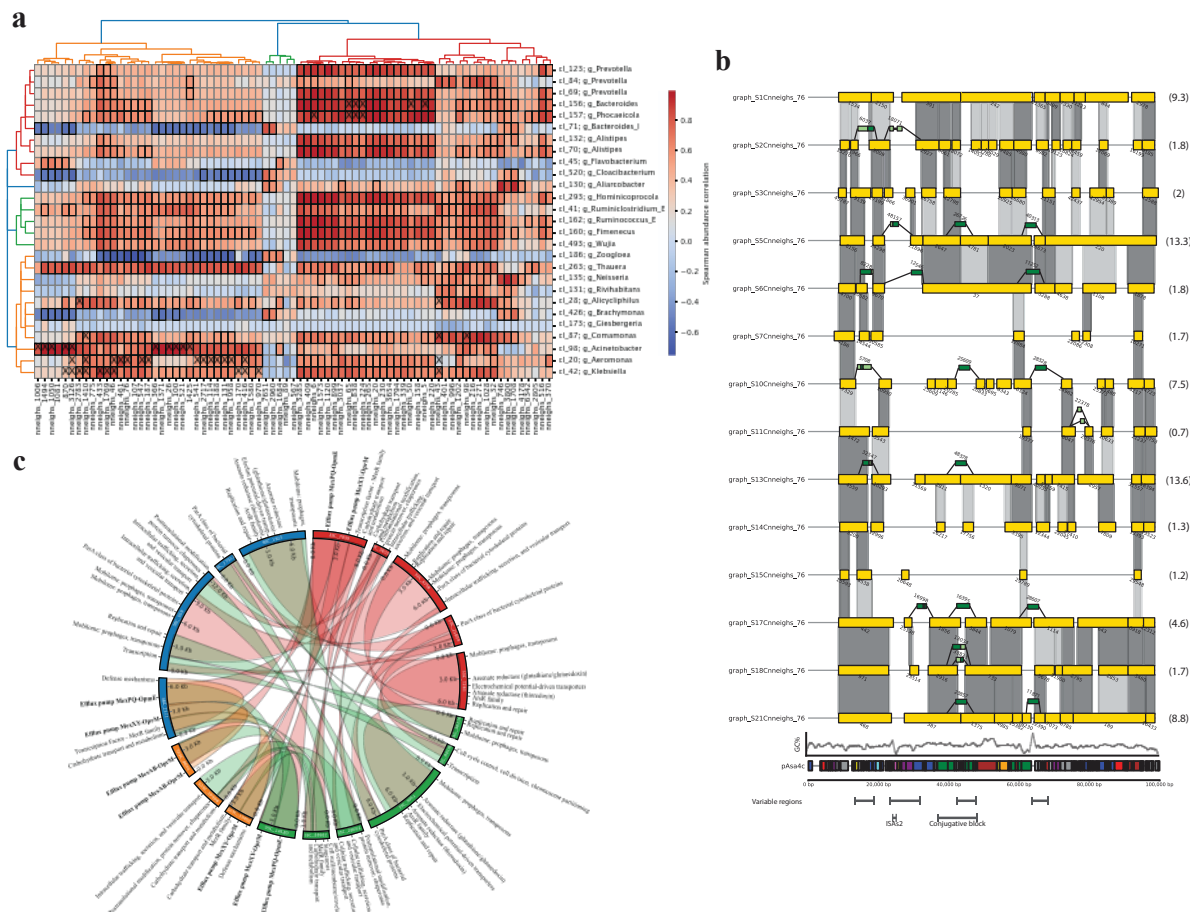
When applying PlasMAAG to a real dataset with thousands of bins and no ground truth, we need to define a threshold to determine whether a bin contains a plasmid. This threshold balances precision and recall. To aid in this decision, we aggregated geNomad's contig plasmid scores across all contigs within each bin. With a low threshold of 0.1, PlasMAAG reconstructed

358    93 NC long-read contigs highly confident plasmid based on the metaplasmidomics reads

359    (long-read plasmidomics, LR-P), which represented a loss of 33% compared to not filtering

360    with aggregated geNomad scores (**Figure 4.C**)**.** Using a stricter geNomad plasmid threshold

361    of 0.95 reduced the number of reconstructed LR-P to 8, a decrease of 94% (**Figure 4.D**). This

362    implied that most long-read contigs, where we had experimental plasmid evidence, were

363    predicted by geNomad to be of virus or chromosomal origin, as they had assigned a relatively

364    low plasmid score (**Figure 4.B, Figure 4.C, Supplementary Figure 19**)**.** By comparing

365    aggregated geNomad scores with experimental plasmid evidence, we found that this

366    mismatch mainly occurred where plasmid evidence was strong but not definitive

367    (**Supplementary Note 6**, **Supplementary Figure 19-21**). This contrasted with the consistency

368    observed in synthetic benchmarks, where geNomad generally demonstrated strong plasmid

369    predictive performance (**Supplementary Figure 16**). Finally, we investigated the effect of this

370    on cellular genomes and when applying a geNomad threshold of 0.95, the PlasMAAG density-

371    based bins, which are the ones not classified as plasmid, were evaluated with CheckM2. We

372    found 18 NC organisms, 3 more than SemiBin2, the second best binner on this dataset, and 8

373    more than VAMB. We found noticeable that PlasMAAG offered a better performance

374    compared to SemiBin2, even though SemiBin2 leveraged single-copy genes whereas

375    PlasMAAG did not. We conclude that PlasMAAG's has state of the art performance on real

376    datasets, both for reconstructing plasmids and cellular genomes.

377    **PlasMAAG enabled host-plasmids exploration from hospital sewage environments**

378    By reconstructing plasmid and cellular genomes from the same samples, PlasMAAG enables

379    an integrated analysis. We investigated host-plasmid abundance correlations of 24 hospital

380    sewage samples collected in Spain (see Methods). PlasMAAG produced 27,954 candidate

381    plasmid bins, and 213,431 non-plasmid candidate bins. PlasMAAG plasmid bins were

382    aggregated into 13,912 cross-sample clusters, and bacterial hosts per plasmid cluster were

383    inferred from PLSDB (see Methods). We identified 323 High quality cellular organism bins (HQ,

384    completeness ≥ 70%, contamination ≤ 10%) and aggregated these using PlasMAAG cross

385    sample cluster information. We found several significant positive correlations between

386    candidate plasmid and cellular organism bins, for example, cluster *cl_20*, annotated as

387    belonging to the *Aeromonas* genus, correlated with up to 41 plasmid clusters (adjusted p-

388    value < 0.05), 12 of which were previously reported as known host-plasmid associations in the

389    PLSDB database (**Figure 5A**). On the other hand, cluster *cl_293*, annotated as *Ruminococcus_E*

390    genus, correlated with 43 plasmid clusters, none of them previously reported in PLSDB (**Figure**

391    **5A**).



392

393

394    **Figure 5. PlasMAAG enables host-plasmid association studies and exploration of intra-plasmid variation across environments,**

395    **demonstrated using 24 hospital sewage samples. a**. Spearman correlation between PlasMAAG high-quality (HQ) cellular clusters and

396    PlasMAAG plasmid clusters with an aggregated geNomad plasmid score above 0.75. Highlighted cells with bold rectangles indicate

397    significant correlations after Benjamini-Hochberg FDR correction. Cells marked with "X" represent plasmid-organism associations

398    previously reported in PLSDB. The organism cluster dendrogram was generated using GTDB-tk taxonomic annotations, while the plasmid

399    cluster dendrogram was based on abundance correlations. **b**. PlasMAAG plasmid cluster *nneighs_416* bins. Each row represents a bin from

400    one sample, and numbers within parenthesis indicate median bin depth. Yellow blocks denote contigs aligned to *pAsa4c*, sorted by

401    alignment position. Dark green blocks represent contigs not mapping to *pAsa4c* (see Methods), with their positions inferred from matches

402    to other PLSDB plasmid accessions. Light green sections withing dark green blocks indicate alignment segments to *pAsa4c*. Dark grey

403    areas indicate alignment graph edges, and light grey areas represent non-restrictive alignment matches (see Methods). GC%: Average GC

404    content computed using a 1000 kb window. Colour code for pAsa4c regions: Blue (Replication and maintenance), Green (Conjugative

405    transfer), Purple (Recombination and DNA repair), Orange (Secretion and surface structures), Red (Metabolism), Yellow (Enzymes), Cyan

406    (Regulatory proteins and transcription factors), Brown (Transposases and mobile genetic elements), Gray (Hypothetical or unclassified). **c**.

407    PlasMAAG plasmid cluster *nneighs_76*, composed of contigs from sample 6 (blue), sample 3 (red), sample 5 (green), and sample 23 (orange).

408    Links represent alignment regions, coloured according to the sample of origin. Bold annotations indicate functions associated to
409    antimicrobial resistance.

**PlasMAAG revealed intra-plasmid variation across hospital sewage samples**

411    In PlasMAAG, contigs from different samples are projected into a shared latent space, enabling

412    them to be clustered together, and split into per-sample bins thereafter. Aggregation of bins

413    into PlasMAAG clusters enabled investigation of highly related plasmids from different

414    samples. Cluster *nneighs_76* was selected for more in-depth analysis. Plasmid bins from the

415    *nneighs_76* cluster reconstructed a 90 kb region from the plasmid *pAsa4c*, which is reported

416    to be hosted by *Aeromonas salmonicida subsp.* (42) (**Figure 5B**). Despite representing highly

417    overlapping regions of the same accession, bins from the *nneighs_76* cluster exhibited varying

418    degrees of contig fragmentation. For instance, the bin from sample 2 was composed of 20

419    contigs, whereas the bin from sample 1 consisted of 10 contigs, which could be explained by

420    the difference in the contig abundance. We then explored the relationship of the plasmid bins

421    using the alignments from the AAG (**Figure 5B**). We also found that some bins in *nneighs_76*

422    contains contigs that did not align to *pAsa4c*. Some of these unaligned contigs were found in

423    multiple bins and were syntenic across bins aligned to each other, suggesting that we found

424    true plasmid variation, and not an error in binning (**Figure 5B**). Using synteny, we could find

425    four approximate locations on the reference sequence where these contigs belonged to. Three

426    of four regions had hallmarks of recombination hotspots, including an ISAs2 insertion site, a

427    known conjugative block and a segment with distinct GC content (42) (**Figure 5B**). Furthermore,

428    14 of 19 contigs not mapping to *pAsa4c* aligned to plasmid accessions reported to be hosted

429    by organisms from the *Aeromonas* genus. Additionally, PlasMAAG clusters, together with the

430    assembly-alignment graph, enable the exploration of diversity among similar plasmids across

431    samples without PLSDB support. As an example, bins from the plasmid cluster *nneighs_416*

432    exhibited a high degree of sequence similarity despite variations in contig fragmentation

433    (**Figure 5C**). PlasMAAG facilitates the tracking of highly similar plasmids across different

434    environments, allowing for the capture of their composition variations.

**DISCUSSION**

436    Plasmids are pivotal in horizontal gene transfer, playing an influential role in shaping microbial

437    communities. Their prevalence across microbial ecosystems highlights their importance, yet

438  studying plasmids from environmental samples has been challenging due to their dynamic
439  and unstable composition. This limitation has hindered efforts to bin and identify plasmids
440  accurately, despite their abundance. The recent decrease in sequencing costs has significantly
441  increased the availability of metagenomic samples, presenting an unprecedented opportunity
442  to uncover plasmid diversity. However, the challenges of plasmid binning emphasize the need
443  for a robust and broad-range plasmid binning method.

444  In this study, we introduced PlasMAAG, a novel deep learning-based framework for
445  metagenomic binning of both plasmids and organisms. PlasMAAG leverages a unique feature
446  we developed—assembly-alignment graphs—which enables the aggregation of assembly
447  graphs across multiple samples. This advancement allows PlasMAAG overcome traditional
448  limitations associated with single-sample plasmid assemblers.

449  PlasMAAG outperformed SCAPP, the current state-of-the-art plasmid assembler, on both
450  synthetic and real datasets, delivering superior results for plasmid binning while being
451  significantly faster. Besides producing more plasmid bins, the set of candidate plasmids
452  produced by PlasMAAG achieved a more balanced trade-off between precision and recall,
453  enabling a broader characterization of metagenomic samples. Notably, PlasMAAG's capability
454  to bin all sequences, including plasmids and organisms, offers a comprehensive approach to
455  metagenomic analysis. PlasMAAG achieves organism binning results that are comparable to
456  leading organism binners on synthetic datasets while demonstrating superior performance in
457  understudied, real-world environments.

458  PlasMAAG's holistic approach enables integrated studies, such as the exploration of plasmid-
459  host associations. Using its comprehensive binning capabilities, we gathered correlation-
460  abundance-based evidence for 773 plasmid-host associations, with only 7% previously
461  reported in the PLSDB database. Furthermore, PlasMAAG's assembly-alignment graph-based
462  clustering revealed intra-plasmid variation across samples, enabling the study of plasmid
463  sequence variation across environments.

464  We demonstrated that geNomad plasmid predictions were significantly enhanced when
465  aggregated across PlasMAAG communities, underscoring the value of binning for refining
466  plasmid sequence identification. However, we also saw that geNomad was inaccurate when
467  applied to understudied environments, as validated by experimental paired metaplasmidomics,

468 long, and short-read data. These discrepancies highlight the need for more robust plasmid

469 sequence identifiers capable of handling complex or uncharted environments.

470 The success of PlasMAAG is largely attributable to the assembly-alignment graph, a feature

471 that complements assembly graph signals across samples in a multi-sample framework. This

472 innovation not only enhances binning accuracy but also facilitates the inference of

473 compositional similarities between samples. Moreover, assembly-alignment graphs also

474 improve binning of contigs from the same sample, through indirect links to contigs of other

475 samples.

476 Another notable innovation in PlasMAAG is its use of contrastive loss to integrate traditional

477 binning features like k-mer composition and contig abundances, with the assembly-alignment

478 graph. This approach could be extended to incorporate other graph-like data in the binning

479 process, such as Hi-C data. As sequencing technologies advance and contigs become

480 decreasingly fragmented, particularly in long-read datasets, the utility of using cross-sample

481 alignments to bridge gaps in the assembly graphs will grow, covering larger genome fractions

482 and providing richer insights.

483 Despite the advances introduced by PlasMAAG, plasmid binning remains a significant

484 challenge, as evidenced by the lack of groundbreaking plasmid binners in recent years. This

485 underscores the necessity of innovative approaches, like PlasMAAG, that address the

486 complexities of plasmid diversity and recombination. By enabling the study of plasmids

487 alongside organisms from highly complex samples, PlasMAAG expands our ability to explore

488 microbial communities comprehensively. Its focus on plasmids—an often-overlooked but

489 critical component of microbial ecosystems—enhances our understanding of their role in

490 horizontal gene transfer and microbial community dynamics.

491 In conclusion, PlasMAAG represents a step forward in plasmid and organism binning from

492 metagenomic samples. By incorporating assembly-alignment graphs and contrastive learning,

493 it addresses longstanding challenges in plasmid binning while providing a framework for

494 studying plasmid-host associations and microbial community dynamics. PlasMAAG offers a

495 valuable tool for advancing our understanding of microbial ecosystems, with implications for

496 environmental microbiology, public health, and biotechnology. PlasMAAG

497 **MATERIAL AND METHODS**

498 **Overview of PlasMAAG**

499 The inputs to the PlasMAAG pipeline are a set of reads per sample. Reads are assembled per

500 sample with *metaSPAdes* v3.15.5 (43) creating an assembly graph and contigs for each sample.

501 The contigs across all samples are concatenated together to create the contig catalogue. Reads

502 are mapped to the catalogue with *minimap2* v2.24 (44) and *samtools* v1.18 (45) , creating per-

503 sample BAM files. The alignment graph is generated by aligning the contigs across samples

504 with NCBI *blast* 2.15.0 (46). The assembly- and alignment graphs are merged into the

505 assembly-alignment graph (AAG). *Fastnode2vec* v0.05 (37), an optimized version of node2vec,

506 is used to embed local AAG context of each contig into an embedding space, from which

507 communities of contigs with similar embeddings are extracted. The k-mer composition and

508 abundance features of contigs are embedding using a variational autoender (VAE), where an

509 additional loss term is added which penalizes distance between contigs of the same

510 community. Using the VAE embedding, communities are expanded, merged, and purified. The

511 *geNomad* (38) tool is used to separate plasmid from non-plasmid contigs: Communities of

512 plasmid contigs are extracted as separate bins, whereas the rest contigs are extracted in bins

513 using a clustering algorithm.

514 **Benchmark datasets**

515 We based our benchmark dataset on the existing CAMI2 short-read human microbiome toy

516 dataset, but had to modify the dataset to allow benchmarking of plasmids: First, the original

517 dataset did not provide assembly graphs, so we assembled the reads and mapped the resulting

518 contigs back to the CAMI2 source genomes to determine their origin, using *minimap2* and

519 accepting hits with an identity > 97% and a query coverage > 90%. Because this approach

520 initially led to many unmapped or ambiguously mapping contigs, we re-simulated the reads

521 using *wgsim* (47) with zero sequencing errors, then assembled each sample using *metaSPAdes*

522 without the use of error correction. Second, CAMI2 considered plasmids to be part of their

523 cellular host genome with the same abundance, which would inhibit our abundance-based

524 binning approach. We changed so that plasmids were separate genomes with an abundance

525 proportional to host abundance times a Gaussian random variable, as done in (18). Finally,

526 CAMI2 did not contain reads simulated from across the edges of the underlying circular

527 sequences, which prevents assembly graph cycles and hobbles graph peeling-based

528 approaches like that used by SCAPP. We made sure to include such reads.

529 **Assembly graph edge weighting**

530 Assembly graphs were extracted from the *assembly_graph_after_simplification.gfa* file

531 generated from metaSPAdes and converted into a NetworkX v3.4.2 (48) directed graph, with

532 contigs represented as nodes, and links between segments in contigs represented as edges.

533 To enrich the assembly graph signal for binning, graph edges were weighted with the

534 *normalized linkage* metric, which is dependent on the number of links established between

535 any segments from each pair of contigs, normalized by the length of the contigs. For a pair of

536 contigs $c^i$, $c^j$, the number of links connecting those contigs $n\_links_{ij}$, and the contig lengths $l^c$,

537 normalized linkage is:

$$normalized\ linkage_{c^i c^j} = \frac{n\_links}{\min(l^{c^i}, l^{c^j})}$$

538

539 **Alignment graph edge weighting**

540 After assembly, contigs shorter than 2000 bp were discarded as done in (26). Contigs were

541 aligned all against all using NCBI blast using *blastn* command with *-perc_identity 95*, only

542 keeping between-sample hits, alignment identity ≥ 98.0% and an alignment ≥ 500 bp. We also

543 removed alignments between sequences that contained large sections that did not align due

544 to sequence diversity, as we wanted the alignments to represent shared sequences across

545 samples. The remaining set of alignments after filtering was defined as 'restrictive' alignments.

546 From the aligments we created an alignment graph with contigs as nodes and alignments as

547 edges. Edges were weighted with the *normalized alignment* metric to reflect the alignment

548 certainty. For a pair of contigs $c^i$, $c^j$, alignment identity *id*, alignment length *L*, and contig length

549 $l^c$:

$$normalized\ alignment_{c^i c^j} = \frac{id}{100} \frac{\min(L, l^{c^i}, l^{c^j})}{\min(l^{c^i}, l^{c^j})}$$

550

551 **Assembly-alignment graph community extraction with node2vec**

552 Assembly and alignment graphs share no edges, since their edges connect only within-sample

553 and between-sample contigs, respectively. This allowed us to trivially merge the graphs by

554 adding the edges from one graph into the other, thus creating the AAG. To extract

555 communities from the AAG, we first ran *fastnode2vec* on the AAG to obtain contig embeddings.

556 We created a new graph by linking contigs within a cosine distance of 0.1 in embedding space,

557 after which we defined each connected component to be a contig community. We optimized

558 the *fastnode2vec* hyperparameters and clustering radius to generate pure communities at

559 genome level, running a small grid search over the re-simulated CAMI2 Airways dataset. The

560 embedding dimensions, walk length, number of walks, window size, p, and q parameters from

561 fastnode2vec were set to 32, 10, 50, 10, 0.1, and 2.0. The embedding clustering cosine distance

562 radius was set to 0.1.

563 **Contrastive-VAMB for community merging and expansion**

564 Contrastive-VAMB is a variation of the original VAMB model, with a modification on the loss

565 function to account for the communities extracted from the fastnode2vec embeddings.

566 Contrastive-VAMB is composed of an encoder, latent representation layer m, and a decoder.

567 Each contig represented by the concatenation of the contig co-abundances along samples $\mathbf{A}_{in}$,

568 the tetranucleotide frequencies $\mathbf{T}_{in}$, and the unnormalized contig abundances $\mathbf{C}_{in}$ and passed

569 to the encoder. The encoder projects the contigs into a latent normal $N(\mu, I)$ distribution

570 parametrized by the m layer, from which the decoder samples. The decoder is optimized to

571 reconstruct $\mathbf{A}_{in}$, $\mathbf{T}_{in}$, and $\mathbf{C}_{in}$ from the instances sampled from $N(\mu, I)$, decrease the latent cosine

572 distance between contigs with closely related node2vec graph embeddings, and decrease the

573 deviance between the latent normal distribution $N(\mu, I)$ parametrized by the $\mu$ layer and the

574 standard normal distribution used as prior $N(0, I)$.

575 **Loss functions**

576 The contrastive-VAMB loss can be decomposed in three terms: reconstruction loss, contrastive

577 loss, and regularization loss. The reconstruction loss ($L_{rec}$) penalizes the reconstruction error of

578 $\mathbf{A}_{in}$, $\mathbf{T}_{in}$, and $\mathbf{C}_{in}$. In the same way than the original VAMB reconstruction loss, cross entropy (CE)

579 and sum of squared errors (SSE) losses were set for the reconstruction of the $\mathbf{A}_{in}$ and $\mathbf{T}_{in}$,

580 respectively, whereas SSE loss was set for the $\mathbf{C}_{in}$ loss. These three terms are weighted with

581 hyperparameters $w_A$, $w_T$, and $w_C$.

582
$$L_{rec} = w_A CE(A_{in}, A_{out}) + w_T SSE(T_{in}, T_{out}) + w_C SSE(C_{in}, C_{out})$$

583    The contrastive loss ($L_{contr}$) penalizes the cosine distance between the VAMB latent

584    representations of the contigs and contigs highly related in node2vec embedding space, when

585    such cosine distance overcomes a predefined margin $m$, $m$ being a hyperparameter. For a

586    contig $ci$ and highly related fastnode2vec embedding space contigs $H^{ci}$ ={$n_{0,...}$, $n_n$}:

587    $$L_{contr} = \max \left( \frac{\sum_{n_i \in H^{ci}} cosine\ distance(\mathrm{m}^{ci}, \mathrm{m}^{n_i})}{|H^{ci}|} - m, 0 \right)$$

588    The regularization loss ($L_{reg}$) penalizes the deviance between the latent normal distribution N($\mu$,

589    I) parametrized by the $\mu$ layer, and the standard normal distribution used as prior N(0, I) with

590    the Kullback-Leibler divergence, which since the standard deviation is set to 1, simplifies to:

591    $$L_{reg} = \frac{1}{2} + \sum \mu^2$$

592    Finally, the model total loss (L) was aggregated with weighting hyperparameters $w_{Lreg}$, and

593    $w_{Lcontr}$:

594    $$L = L_{rec} + w_{L_{reg}} L_{reg} + w_{L_{contr}} + L_{contr}$$

**Clustering plasmid/organism candidates with geNomad**

596    Two parallel strategies were implemented to cluster the latent space tailored to extract

597    plasmids and non-plasmids, respectively. The plasmid clustering strategy is composed of two

598    phases: clustering community-based and clustering iterative medoid based, both based on

599    latent space cosine distances. The clustering community-based works in five steps

600    (**Supplementary Figure 9**): (1) for each community extracted from the node2vec embeddings,

601    link contigs belonging to the same community, and remove links between contigs with a VAE

602    embedding cosine distance > 0.2. (2) Contigs are recruited into the community if within 0.01

603    cosine distance to any community member. If the recruited contig is part of a community, the

604    two communities are merged. (3) The expanded communities are extracted from the latent

605    space as bins, and remaining contigs are clustered with the original medoid based VAMB

606    clustering algorithm, (4) self-circularized contigs are extracted based upon mapping read-pairs

607    where mates map to opposite contig ends within 50 bps from the contig end, and extracted

608    from the clusters, (5) Plasmid score is defined for each cluster by aggregating the geNomad

609    plasmid contig scores with a contig length weighted mean, defining plasmid candidates when

610    cluster scores are larger than the defined threshold. When geNomad plasmid threshold is

611   larger than 0.5, a fixed geNomad plasmid threshold of 0.5 is applied to the circular contigs,

612   accounting for the circular evidence relatable to plasmids. The non-plasmid clustering strategy

613   consists in 2 steps: (1) Cluster the VAMB-latent space with the iterative medoid clustering

614   algorithm from VAMB. (2) Extract contigs belonging to any plasmid candidate cluster defined

615   by the plasmid prone clustering strategy.

616   **Binning benchmarking – CAMI2 reassembled**

617   We compared the plasmid and organism binning performance of PlasMAAG, VAMB v4.1.3,

618   MetaBAT2 v2.12.1, SemiBin2 v2.1.0, Comebin v1.0.4, MetaDecoder v1.0.19, and SCAPP v0.1.4

619   over the re-simulated CAMI2 datasets. Binning performance was evaluated in terms of

620   genomes recovered with precision ≥ 95% and recall ≥ 90%, so-called "NC genomes". Since

621   PlasMAAG, and VAMB, MetaBAT2, SemiBin2, Comebin, MetaDecoder perform the binning

622   after assembling the contigs, precision and recall of the bins were obtained from the contig

623   references, using BinBencher v0.3.0 (49). On the other hand, SCAPP and MetaPlasmidSPAdes

624   v3.15.3 assemble their own contigs. Here, we produced a ground truth by aligning the output

625   bins to the origin genomes using NCBI blast 2.15.0 accepting hits with an identity > 97% and

626   a query coverage > 90%, after which we benchmarked using BinBencher.

627   **Sample benchmarking CAMI2 reassembled**

628   Precision, recall, and F1 was computed for each set of plasmid candidates, reflecting the

629   plasmid characterisation at the sample level, not at the bin level. Given a sample (*s*), a set of

630   plasmid candidates (*candidates*), binning precision and binning recall thresholds (*pre*, *rec*), and

631   the set of true plasmids present in the sample (*plasmids*):

632   $$Sample\ precision_{candidates,pre,rec} = \frac{\#\ candidates > (pre, rec)}{\#\ candidates}$$

633   $$Sample\ recall_{candidates,plasmids,pre,rec} = \frac{\#\ candidates > (pre, rec)}{\#\ plasmids}$$

634   Enabling to compare the number of bins classified as plasmid, compared to the total number

635   of plasmid genomes at specific binning precision and recall thresholds.

636   **Hospital sewage samples sequence datasets**

637 Two datasets were used in this study to assess the quality of plasmid binning. Urban sewage
638 samples (UWS) samples were collected from comparable UWSs from Denmark and Spain
639 located in Odense and Santiago de Compostela, as previously described (41). In this study,
640 only hospital sewage samples from each location were used. Sewage samples were collected
641 in the winter and summer of 2018 using ISCO automatic samplers for 24-hour flow (50 mL per
642 5 min) in Denmark, while 24-hour-time proportional samples in SP (mixing hourly samples
643 according to flow information) (**Supplementary Table 8**). Three replicates per site and season
644 were collected on three consecutive days without rain events. All samples were initially cooled
645 with ice on-site, then 100 mL of each sample was centrifugated at 10,000 g for 8 min at 4 °C
646 in the laboratory. After removing supernatant, pellets were resuspended in 20 % of glycerol
647 stock to reach a final volume of 10 mL for storage at −80 °C. In total, environmental DNA was
648 extracted from all samples using NucleoSpin Soil kit (Macherey & Nagel, Dürein, DE) using
649 500µl of glycerol stock material for direct shotgun metagenomic using Illumina NovaSeq using
650 2x150bp paired-end mode (all samples) and PacBio Sequel2e (5 samples from Denmark).
651 PacBio libraries were built from the same DNA extracts using libraries using SMRTbell express
652 template 2.0 kit and Sequel II Binding Kit 3.2 (Pacific Bioscience, CA, USA) and barcoded using
653 SMRTbell Barcoded Adapter Plate 3.0 (Pacific Bioscience, CA, USA). Two libraries per 8M
654 SMRTcell (Pacific Bioscience, CA, USA) were pooled and sequenced on a PacBio Sequel2e
655 instrument at University of Copenhagen.

656 For plasmids enriched samples, we used specific methods to deplete non-plasmid DNA as
657 described previously (50, 51). Briefly, hospital sewage samples were pretreated by filtration,
658 vortex and sonication and resuspended in TE buffer. Afterwards, a pre-lysis cocktail of cell-wall
659 degrading enzymes: lysozyme, mutanolysin, and lysostaphin was used to facilitate lysis of
660 Gram-positive bacteria during alkaline lysis. Pre-lysis was followed by alkaline lysis to remove
661 chromosomal DNA (52), followed by Plasmid-Safe™ ATP-Dependent DNase (Lucigen, UK)
662 digestion. Plasmid-Safe DNase will digest any fragments of dsDNA with open 3' or 5' termini,
663 hence removing fragmented chromosomal DNA. The purified plasmid DNA was then quality-
664 checked, libraries prepared and sequenced on an Illumina NextSeq platform with a v2.5
665 sequencing kit (Illumina, San Diego, CA, USA) in paired-end mode.

666 **Binning benchmarking – hospital sewage**

667  We compared the binning performance of PlasMAAG, VAMB, MetaBAT2, SemiBin2, Comebin,
668  MetaDecoder, metaplasmidSPAdes, and SCAPP over the 5 hospital sewage samples.
669  Performance evaluation was based on the long-read sequences generated from the same
670  samples and defined by the long-read contigs recovered with precision ≥ 95% and recall ≥
671  90%, so-called "NC long-read assemblies". To evaluate the overall binning performance, the
672  entire set of long-read contigs was used to build the reference. Whereas to evaluate the
673  plasmid binning performance, only the long-read contigs either circular or with
674  metaplasmidomics reads coverage > 50% were used to build the reference. To build the
675  references, we mapped the short-read contigs to either set of long-read contigs to determine
676  their origin, using minimap2 v2.24 and accepting hits with an identity > 97% and a query
677  coverage > 90%, and used Binbencher for the benchmarking. To account for plasmid circularity,
678  2 copies of each long-read contig were concatenated before mapping the short-read contigs.
679  adovNC organisms were estimated with CheckM2 v0.1.3.

680  **Host-plasmid and intra-plasmid diversity exploration**

681  PlasMAAG was used to bin the contig sequences from 24 hospital sewage samples from
682  hospitals in Spain. PlasMAAG bins were aggregated into PlasMAAG clusters and classified as
683  plasmids if the aggregated geNomad plasmid score exceeded 0.75, defining them as plasmid
684  clusters. Only plasmids clusters with more than 150 kb were considered for the host-plasmid
685  association. Organism's bin quality was estimated with CheckM2 v0.1.3, and only high-quality
686  (completeness ≥ 70% and precision ≥ 90%) (HQ) bins were kept. GTDBtk v2.4.0 (53) was used
687  to estimate taxonomy for the HQ bins, with cluster taxonomy assigned based on majority vote.
688  Abundance correlation analysis was only conducted for plasmids and organism's clusters with
689  non-zero abundance over at least 18 overlapping samples. Spearman correlation coefficients
690  and p-values were computed using *scipy.stats.spearmanr*. To account for multiple testing, p-
691  values were corrected using the Benjamini-Hochberg (FDR) correction implemented in the
692  *statsmodels.stats.multitest.multipletests* package. Plasmid cluster hosts were inferred from
693  PLSDB when aligning to any PLSDB entry with >80% identity and >80% coverage. Functional
694  annotations of contigs were performed with *anvi'o* v8 software, using the 'anvi-run-workflow
695  -w contigs' command.

696  **Resource usage**

697 We evaluated computational resource usage of all methods using the Airways CAMI2 re-

698 assembled dataset and five samples from the hospital sewage dataset. For the Airways dataset,

699 PlasMAAG used 46 minutes, 8 threads, and 16 GB of RAM. In contrast, SCAPP, excluding the

700 BAM file generation step, took 192 minutes, utilized 16 threads, and required 24 GB of RAM

701 (**Supplementary Table 5**). Among the other binners, PlasMAAG was slower than VAMB,

702 MetaDecoder, and MetaBAT2. For example, VAMB completed the task in just 8 minutes while

703 using 8 threads and 16 GB of RAM. However, we observed a different trend when evaluating

704 performances on the five hospital sewage samples. When accounting for the additional steps

705 of read assembly and read mapping required to compute abundances, PlasMAAG exhibited

706 similar runtimes to most binners, except for SCAPP, which required significantly more time.

707 Specifically, PlasMAAG took 3,575 minutes, VAMB took 3,435 minutes, ComeBin required 4,911

708 minutes, and metaplasmidSPAdes took 4,430 minutes (**Supplementary Table 6**). In contrast,

709 SCAPP required 116,965 minutes—32 times longer than PlasMAAG. This difference in runtime

710 remained consistent even when excluding the read assembly steps (**Supplementary Table 6**).

711

712 **DATA AVAILABILITY**

713 Reads, contigs, and contig annotations for the re-assembled CAMI2 datasets are available here:

714 https://erda.ku.dk/archives/826fe4d8889f88db2ec20058f9eaa015/published-archive.html and

715 https://erda.ku.dk/archives/fb2c6dd2a8e002becb58233bd4388f7c/published-

716 archive.html. The metagenomic short reads, metaplasmidomic short reads, and metagenomic long

717 reads from the 5 Danish hospital sewage samples, as well as the metagenomic short reads from the

718 24 Spanish hospital sewage samples, are available in the European Nucleotide Archive under

719 BioProject PRJEB85938, whereas the assemblies for all samples are available here:

720 https://erda.ku.dk/archives/e87f0d5e12ca4c1204379d4932c3ae59/published-

721 archive.html (**Supplementary Table 8**).

722

723 **SUPPLEMENTARY DATA**

724 Supplementary Data are available online.

725

**AUTHOR CONTRIBUTIONS**

S.R. and S.J.S. conceived the study. S.R., J.N.N, and P.P.L guided the analysis. P.P.L. developed PlasMAAG, wrote the software, and performed the analysis. Additionally, J.N.N., and L.S.D. also wrote the software. J.N.N. also performed analyses. S.J.S, J.N, M.P.A, and L.J.J. provided input for the analysis. I.K and J.N generated the sewage data. P.P.L., J.N.N., and S.R. wrote the manuscript with contributions from all co-authors. All authors read and approved the final version of the manuscript.

**ACKNOWLEDGEMENTS**

**CONFLICT OF INTEREST**

S.R. is the founder and owner of the Danish company BioAI and has performed consulting for Sidera Bio ApS. The remaining authors declare no conflict of interest.

**CODE AVAILABILITY**

PlasMAAG is freely available at https://github.com/RasmussenLab/vamb/tree/vamb_n2v_asy/workflow_PlasMAAG.

**REFERENCES**

752    1. Schwengers,O., Barth,P., Falgenhauer,L., Hain,T., Chakraborty,T. and Goesmann,A. (2020)
753        Platon: identification and characterization of bacterial plasmid contigs in short-read
754        draft assemblies exploiting protein sequence-based replicon distribution scores.
755        *Microb Genom*, **6**.

756    2. Rodríguez-Beltrán,J., DelaFuente,J., León-Sampedro,R., MacLean,R.C. and San Millán,Á. (2021)
757        Beyond horizontal gene transfer: the role of plasmids in bacterial evolution. *Nat. Rev.*
758        *Microbiol.*, **19**, 347–359.

759    3. Sherratt,D.J. (1974) Bacterial plasmids. *Cell*, **3**, 189–195.

760    4. San Millan,A. and MacLean,R.C. (2017) Fitness Costs of Plasmids: a Limit to Plasmid
761        Transmission. *Microbiol Spectr*, **5**.

762    5. Thomas,C.M. and Nielsen,K.M. (2005) Mechanisms of, and Barriers to, Horizontal Gene
763        Transfer between Bacteria. *Nat. Rev. Microbiol.*, **3**, 711–721.

764    6. Norman,A., Hansen,L.H. and Sørensen,S.J. (2009) Conjugative plasmids: vessels of the
765        communal gene pool. *Philos. Trans. R. Soc. Lond. B Biol. Sci.*, **364**, 2275–2289.

766    7. Zhang,T., Zhang,X.-X. and Ye,L. (2011) Plasmid metagenome reveals high levels of antibiotic
767        resistance genes and mobile genetic elements in activated sludge. *PLoS One*, **6**, e26041.

768    8. Stockdale,S.R. and Hill,C. (2023) Incorporating plasmid biology and metagenomics into a
769        holistic model of the human gut microbiome. *Curr. Opin. Microbiol.*, **73**, 102307.

770    9. NCBI Reference Sequence (RefSeq) Database, Release 226, September 9, 2024.

771    10. Hugenholtz,P. and Tyson,G.W. (2008) Microbiology: metagenomics. *Nature*, **455**, 481–483.

772    11. Handelsman,J. (2004) Metagenomics: application of genomics to uncultured
773        microorganisms. *Microbiol. Mol. Biol. Rev.*, **68**, 669–685.

774    12. Yu,M.K., Fogarty,E.C. and Eren,A.M. (2024) Diverse plasmid systems and their ecology across
775        human gut metagenomes revealed by PlasX and MobMess. *Nat Microbiol*, **9**, 830–847.

776    13. Finks,S.S. and Martiny,J.B.H. (2023) Plasmid-encoded traits vary across environments. *MBio*,
777        **14**, e0319122.

778    14. Pesesky,M.W., Tilley,R. and Beck,D.A.C. (2019) Mosaic plasmids are abundant and unevenly
779        distributed across prokaryotic taxa. *Plasmid*, **102**, 10–18.

780    15. Bouchot,J.-L., Trimble,W.L., Ditzler,G., Lan,Y., Essinger,S. and Rosen,G. (2014) Advances in
781        machine learning for processing and comparison of metagenomic data. In
782        *Computational Systems Biology*. Elsevier, pp. 295–329.

783    16. Rodríguez-Beltrán,J., Tourret,J., Tenaillon,O., López,E., Bourdelier,E., Costas,C., Matic,I.,
784        Denamur,E. and Blázquez,J. (2015) High recombinant frequency in extraintestinal
785        pathogenic Escherichia coli strains. *Mol. Biol. Evol.*, **32**, 1708–1716.

786    17. Fernandez-Lopez,R., Redondo,S., Garcillan-Barcia,M.P. and de la Cruz,F. (2017) Towards a
787        taxonomy of conjugative plasmids. *Curr. Opin. Microbiol.*, **38**, 106–113.

788 18. Pellow,D., Zorea,A., Probst,M., Furman,O., Segal,A., Mizrahi,I. and Shamir,R. (2021) SCAPP:
789     an algorithm for improved plasmid assembly in metagenomes. *Microbiome*, **9**, 144.

790 19. Antipov,D., Raiko,M., Lapidus,A. and Pevzner,P.A. (2019) Plasmid detection and assembly in
791     genomic and metagenomic data sets. *Genome Res.*, **29**, 961–968.

792 20. Rozov,R., Brown Kav,A., Bogumil,D., Shterzer,N., Halperin,E., Mizrahi,I. and Shamir,R. (2017)
793     Recycler: an algorithm for detecting plasmids from de novo assembly graphs.
794     *Bioinformatics*, **33**, 475–482.

795 21. Wajid,B. and Serpedin,E. (2012) Review of general algorithmic features for genome
796     assemblers for next generation sequencers. *Genomics Proteomics Bioinformatics*, **10**,
797     58–73.

798 22. Ayling,M., Clark,M.D. and Leggett,R.M. (2020) New approaches for metagenome assembly
799     with short reads. *Brief. Bioinform.*, **21**, 584–594.

800 23. Pellow,D., Mizrahi,I. and Shamir,R. (2020) PlasClass improves plasmid sequence
801     classification. *PLoS Comput. Biol.*, **16**, e1007781.

802 24. Lander,E.S. and Waterman,M.S. (1988) Genomic mapping by fingerprinting random clones:
803     a mathematical analysis. *Genomics*, **2**, 231–239.

804 25. Quince,C., Nurk,S., Raguideau,S., James,R., Soyer,O.S., Summers,J.K., Limasset,A., Eren,A.M.,
805     Chikhi,R. and Darling,A.E. (2021) STRONG: metagenomics strain resolution on assembly
806     graphs. *Genome Biol.*, **22**, 214.

807 26. Nissen,J.N., Johansen,J., Allesøe,R.L., Sønderby,C.K., Armenteros,J.J.A., Grønbech,C.H.,
808     Jensen,L.J., Nielsen,H.B., Petersen,T.N., Winther,O., *et al.* (05/2021) Improved
809     metagenome binning and assembly using deep variational autoencoders. *Nat.*
810     *Biotechnol.*, **39**, 555–560.

811 27. Kang,D.D., Li,F., Kirton,E., Thomas,A., Egan,R., An,H. and Wang,Z. (2019) MetaBAT 2: an
812     adaptive binning algorithm for robust and efficient genome reconstruction from
813     metagenome assemblies. *PeerJ*, **7**, e7359.

814 28. Wu,Y.-W., Simmons,B.A. and Singer,S.W. (2016) MaxBin 2.0: an automated binning
815     algorithm to recover genomes from multiple metagenomic datasets. *Bioinformatics*, **32**,
816     605–607.

817 29. Pan,S., Zhao,X.-M. and Coelho,L.P. (2023) SemiBin2: self-supervised contrastive learning
818     leads to better MAGs for short- and long-read sequencing. *bioRxiv*,
819     10.1101/2023.01.09.523201.

820 30. Pan,S., Zhu,C., Zhao,X.-M. and Coelho,L.P. (2022) A deep siamese neural network improves
821     metagenome-assembled genomes in microbiome datasets across different
822     environments. *Nat. Commun.*, **13**, 2326.

823 31. Líndez,P.P., Johansen,J., Kutuzova,S., Sigurdsson,A.I., Nissen,J.N. and Rasmussen,S. (2023)
824     Adversarial and variational autoencoders improve metagenomic binning. *Commun.*
825     *Biol.*, **6**, 1073.

32. Lamurias,A., Sereika,M., Albertsen,M., Hose,K. and Nielsen,T.D. (2022) Metagenomic binning with assembly graph embeddings. *Bioinformatics*, **38**, 4481–4487.

33. Kutuzova,S., Piera,P., Nor Nielsen,K., Olsen,N.S., Riber,L., Gobbi,A., Forero-Junco,L.M., Dougherty,P.E., Westergaard,J.C., Christensen,S., *et al.* (2024) Binning meets taxonomy: TaxVAMB improves metagenome binning using bi-modal variational autoencoder. *bioRxiv*, 10.1101/2024.10.25.620172.

34. Johansen,J., Plichta,D., Nissen,J.N., Jespersen,M.L., Shah,S.A., Deng,L., Stokholm,J., Bisgaard,H., Nielsen,D.S., Sørensen,S., *et al.* (2021) Genome binning of viral entities from bulk metagenomics data Genomics.

35. Kutuzova,S., Nielsen,M., Piera,P., Nissen,J.N. and Rasmussen,S. (2024) Taxometer: Improving taxonomic classification of metagenomics contigs. *Nat. Commun.*, **15**, 8357.

36. Mattock,J. and Watson,M. (2023) A comparison of single-coverage and multi-coverage metagenomic binning reveals extensive hidden contamination. *Nat. Methods*, **20**, 1170–1173.

37. Abraham,L. (2020) louisabraham/fastnode2vec version-0.0.5 Zenodo.

38. Camargo,A.P., Roux,S., Schulz,F., Babinski,M., Xu,Y., Hu,B., Chain,P.S.G., Nayfach,S. and Kyrpides,N.C. (2023) Identification of mobile genetic elements with geNomad. *Nat. Biotechnol.*, 10.1038/s41587-023-01953-y.

39. Mallawaarachchi,V., Wickramarachchi,A. and Lin,Y. (2020) GraphBin: refined binning of metagenomic contigs using assembly graphs. *Bioinformatics*, **36**, 3307–3313.

40. Mallawaarachchi,V., Wickramarachchi,A., Xue,H., Papudeshi,B., Grigson,S.R., Bouras,G., Prahl,R.E., Kaphle,A., Verich,A., Talamantes-Becerra,B., *et al.* (2024) Solving genomic puzzles: computational methods for metagenomic binning. *Brief. Bioinform.*, **25**.

41. Yu,Z., He,W., Klincke,F., Madsen,J.S., Kot,W., Hansen,L.H., Quintela-Baluja,M., Balboa,S., Dechesne,A., Smets,B., *et al.* (2024) Insights into the circular: The cryptic plasmidome and its derived antibiotic resistome in the urban water systems. *Environ. Int.*, **183**, 108351.

42. Tanaka,K.H., Vincent,A.T., Trudel,M.V., Paquet,V.E., Frenette,M. and Charette,S.J. (2016) The mosaic architecture of Aeromonas salmonicida subsp. salmonicida pAsa4 plasmid and its consequences on antibiotic resistance. *PeerJ*, **4**, e2595.

43. Nurk,S., Meleshko,D., Korobeynikov,A. and Pevzner,P.A. (05/2017) metaSPAdes: a new versatile metagenomic assembler. *Genome Res.*, **27**, 824–834.

44. Li,H. (2018) Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*, **34**, 3094–3100.

45. Li,H., Handsaker,B., Wysoker,A., Fennell,T., Ruan,J., Homer,N., Marth,G., Abecasis,G., Durbin,R. and 1000 Genome Project Data Processing Subgroup (2009) The Sequence Alignment/Map format and SAMtools. *Bioinformatics*, **25**, 2078–2079.

46. Altschul,S.F., Gish,W., Miller,W., Myers,E.W. and Lipman,D.J. (1990) Basic local alignment search tool. *J. Mol. Biol.*, **215**, 403–410.

865     47. Li,H. (2011) wgsim-Read simulator for next generation sequencing. *Github repository*.

866     48. Hagberg,D. and Pieter,J. Exploring network structure, dynamics, function using NetworkX".

867     49. Nissen,J.N., Lindéz,P.P. and Rasmussen,S. (2024) BinBencher: Fast, flexible and meaningful
868          benchmarking suite for metagenomic binning. *bioRxiv*, 10.1101/2024.05.06.592671.

869     50. Alanin,K.W.S., Jørgensen,T.S., Browne,P.D., Petersen,B., Riber,L., Kot,W. and Hansen,L.H.
870          (2021) An improved direct metamobilome approach increases the detection of larger-
871          sized circular elements across kingdoms. *Plasmid*, **115**, 102576.

872     51. Kav,A.B., Sasson,G., Jami,E., Doron-Faigenboim,A., Benhar,I. and Mizrahi,I. (2012) Insights
873          into the bovine rumen plasmidome. *Proc. Natl. Acad. Sci. U. S. A.*, **109**, 5452–5457.

874     52. Green,M.R. and Sambrook,J. (2012) Molecular Cloning 4th ed. Cold Spring Harbor
875          Laboratory Press, New York, NY.

876     53. Chaumeil,P.-A., Mussig,A.J., Hugenholtz,P. and Parks,D.H. (2022) GTDB-Tk v2: memory
877          friendly classification with the genome taxonomy database. *Bioinformatics*, **38**, 5315–
878          5316.