

1 Binning meets taxonomy: TaxVAMB improves
2 metagenome binning using bi-modal variational
3 autoencoder

Svetlana Kutuzova^{1,2}, Pau Piera², Knud Nor Nielsen^{2,3,6},
 Nikoline S. Olsen³, Leise Riber³, Alex Gobbi^{3,4},
 Laura Milena Forero-Junco³, Peter Erdmann Dougherty³,
 Jesper Cairo Westergaard³, Svend Christensen³,
 Lars Hestbjerg Hansen³, Mads Nielsen¹, Jakob Nybo Nissen^{2*†},
 Simon Rasmussen^{2,5*†}

¹² The Novo Nordisk Foundation Center for Basic Metabolic Research,
¹³ University of Copenhagen, Blegdamsvej 3A, Copenhagen, 2200,
¹⁴ Denmark.

¹⁵ ³Department of Plant and Environmental Sciences, University of
¹⁶ Copenhagen, Thorvaldsensvej 40, Copenhagen, 1871, Denmark.

¹⁷ European Food and Safety Authority (EFSA), Via Carlo Magno 1A,
¹⁸ 43126, Parma, Italy.

¹⁹ ²⁰ ⁵The Novo Nordisk Foundation Center for Genomic Mechanisms of Disease, Broad Institute of MIT and Harvard, Cambridge, 02142, MA

²¹ USA.
²² ²³ ⁶The Novo Nordisk Foundation Center for Biosustainability, Technical University of Denmark, Kemitorvet 220, Lyngby, 2800, Denmark.

24 *Corresponding author(s). E-mail(s): jakob.nissen@sund.ku.dk;
25 srasmuss@sund.ku.dk;

26 †These authors contributed equally to this work.

27 **Abstract**

28 A common procedure for studying the microbiome is binning the sequenced
29 contigs into metagenome-assembled genomes. Currently, unsupervised and self-
30 supervised deep learning based methods using co-abundance and sequence based
31 motifs such as tetranucleotide frequencies are state-of-the-art for metagenome
32 binning. Taxonomic labels derived from alignment based classification have not
33 been widely used. Here, we propose TaxVAMB, a metagenome binning tool
34 based on semi-supervised bi-modal variational autoencoders, combining tetra-
35 nucleotide frequencies and contig co-abundances with contig annotations returned
36 by any taxonomic classifier on any taxonomic rank. TaxVAMB outperforms all
37 other binners on CAMI2 human microbiome datasets, returning on average 40%
38 more near-complete assemblies than the next best binner. On real long-read
39 datasets TaxVAMB recovers on average 13% more near-complete bins and 14%
40 more species. When used in a single-sample setup, TaxVAMB on average returns
41 83% more high quality bins than VAMB. TaxVAMB bins incomplete genomes
42 drastically better than any other tool, returning 255% more high quality bins
43 of incomplete genomes than the next best binner. Our method has immediate
44 research and industrial applications, as well as methodological novelty which
45 can be translated to other biological problems with semi-supervised multimodal
46 datasets.

47 **Keywords:** metagenomics binning, variational autoencoders, semi-supervised
48 learning, deep learning, metagenomics, hierarchical loss

49 Shotgun metagenome sequencing is an accessible technology that enables high-
50 throughput analysis of complex microbial communities for both taxonomic profiling
51 and metagenome assembly tasks. The field is currently dominated by short-read (com-
52 monly 100–300 bp) technologies¹, however, recently long-read sequencing has gained
53 prominence, as it allows the recovery of even more individual genomes with higher
54 accuracy^{2–4}. When working with environmental samples in the absence of cultured
55 isolates, the assembled contigs are grouped together during the process of metagenome
56 binning⁵.

57 Most metagenome binning tools^{6–11} are based on analysing both contig composition,
58 commonly represented as k-mer frequencies vectors such as tetranucleotide
59 frequencies (TNFs)¹², and contig co-abundances across multiple samples. Besides the
60 information contained in contigs, some tools additionally rely on assembly graphs^{13–16},
61 codon usage¹⁷, GC content⁶, single-copy genes^{18–22} and contig-level taxonomy pro-
62 filing^{20,23–25}. Furthermore, ensemble tools leverage the binning results created by
63 multiple approaches^{26–28}. Most metagenome binning tools have been optimized for
64 short-read sequences and their performance on long-read datasets has not been thor-
65oughly evaluated. Recently, several tools such as GraphMB^{15?}, SemiBin2²¹ and
66 LRBinner²⁹ have been developed specifically for long-read sequencing data. In general,
67 the large amounts and complexity of metagenomic data make it a suitable application
68 for deep learning algorithms^{11,15,16,20–22,28}.

69 For the purpose of this study, we emphasise a rough distinction between the intrin-
70 sic features³⁰ derived purely from a given set of reads and their corresponding contigs
71 (k-mer frequencies, GC-content, co-abundances) and the annotation features that
72 require searching external databases (e.g. single-copy genes, taxonomic labels from
73 sequence alignment). A taxonomic label is an example of an annotation feature, which
74 can be extracted from a read or a contig using taxonomic profiling tools^{31–38}. These
75 annotations are often incomplete since not all the contigs can be successfully mapped
76 to a reference sequence. The annotation might also be biased towards better-studied
77 organisms that will be more prevalent in databases.

78 Recently, single-copy genes (SCGs) have been used as a key clustering feature
79 by the SemiBin2²¹ and Comebin²² methods. Traditionally, SCGs have been used for
80 evaluating metagenome assembled genomes (MAGs), as in the popular metagenomic
81 binning evaluation tools CheckM and CheckM2^{39,40}. While missing or duplicated
82 single-copy genes are indeed a strong signal of the MAG quality, one might be cau-
83 tious about using these both as input to binning and as an evaluation of the produced
84 bins. This turns the evaluation metric into a training target, an observation sometimes

85 referred to as Goodhart's law: "When a measure becomes a target, it ceases to be a
86 good measure" ⁴¹.

87 For instance, SemiBin2 performs self-supervised deep learning using only intrinsic
88 features (TNFs and co-abundances) followed by a reclustering step using SCGs as the
89 optimization target ²¹. The reclustering step includes a greedy search over multiple
90 possible clusterings, selecting clustering with the best contamination and completeness
91 values calculated based on the presence of single-copy genes. In principle, reclustering
92 does not have to be done on the embeddings produced by a deep learning model, but
93 could more efficiently be applied directly to the input features, foregoing the embed-
94 ding entirely. The authors of SemiBin2 reported that the performance of reclustering
95 the input features of simulated short-read datasets was only 17.7% lower than that
96 of reclustered embeddings. This suggests that reclustering, which introduces the eval-
97 uation metrics as an input to the model, is by far the most informative part of the
98 SemiBin2 algorithm. Therefore, investigation of binning performance with and without
99 single-copy genes is needed to provide unbiased benchmarks of different methods.

100 Incorporation of taxonomic information presents a computational challenge due to
101 its hierarchical nature. The taxonomic labels used to classify the hierarchical phylogeny
102 of microorganisms are organised into the seven classical taxonomic ranks from kingdom
103 to species. Lower taxonomic ranks provide more precise information about the contig
104 phylogenetic placement, but are more often mislabeled or missing. As demonstrated
105 in the Taxometer tool ⁴², a hierarchical loss allows training on the labels acquired on
106 all the taxonomic ranks (e.g. phylum or genus) without requiring annotations on a
107 particular taxonomic rank (e.g. species). Previously, both SemiBin ²⁰ and SolidBin ²⁵
108 used taxonomic labels to generate cannot-link constraints in the loss functions of self-
109 supervised deep learning algorithms but neither labels themselves nor their hierarchical
110 structure were a part of the training data.

111 A key feature of semi-supervised machine learning is that the models can be trained
112 using both annotated and un-annotated samples. Analogous to the standard unsuper-
113 vised VAEs, the semi-supervised multi-modal VAEs exhibit generative capabilities and

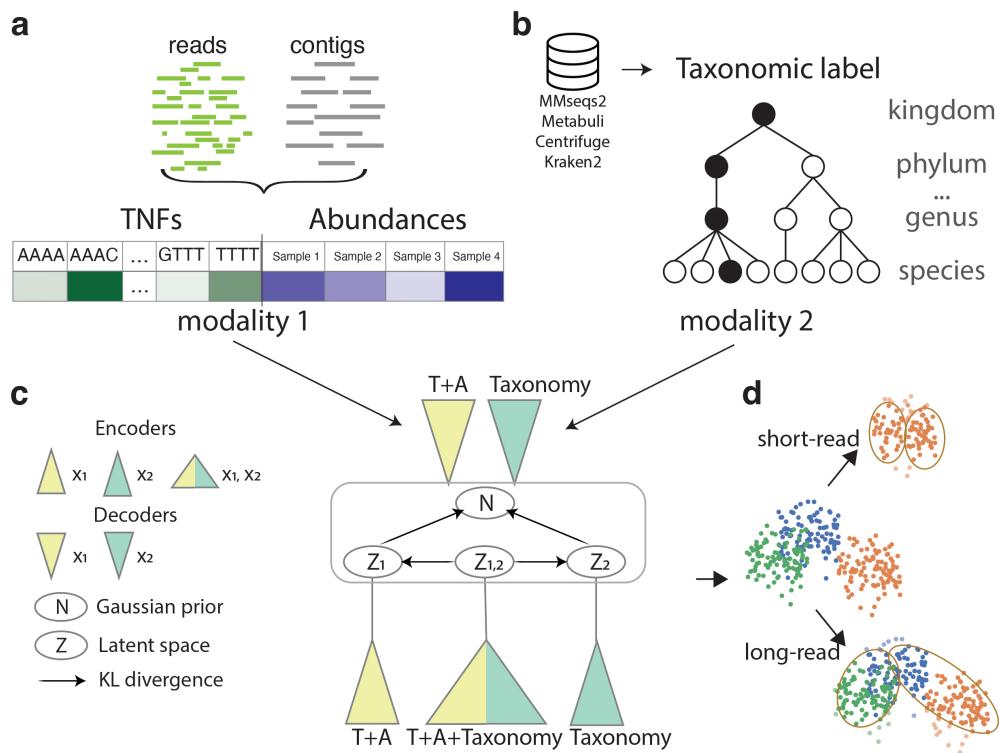


Fig. 1 TaxVAMB workflow. **a.** Tetranucleotide frequencies and contig abundances across samples are extracted from reads and their assemblies. **b.** Contigs are annotated with taxonomic labels by a taxonomic classifier, and the labels are refined by the Taxometer tool, resulting in higher quality annotations. The taxonomic label is represented by a binary vector where each element encodes a taxon. Taxa from all the taxonomic ranks are present. **c.** We consider a concatenated vector of TNFs and abundances to be the first modality and the taxonomy label to be the second modality. Bi-modal variational autoencoder is trained on the two modalities. For each sample, three observations are created: 1) modality one; 2) modality two; 3) a concatenation of modality one and modality two. Each observation is encoded with a corresponding encoder. Each modality is decoded with its own decoder. The loss function has KL-divergence terms to ensure the shared representation regardless of the modality. See Methods for details. **d.** After training, clustering is performed on the resulting embedded vectors. The clustering method is based on iterative clustering as is used in VAMB. Optionally, a reclustering step using single-copy genes can be applied. Here, k-means based reclustering is used when the input is short-read data, and DBSCAN based reclustering is used if the input is long-read data.

114 produce embeddings for downstream tasks that combine the information from two or
 115 more modalities^{43–51}. Therefore, unlike other popular DL based methods with multi-
 116 modal capabilities like stacked autoencoders or siamese networks⁵², most multi-modal
 117 VAEs do not require the dataset to be fully labeled.

118 Here, we introduce TaxVAMB, which combines the strengths of intrinsic and anno-
119 tation features to create high-quality MAGs that cover more taxonomic diversity than
120 any other binning tool (Figure 1). It outperforms all other binners in the number of
121 near-complete assemblies for both short- and long-read datasets. We demonstrate that
122 using TaxVAMB is especially beneficial for datasets with fewer than 100 samples, and
123 that it bins incomplete genomes drastically better than any other tool. TaxVAMB also
124 runs sufficiently fast to be one of the few binning tools that can process large-scale
125 experiments with as many as 1000 samples. We demonstrate the model performance
126 on several short- and long-read datasets from various environments and found that
127 TaxVAMB was 40.2% and 25% better compared to existing methods. TaxVAMB and
128 the source code is freely available at <https://github.com/RasmussenLab/vamb>

129 Results

130 TaxVAMB is a semi-supervised deep learning method that consists of two neural net-
131 works. First, we bridge intrinsic and annotation features by predicting taxonomy labels
132 for unannotated contigs and refining the annotated ones using Taxometer trained on
133 the contigs of this dataset⁴². Second, we use labels on all the taxonomic ranks and
134 the intrinsic features as two modalities to train a bi-modal variational autoencoder
135 to construct the joint latent space. Bi-modal variational autoencoders are a family of
136 VAE-based methods adapted for semi-supervised learning (Supplementary Figure S1).
137 We use an architecture with three encoders (TNFs and abundances; taxonomic labels;
138 TNFs, abundances and taxonomic labels) and two decoders (TNFs and abundances;
139 taxonomic labels). The loss function is designed in a way to ensure that the similar-
140 ity between the latent vectors produced by either of the three encoders is preserved if
141 they are generated from the same contig. Following the strategy previously introduced
142 in the Taxometer method⁴², we use a flat softmax hierarchical loss⁵³ to train on all
143 the taxonomic labels returned by a taxonomic classifier. The latent space is clustered
144 using the original VAMB clustering algorithm. Optionally, the latent space can also be

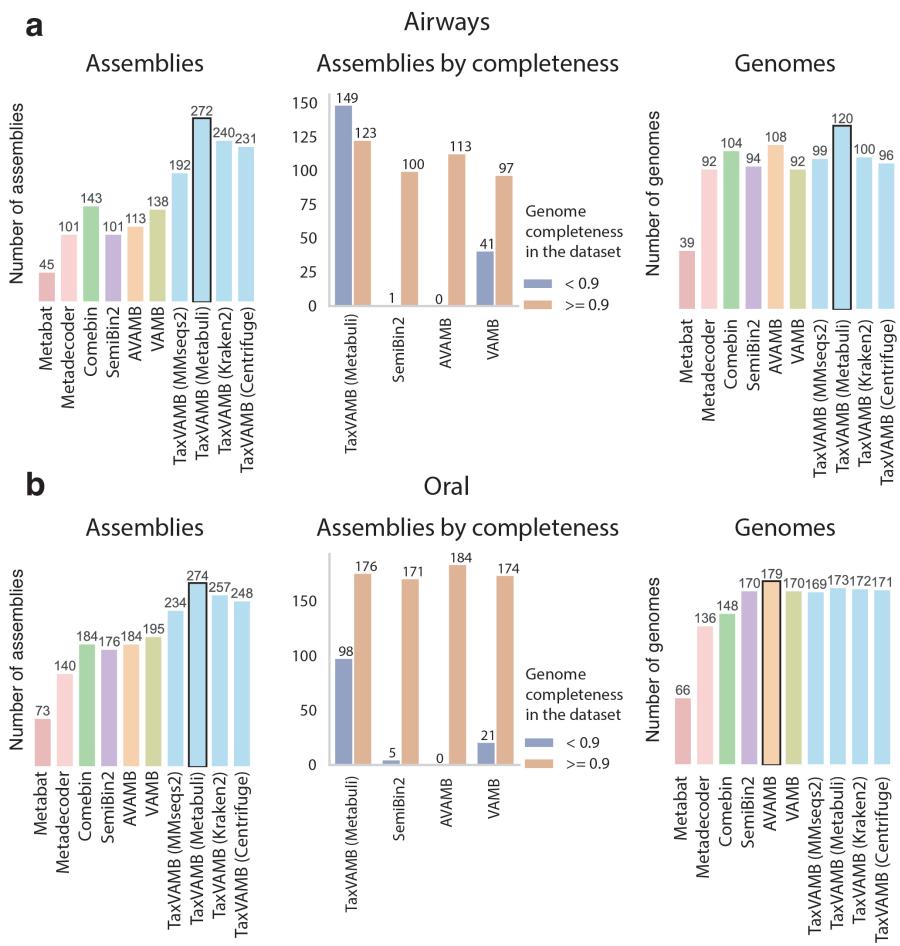


Fig. 2 CAMI2 human microbiome benchmarks. Metagenome binning benchmarks on CAMI2 human microbiome datasets, with TaxVAMB using four different taxonomic classifiers. The bars show the number of near complete assemblies or genomes (recall more than 0.9 and precision more than 0.95). The 'assembly' and 'genome' metrics differ in how the recall is measured. The assembly metrics measures recall in respect to the part of the genome that was provided as input to the biner. The full genome metrics measures recall in respect to the full bacterial genomes even though parts of the genome can be missing from the assembly. Assemblies by completeness show the performance of the binnings stratified by whether the genomes had an assembled share (contigs of that genome provided as input to the biner) less than 0.9, or equal/more than 0.9. SemiBin2, VAMB, AVAMB and TaxVAMB results are shown after applying the k-means based reclustering step. The datasets used are **a**) CAMI2 Airways; **b**) CAMI2 Oral.

145 reclustered using the method adapted from SemiBin2, which is based on single-copy
 146 genes (Figure 1, Supplementary Figure S2, Supplementary Table 1).

¹⁴⁷ **TaxVAMB outperformed all other binners on CAMI2 Toy
148 datasets**

¹⁴⁹ To evaluate TaxVAMB's performance on human microbiome datasets that include
¹⁵⁰ truth annotations, we benchmarked TaxVAMB against six other binners on the
¹⁵¹ synthetic CAMI2 toy human microbiome short-read datasets. We used BinBencher
¹⁵² v0.3.0⁵⁴ to compute two distinct metrics relative to the known ground truth: The
¹⁵³ number of near-complete genomes, and the number of near-complete assemblies (see
¹⁵⁴ Methods). Measured in the number of near-complete assemblies, TaxVAMB outper-
¹⁵⁵ formed all datasets with improvement over the second best binner of 90% for Airways
¹⁵⁶ (272 over 138 from VAMB), 24% for Urogenital (158 over 127 from AVAMB), 9%
¹⁵⁷ for Gastrointestinal (178 over 163 from SemiBin2), 38% for Skin (272 over 184 from
¹⁵⁸ Comebin), and 40% for the Oral dataset (274 over 195 from VAMB) (Figure 2,
¹⁵⁹ Supplementary Figure S3, Supplementary Figure S4). Measured in the number of near-
¹⁶⁰ complete genomes, TaxVAMB demonstrated state-of-the-art performance on 4 out of 5
¹⁶¹ datasets. We found that the improvement of using TaxVAMB compared to the second
¹⁶² best binner were 11% for Airways, 4% for Urogenital, 3% for Gastrointestinal, and 9%
¹⁶³ for Skin, whereas for the Oral dataset, the AVAMB binner was 3% better. The results
¹⁶⁴ showed the largest boost in performance metrics was when the recall was calculated
¹⁶⁵ using the contigs that were provided as an input to the binner, as opposed to the full
¹⁶⁶ genome. This indicates that TaxVAMB was drastically better at binning contigs that
¹⁶⁷ originated from incomplete genomes. For instance, TaxVAMB reconstructed 149 and
¹⁶⁸ 96 assemblies that had less than 0.9 of the total genome present in the input data for
¹⁶⁹ the Airways and Oral datasets, respectively. In comparison, SemiBin2 reconstructed
¹⁷⁰ 1 and 5 assemblies, respectively. For genomes that were almost completely present in
¹⁷¹ the input data the other binners had nearly as good performance. We conclude that
¹⁷² TaxVAMB reaches the state-of-the-art binning performance.

173 MAGs quality depends on the annotation tool

174 Because TaxVAMB used taxonomic annotations as input we evaluated how taxonomic
175 annotations from different tools impacted the performance. Since the quality of the
176 binning depended on noisy and incomplete taxonomic labels, the flexibility of Tax-
177 VAMB is beneficial in case the input dataset is labeled inconsistently by different
178 classifiers. In this case it is possible to run TaxVAMB independently for each taxo-
179 nomic classifier and compare the quality of the resulting bins, even if the labels came
180 from different databases such as GTDB or NCBI annotations. When investigating
181 annotations provided by MMseqs2, Metabuli, Kraken2 and Centrifuge, we found that
182 Metabuli resulted in the highest quality bins for 3 out of 5 the CAMI2 datasets pro-
183 viding 5%-20% more genomes. For the Gastrointestinal and Oral datasets the number
184 of NC genomes differed by only 1 bin for Metabuli and MMseqs2 (Figure S4 - Gas-
185 trointestinal, Oral). Additionally, as Metabuli provided labels at the subspecies level,
186 we included an additional benchmark where these annotations were used as bin iden-
187 tifiers. Here we found that TaxVAMB using Metabuli labels outperformed Metabuli
188 itself with a 16% improvement on average for the five CAMI2 datasets measured in
189 the number of NC genomes and with a 13% improvement on average measured in NC
190 assemblies (Supplementary Figure S5).

**191 TaxVAMB outperformed SCG-based binners without using
192 SCGs**

193 Because single-copy genes are powerful features for binning we wanted to provide
194 an unbiased estimate of performance across binners. To do this, we investigated the
195 number of NC genomes reconstructed by VAMB, SemiBin2 and TaxVAMB before
196 and after SCG-based reclustering. We observed that VAMB, which only uses intrinsic
197 features, outperformed SemiBin2 before reclustering for all five datasets (Figure 3).
198 Improvements of TaxVAMB compared to SemiBin2 were 28% for Airways, 11% for
199 Urogenital, 1.7% for Oral, and 24% for Skin, whereas for the Gastrointestinal dataset,

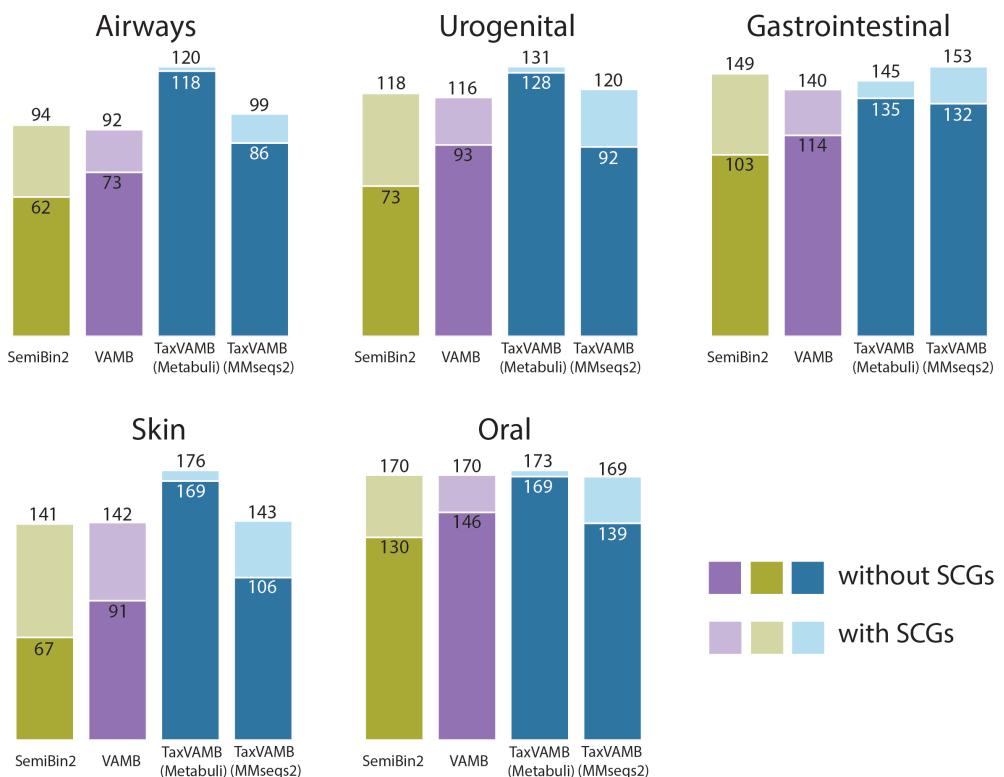


Fig. 3 The effect of reclustering using single-copy genes. Evaluating the effect of reclustering using single-copy genes for SemiBin2, VAMB and TaxVAMB. The darker colors represent binning results without SCGs and the lighter colors represent the results using k-means based SCG-reclustering.

the SemiBin2 binner was 2.6% better. This suggests that a main factor of performance gain was from using single-copy genes for re-clustering of bins. This opposes the claim made in the SemiBin2 study that self-supervised contrastive learning using intrinsic features led to better MAGs. Conversely, we found that the performance of TaxVAMB performance was not that affected by reclustering using SCGs. Here we found that reclustering only resulted in 1.6%-4.1% more genomes when applied to TaxVAMB compared to 30%-110% more genomes when applied to SemiBin2. We conclude that, even when reclustering was applied, the performance of TaxVAMB was less driven by SCGs compared to SemiBin2.

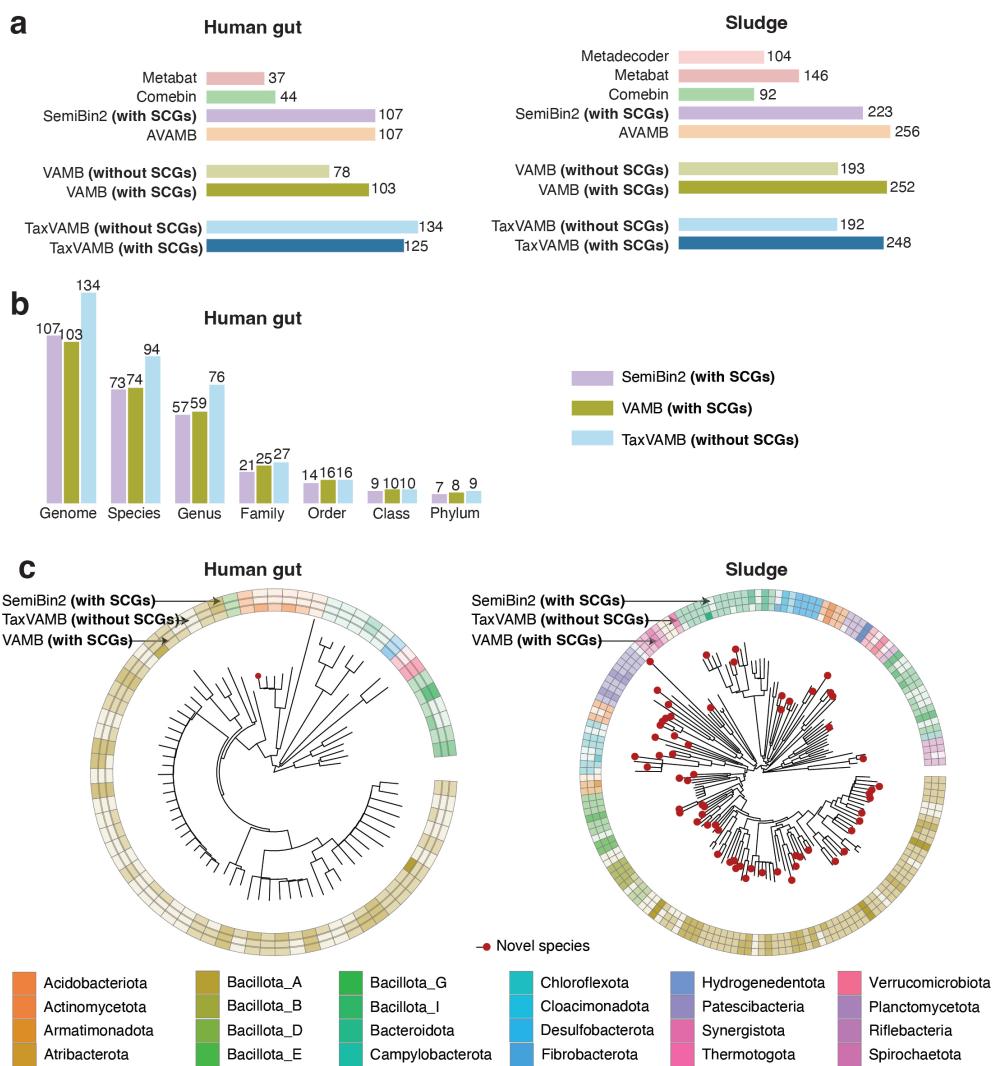


Fig. 4 Long-read datasets benchmarks. **a.** Human gut and sludge dataset benchmarks for different metagenomic binners. The performance is measured as the number of 'near-complete bins', i.e. bins evaluated by CheckM2 to have > 90% completeness and < 5% contamination. VAMB and TaxVAMB are presented with and without SCG-based DBSCAN reclustering. **b.** The phylogenetic diversity of the SemiBin2, VAMB and TaxVAMB near-complete bins, using GTDB-tk placement, as the unique number of taxa on each taxonomic rank. SemiBin2 and VAMB were run using SCG-based DBSCAN reclustering, whereas TaxVAMB was run without. **c.** Visualisation of GTDB-tk placement for SemiBin2, VAMB and TaxVAMB down to the species level, annotated by the color on the phylum level. The darker color in the annotation indicates that more near-complete bins were recovered for this phylum. Red dot indicates a novel species.

**209 TaxVAMB outperformed all other binners on a human gut
210 long read dataset**

211 To test the hypothesis that the performance gains of TaxVAMB would be higher
212 for better-quality taxonomic labels, we benchmarked TaxVAMB on two long read
213 datasets, contrasting a well-studied environment (human gut, three samples) with a
214 poorly studied environment (sludge from an anaerobic digester, four samples). Given
215 that the human gut is more well-studied compared to digested sludge, we hypothe-
216 sized that the taxonomic classifiers would return more complete and correct labels
217 for the human gut samples compared to the sludge samples. As expected, for the
218 human gut dataset TaxVAMB outperformed SemiBin2 and AVAMB, reconstructing
219 25% more near-complete bins (Figure 4a). When applied to the sludge dataset, we
220 found that TaxVAMB performed similarly to VAMB and AVAMB, reconstructing 248,
221 252 and 248 NC genomes, respectively. In comparison, SemiBin2 reconstructed 223 NC
222 genomes from the sludge data. This confirms that a potentially noisy and incomplete
223 taxonomy did not carry sufficient signal to improve or degrade binning performance
224 of TaxVamb. As with the short-read datasets above, we evaluated the effect of reclus-
225 tering using single-copy genes (see Methods subsection 6). For the sludge dataset,
226 the application of SCG based reclustering improved performance for both VAMB and
227 TaxVAMB. For the human gut dataset, TaxVAMB generated the most near-complete
228 bins without using single-copy genes, and the performance was even reduced when
229 using reclustering (Figure 4a). We then investigated the phylogenetic diversity of the
230 NC bins reconstructed by SemiBin2, VAMB, and TaxVAMB from the human gut
231 dataset⁵⁵. (Figure 4b). At phylum level TaxVAMB recovered 9 phyla, VAMB recov-
232 ered 8 and SemiBin2 recovered 7. Additionally, TaxVAMB recovered 27% more species
233 than SemiBin2 at the species level. Most of the MAGs of the sludge dataset were
234 assigned to a novel species whereas novel species were rare in the human gut dataset
235 (Figure 4c). The results suggested that TaxVAMB, using the high-quality taxonomy
236 as input to binning, recovered more diverse MAGs.

**237 Bi-modal variational autoencoder outperformed stacked
238 autoencoder in the number of high quality bins**

239 To ensure that the semi-supervised architecture of the bi-modal VAE was beneficial for
240 binning we empirically evaluated it compared to a stacked autoencoder (Figure S1).
241 Using the short-read and long-read datasets presented above we found that TaxVAMB
242 and the stacked VAE had similar performance with an average 2.3% absolute differ-
243 ence in performance across the short read datasets. The bi-modal VAE outperformed
244 the stacked VAE 6 out of 15 times for MMseqs2, Kraken2 and Centrifuge classifiers.
245 When Metabuli taxonomic labels were used, the bi-modal VAE always outperformed
246 the stacked VAE with an average gain of 16% across all CAMI2 datasets (Supplemen-
247 tary Figure S6a). For the long-read datasets, TaxVAMB outperformed the stacked
248 autoencoder architecture across all datasets and classifiers, with an average gain of
249 13% more high quality genomes for the Human Gut dataset and 18% more high qual-
250 ity genomes for the Sludge dataset (Supplementary Figure S6b). This indicate that
251 binning the latent vectors provided by TaxVAMB using a bi-modal VAE resulted in
252 better overall performance for the task of metagenomics binning compare to the same
253 workflow that uses a stacked VAE.

**254 Taxonomy primarily improved binning at low number of
255 samples**

256 The benefit of using co-abundance for binning positively correlates with the number
257 of samples in the experiment⁵⁶. We, therefore, investigated how taxonomic labels
258 improved binning as a function of the size of the dataset. To achieve this, we varied
259 the number of samples from 1-1000 when using human gut microbiome samples from
260 Almeida et al.⁵⁷ comparing VAMB to TaxVAMB ((Figure 5a, see Methods for details)).
261 We found that when the dataset included 1,000 samples TaxVAMB only recovered 3%
262 more near-complete MAGs than VAMB. However, when we subsampled the datasets to
263 100 samples, the performance improvement increased to 16%, while for datasets of 10

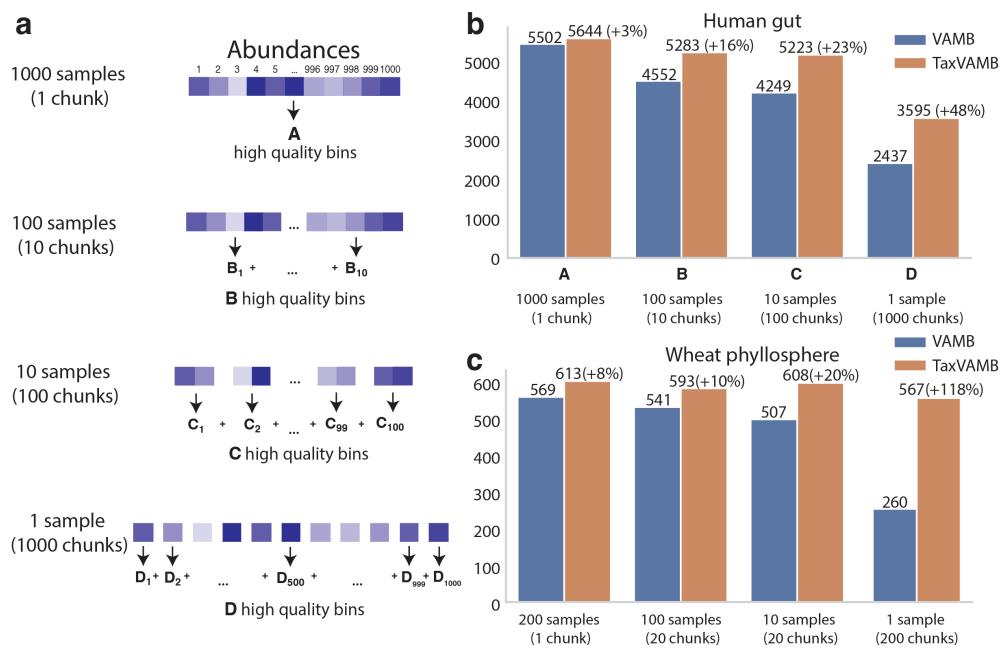


Fig. 5 Effect of abundance vector and taxonomy information. **a.** For 1,000 samples, a single TaxVAMB and VAMB run was performed using all contigs and the entire abundance vector. 10 runs were performed on chunks of 100 samples and their corresponding contigs and abundances, 100 runs with 10 samples and 1,000 runs with 1 sample. The number of near-complete bins for all chunks was summed for each set of 1,000 samples. **b.** The results for the human gut dataset of Almeida et al. using TaxVAMB and VAMB. **c.** The results for the wheat phyllosphere dataset using TaxVAMB and VAMB.

samples, the improvement was 23%. Finally, when performing single-sample binning, i.e. one sample as input, the performance increased by 48% when using TaxVAMB compared to VAMB (Figure 5b). In a similar experiment using a wheat phyllosphere dataset we found that the gains compared to single-sample binning were even larger, resulting in 118% more bins in TaxVAMB compared to VAMB. Interestingly, for the experiment using datasets of 10 samples, we found that TaxVAMB increased the number of near-complete bins to the level of running VAMB on chunks of 100 samples. Therefore, in experiments with few samples, TaxVAMB was able to compensate for a less expressive abundance vector by using the taxonomy label modality. We conclude that TaxVAMB delivered larger gains compared with VAMB on datasets with fewer than 100 samples.

275 TaxVAMB provided consistent bin annotations

276 A key step of TaxVAMB is to predict taxonomic assignments of contigs without taxo-
277 nomic labels. This is done using a neural network (Taxometer) resulting in taxonomic
278 labels for all binned contigs (Supplementary Figure S7a). We, therefore, investigated if
279 majority voting of these could be used as a taxonomic classification of the bins. Using
280 Kraken2 classification of the CAMI2 human microbiome dataset we isolated the sub-
281 set of high-quality MAGs. We found that 95-98% of bins per dataset were correctly
282 annotated down to the species level, with the GTDBtk classifier⁵⁵ correctly annotat-
283 ing 98-99% bins (Supplementary Figure S7b, Supplementary Figure S8). TaxVAMB
284 annotations have the advantage over GTDBtk of not requiring any additional run-
285 time and not being limited by prokaryotes in providing the annotations. Therefore, we
286 argue that TaxVAMB provide high-quality taxonomic annotations for the bins created
287 from data from well-studied environments without the need to use additional MAG
288 taxonomic classification tools.

**289 TaxVAMB uncovers both bacterial and fungal MAGs in wheat
290 phyllosphere dataset**

291 Finally, we investigated the results of applying TaxVAMB to the short-read wheat
292 phyllosphere dataset. The dataset consisted of 211 samples from the surface of wheat
293 flag leaves at nine time points during the end of growth season of 2022 from a single
294 field in Denmark (see Methods). From the dataset we reconstructed 614 high quality
295 and 647 medium quality bacterial MAGs across five phyla (Actinomycetota, Bacillota,
296 Bacteroidota, Deinococcota, Pseudomonadota) (Figure 6a, Supplementary Figure S9).
297 We found that the binned MQ and HQ MAGs explained 13.4%-98.4% of the total
298 reads across the samples with a mean of 49.2% of the reads (Figure 6c). When inves-
299 tigating the prevalence of the species, we found that the five most prevalent species
300 measured in the number of MAGs assigned were present in 30-60% of all the samples
301 (*Pseudomonas poae*, *Frigoribacterium sp001421165*, *Pseudomonas graminis*, *Pantoea*

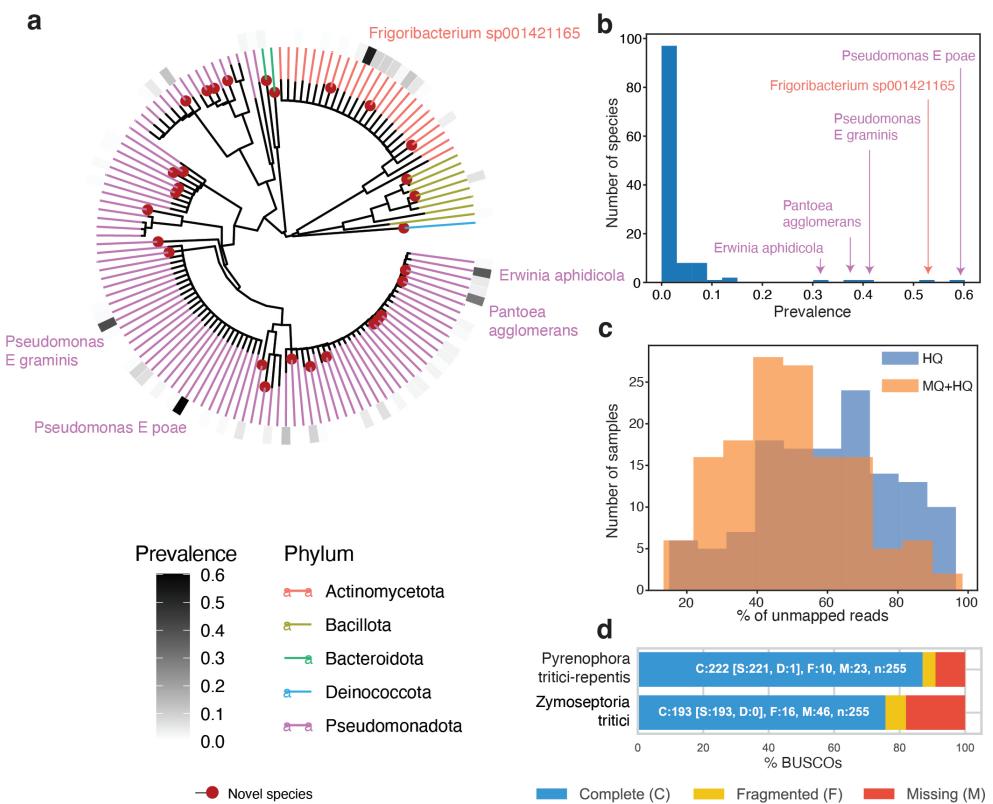


Fig. 6 Wheat phyllosphere MAGs. **a.** Phylogenetic tree of high quality bacterial MAGs indicating the most prevalent species in terms of high quality MAGs per sample. **b.** Distribution of prevalences for all species. The top 5 most prevalent species are annotated with labels. **c.** Distributions of shares of unmapped reads for each sample with MQ or HQ MAGs. Blue color (HQ) is the share of unmapped reads when only mapping to HQ MAGs (completeness > 90%, contamination < 5%). Orange color (HQ+MQ) is the share of unmapped reads when mapping to HQ and MQ MAGs (completeness ≥ 50%, contamination < 10%). **d.** BUSCO results for two fungal MAGs, annotated by TaxVAMB as *Zymoseptoria tritici* and *Pyrenophora tritici-repentis*.

302 *agglomerans*, *Erwinia aphidicola*) and had previously been described in literature as
 303 part of the wheat phyllosphere^{58–63} (Figure 6b, Supplementary Table 2, Supplemen-
 304 tary Table 3). In addition to these, TaxVAMB reconstructed a novel species of genus
 305 *Sphingomonas* that was present in 12% of samples with an abundance of >1% of
 306 the mapped reads. Additionally, we discovered that the *Pantoea agglomerans* species
 307 was more prevalent as the plants senesced (Mann-Whitney U test, p-value = 2e-18)
 308 (Supplementary Figure S9). We also tested the ability of TaxVamb to recover fun-
 309 gal bins by investigating the bins that were annotated by TaxVAMB as Eukaryotes.

310 Two of such bins were larger than 27Mb, with 99.9% and 20% of their contigs anno-
311 tated with fungal species after the Taxometer step. Thus, we recovered *Zymoseptoria*
312 *tritici* and *Pyrenophora tritici-repentis* MAGs with corresponding BUSCO⁶⁴ com-
313 pleteness scores of 75% and 87% (Figure 6d). Both of these fungal species are known
314 wheat pathogens^{65,66}. We conclude that TaxVAMB recovered a large variety of novel
315 MAGs of medium and high quality, providing insights into the bacterial and fungal
316 composition of the wheat phyllosphere.

317 Discussion

318 In summary, we present TaxVAMB, a method for combining intrinsic features (TNFs
319 and abundances) with annotation features (taxonomic labels). We utilize the full hier-
320 archical structure of taxonomic labels with a deep hierarchical loss, allowing us to
321 train the model even on contig annotations from higher taxonomic ranks. We justified
322 the network architecture choice by empirically evaluating a competing stacked VAE
323 model. The labels were also used to assign the preliminary taxonomic annotations
324 to MAGs down to the species level, matching GTDBtk tool in performance for the
325 high quality bins when tested on CAMI2. In the number of high quality bins, Tax-
326 VAMB exceeds or matches the performance of other binners for the genomes that are
327 almost completely present in the input, and shows unmatched performance in binning
328 incomplete genomes.

329 We identified two conditions where TaxVAMB exhibits the biggest gains com-
330 pared to previous state-of-the-art: a) sufficiently high quality taxonomic labels which
331 occurred for the datasets from more well studied environments such as human gut;
332 b) a relatively low number of samples where the signal from the co-abundance vector
333 was less strong (< 100 samples judging from the Figure 5). For example, we expect
334 that one use case of TaxVAMB is application to studies of the human microbiome
335 where few samples are available. Similarly, TaxVAMB also shows unmatched perfor-
336 mance for binning incomplete genomes, which, just like datasets with a small number

337 of samples, produce low quality abundance vectors. We notice that using TaxVAMB in
338 poorly studied and much more complex environments such as anaerobic digester sludge
339 does not degrade the performance compared to the unsupervised binning provided by
340 VAMB. We also demonstrated that TaxVAMB does not rely on single-copy genes for
341 reaching optimal performance, which allows it to better bin incomplete genomes and
342 non-bacterial entities.

343 Predicting annotations for the full dataset based on pre-annotated contigs might
344 increase the bias toward well-studied taxa, since these taxa are more likely to be
345 returned by a taxonomic classifier such as MMseqs2. The qualities of predictions
346 depend on the share of annotated contigs and taxonomic diversity of the samples.
347 We address the possible bias in two ways. First, the taxonomy predictions come as
348 probability vectors for each taxonomic rank, allowing to set a confidence threshold for
349 predictions. A contig with an ambiguous prediction will remain unannotated on this
350 taxonomic rank and below. Second, unsupervised learning still allows correct binning
351 of contigs without any good match in the database, but which share intrinsic features
352 such as TNFs and abundances with each other.

353 The databases of reference genomes are constantly updated (GTDB increased by
354 around 30% in terms of bacterial species clusters from v207 to v214⁶⁷, NCBI estimates
355 annual growth in terms of the number of genomes at 15%⁶⁸). The bias introduced by
356 taxonomic annotation of a subset of contigs in a dataset will continue to reduce as the
357 number of diverse genomes in databases grows.

358 After reclustering using single-copy genes, unsupervised (VAE) and self-supervised
359 (contrastive learning) models demonstrate similar performance on short- and long-
360 read datasets. While we flag optimizing for single-copy genes metrics during training
361 as a potential source of error when also used as an evaluation metric, we also interpret
362 the large gains from single-copy gene-guided reclustering as the indicator that future
363 research in the field of metagenomic binning would find the biggest gains in integrating
364 new data modalities with intrinsic features, and not solely in applying new algorithms
365 to TNFs and abundances.

366 Multi-omics data integration is a powerful technique for understanding complex
367 biological systems^{69–71}. Semi-supervised multimodal VAEs can be easily adapted to
368 learn from weakly labeled heterogeneous multi-omics datasets in other fields beyond
369 metagenomics binning. The same applies to the hierarchical loss, since many biological
370 data types follow a hierarchical structure.

371 Improving over previous attempts to incorporate taxonomic labels to the
372 metagenome binning process, we conclude that it greatly improves binning when the
373 full hierarchical structure is in use, and is a great additional input when the abundance
374 vector is not informative enough. As the quality of reference databases will improve
375 over time, the impact of using TaxVAMB will be potentially even more powerful in
376 the future.

377 Methods

378 Bi-modal variational autoencoder

379 VAE is a generative model performing variational inference over the latent variable
380 z . The model is formally defined as $p(x, z) = p(z)p(x|z)$. The intractable posterior
381 $q(z|x)$ and the conditional distribution $p(x|z)$ are approximated by neural networks
382 using the ELBO-loss function:

$$\mathcal{L} = E_{q(z|x)}[\log p(x|z)] - D_{KL}(q(z|x)\|\mathcal{N}(0, I)) \quad (1)$$

383 The bi-modal VAE extends the basic VAE by allowing training and inference on
384 the dataset where: a) the input consists of two modalities and b) a modality can be
385 missing for one or more samples. We define the modality here as the part of the data
386 that is missing for part of the dataset. Thus notice that while the VAMB model is
387 trained on both TNFs and abundances, we do not define it as bi-modal for the purpose
388 of this summary, since both TNFs and abundances are present for all samples and

389 can be trivially converted into one modality by concatenating the corresponding input
 390 vectors.

391 While the VAE approximates the posterior $q(z|x)$ with a neural network encoder
 392 that takes x as an input, bi-modal VAE extends this approach by modelling $q(z|x_1, x_2)$,
 393 $q_1(z|x_1)$ and $q_2(z|x_2)$, which replace the single $q(z|x)$. There are two decoders approx-
 394 imating distributions $p(x_1|z)$ and $p(x_2|z)$. Multimodal VAEs differ in 1) the way
 395 they approximate $q(z|x_1, x_2)$, $q_1(z|x_1)$ and $q_2(z|x_2)$ by neural networks and/or 2) the
 396 structure of the loss function.

397 TaxVAMB implements the VAEVAE⁵⁰ model from bi-modal VAE family which
 398 models $q(z|x_1, x_2)$, $q_1(z|x_1)$ and $q_2(z|x_2)$ by corresponding neural networks. The
 399 following ELBO-like loss \mathcal{L} is minimised:

$$\begin{aligned} \mathcal{L}_{paired} = & E_{p_{paired}(x_1, x_2)} [E_{q(z|x_1, x_2)} [\log p_1(x_1|z) + \log p_2(x_2|z)] \\ & - D_{KL}(q(z|x_1, x_2) \| p(z|x_1)) - D_{KL}(q(z|x_1, x_2) \| p(z|x_2))] \\ & + E_{p_{paired}(x_1)} [E_{q(z|x_1)} [\log p_1(x_1|z)] - D_{KL}(q(z|x_1) \| p(z))] \\ & + E_{p_{paired}(x_2)} [E_{q(z|x_2)} [\log p_2(x_2|z)] - D_{KL}(q(z|x_2) \| p(z))] \quad (2) \\ \mathcal{L}_1 = & E_{p_{unpaired}(x_1)} [E_{q(z|x_1)} [\log p_1(x_1|z)] - D_{KL}(q(z|x_1) \| p(z))] \quad (3) \\ \mathcal{L}_2 = & E_{p_{unpaired}(x_2)} [E_{q(z|x_2)} [\log p_2(x_2|z)] - D_{KL}(q(z|x_2) \| p(z))] \quad (4) \\ \mathcal{L} = & \mathcal{L}_{paired} + \mathcal{L}_1 + \mathcal{L}_2 \quad (5) \end{aligned}$$

400 with $D_{KL}(p(x) \| q(x))$ being the Kullback–Leibler divergence between two proba-
 401 bility distributions $p(x)$ and $q(x)$, defined as:

$$D_{KL}(p(x) \| q(x)) = \int_{-\infty}^{\infty} p(x) \log \left(\frac{p(x)}{q(x)} \right) dx \quad (6)$$

402 The training procedure includes constructing the dataset with paired and unpaired
403 samples. Let C be a list of all contigs. Three copies of C , denoted as C_{paired} , C_1 , and
404 C_2 , are independently shuffled. The paired samples are ordered tuples (x_1, x_2) where
405 x_1 is a concatenation of TNF vector and abundance vector (the input of VAMB) and
406 x_2 is a taxonomy label vector described in Subsection 6, and x_1 and x_2 correspond to
407 the same contig from the set C_{paired} . An unpaired TNFs and abundances vector x'_1
408 corresponds to a contig from the list C_1 . An unpaired taxonomy label corresponds to
409 a contig from the list C_2 . The forward pass follows the steps from the Algorithm 1.

Algorithm 1 Loss computation (forward pass)

Require: Paired sample (x_1, x_2) , unpaired sample x'_1 , unpaired sample x'_2

- 1: $z' \sim q(z|x_1, x_2)$
 - 2: $z_{x_1} \sim q_1(z|x_1)$
 - 3: $z_{x_2} \sim q_2(z|x_2)$
 - 4: $d_1 = D_{KL}(q(z'|x_1, x_2) \| q_1(z_{x_1}|x_1)) + D_{KL}(q_1(z_{x_1}|x_1) \| p(z))$
 - 5: $d_2 = D_{KL}(q(z'|x_1, x_2) \| q_2(z_{x_2}|x_2)) + D_{KL}(q_2(z_{x_2}|x_2) \| p(z))$
 - 6: $\mathcal{L}_{paired} = \log p_1(x_1|z) + \log p_2(x_2|z) + \log p_1(x_1|z_{x_1}) + \log p_2(x_2|z_{x_2}) + d_1 + d_2$
 - 7: $\mathcal{L}_{x_1} = \log p_1(x'_1|z_{x_1}) + D_{KL}(q_1(z_{x_1}|x'_1) \| p(z))$
 - 8: $\mathcal{L}_{x_2} = \log p_2(x'_2|z_{x_2}) + D_{KL}(q_2(z_{x_2}|x'_2) \| p(z))$
 - 9: $\mathcal{L} = \mathcal{L}_{paired} + \mathcal{L}_{x_1} + \mathcal{L}_{x_2}$
-

410 **Data preprocessing**

411 The workflow of preprocessing the data is the same as in Taxometer (v.b5fd0ea)⁴²
412 and VAMB¹¹. The synthetic short paired-end reads from each sample were aligned
413 using bwa-mem (v.0.7.15)⁷². BAM files were sorted using samtools (v.1.14)⁷³. Contigs
414 $\leq 2,000$ base pairs (bp) were removed for each dataset. the long-read datasets were
415 both sequenced using Pacific Biosciences HiFi technology. We assembled each sample
416 using metaMDBG (v.b55df39)⁷⁴, mapped reads using minimap2 (v.2.24)⁷⁵ with the
417 '-ax map-hifi' setting, and then continued with the same workflow as with the short
418 reads.

419 Abundances and TNFs

420 The workflow of computing abundances and TNFs is the same as in Taxometer
421 (v.b5fd0ea)⁴² and VAMB¹¹. Computation of abundances and TNFs was done using
422 the VAMB metagenome binning tool¹¹. To determine TNFs, tetramer frequencies of
423 non-ambiguous bases were calculated for each contig, projected into a 103-dimensional
424 orthonormal space and normalized by z-scaling each tetranucleotide across the contigs.
425 To determine the abundances of each sample, we used pycoverm 0.6.0⁷⁶. The abun-
426 dances were first normalized within sample by total number of mapped reads, then
427 across samples to sum to 1. To determine absolute abundance, the sum of abundances
428 for a contig was taken before the normalization across samples. The dimensionality of
429 the feature table was then $N_c \times (103 + N_s + 1)$ where N_c was the number of contigs,
430 N_s was the number of samples.

431 Network architecture and hyperparameters

432 The encoder architectures for the concatenated vector of abundances and TNFs is the
433 same as in Taxometer (v.b5fd0ea)⁴² and VAMB¹¹. The input vector of dimensionality
434 $N_c \times (103 + N_s + 1)$ was passed through 4 fully connected layers ($(103 + N_s + 1) \times 512$,
435 512×512 , 512×512 , 512×512) with leaky ReLU activation function (negative slope
436 0.01), each using batch normalization (epsilon $1e - 05$, momentum 0.1) and dropout
437 ($P = 0.2$).

438 The encoder network for the taxonomy labels had the input dimensions of N_l
439 where N_l was the number of leaves in the taxonomic tree. The input vector was passed
440 through 4 fully connected layers ($N_l \times 512$, 512×512 , 512×512 , 512×512) with
441 leaky ReLU activation function (negative slope 0.01), each using batch normalization
442 (epsilon $1e - 05$, momentum 0.1) and dropout ($P = 0.2$).

443 The encoder network for the concatenation of the two modalities had the input
444 dimensions of $(103 + N_s + 1) + N_l$ where N_s was the number of samples and N_l was the

445 number of leaves in the taxonomic tree. The input vector was passed through through
446 4 fully connected layers $((103 + N_s + 1) \times 512, 512 \times 512, 512 \times 512, 512 \times 512)$ with
447 leaky ReLU activation function (negative slope 0.01), each using batch normalization
448 (epsilon $1e - 05$, momentum 0.1) and dropout ($P = 0.2$).

449 The bi-modal VAE has two decoder networks, one for each modality. Both of them
450 follow the same architectures as the corresponding encoders, with the input vector
451 having the dimensionality of the latent space, and the output having the dimensionality
452 of the corresponding modality.

453 For short-read datasets, the network was trained for 300 epochs with batch size
454 256, latent space dimensionality 32. For long-read datasets, the network was trained
455 for 1000 epochs with batch size 1024, latent space dimensionality 64. All models were
456 using the Adam optimizer with learning rates set via D-Adaptation⁷⁷. The model
457 was implemented using PyTorch (v.1.13.1)⁷⁸, and CUDA (v.11.7.99) was used when
458 running on a V100 GPU.

459 Hierarchical loss

460 The hierarchical loss is the same as in Taxometer (v.b5fd0ea)⁴². A phylogenetic tree
461 was constructed for each dataset from the taxonomy classifier annotations for the set
462 of contigs. Thus, the resulting taxonomy tree T was a subgraph of a full taxonomy
463 and the space of possible predictions was restricted to the taxonomic identities that
464 appeared in the search results. For the above experiments, we used a flat softmax loss.
465 Let N_l be the number of leaves in the tree T . The likelihoods of leaf nodes of the
466 taxonomy tree were obtained from the softmax over the network output layer with
467 dimensionality $1 \times N_l$. The likelihood of an internal node was then a sum of likelihoods
468 of its children and computed recursively bottom-up. The model output was a vector of
469 likelihoods for each possible label. For the backpropagation, the negative log-likelihood
470 loss was computed for all the ancestors of the true node and the true node itself.
471 Predictions were made for all taxonomic levels and for each level, the node descendant

472 with the highest likelihood was selected. If no node descendant had likelihood > 0.5,
473 the predictions from this level and the levels below were not included in the output.

474 Taxonomic classifiers

475 We obtained the taxonomic annotations for contigs of all seven short-read and
476 two long-read datasets from MMseqs2 (v.7e2840)³³, Metabuli (v.1.0.1)⁷⁹, Centrifuge
477 (v.1.0.4)⁸⁰ and Kraken2 (v.2.1.3)⁸¹. For MMseqs2, we used the mmseqs taxonomy
478 command. For Metabuli, we used the metabuli classify command with *-seq-mode 1*
479 flag. For Centrifuge, we used the centrifuge command with *-k 1* flag. For Kraken2, we
480 used the kraken command with *-minimum-hit-groups 3* flag. MMseqs2 and Metabuli
481 were configured to use GTDB v207 as the reference database. Centrifuge, Kraken2
482 and MetaMaps were configured to use NCBI identifiers. All the taxonomic annotations
483 were first refined with Taxometer⁴² (v.b5fd0ea) with the default parameters (epochs
484 100, batch size 1024).

485 Benchmarked binners

486 Metabat (v.2.12.1) *metabat* command with the default parameters is used. Metadecoder
487 (v.1.0.19) *coverage*, *seed* and *cluster* commands are used as described in <https://github.com/liu-congcong/MetaDecoder>. Comebin (v.1.0.3) *run_comebin.sh* command
488 with the default parameters is used. SemiBin2 (v.2.1.0) *multi-easy-bin* command is
489 used with the flags *-engine gpu*, *-separator C*, *-t 20*, *-write-pre-reclustering-bins* and
490 *-self-supervised*. VAMB, AVAMB and TaxVAMB were run as a part of the VAMB
491 codebase (commit 5f2cd7), with the corresponding commands "*vamb bin default*",
492 "*vamb bin avamb*" and "*vamb bin taxvamb*".

494 Reclustering using SCGs

495 Short-read and long-read reclustering algorithms that use single-copy genes are the
496 same as introduced in SemiBin2⁸². The code was adapted from the SemiBin2 code-
497 base (https://github.com/BigDataBiology/SemiBin/blob/main/SemiBin/long_read_
498 `cluster.py` and <https://github.com/BigDataBiology/SemiBin/blob/main/SemiBin/>
499 `cluster.py`) for the TaxVAMB codebase (<https://github.com/RasmussenLab/vamb/>
500 `blob/master/vamb/reclustering.py`). TaxVAMB uses the same 107 single-copy marker
501 genes as used in the SemiBin2 tool to estimate the completeness, contamination, and
502 F1-score of every bin. Completeness for each bin is calculated as $\frac{N}{107}$, contamination
503 as $\frac{G-N}{G}$ and F1-score as $\frac{2*\text{completeness}*(1-\text{contamination})}{\text{completeness}+(1-\text{contamination})}$. N is the number of different
504 single-copy genes in a bin, G is the total number of sequences matching any single-copy
505 gene.

506 For the short-read datasets, k-means based reclustering of TaxVAMB/VAMB
507 clusters is performed. Bins where more than one marker gene of the same kind is
508 present are reclustered with the weighted k-means method using the contigs contain-
509 ing the repeated marker gene as the initial centroids. It results in bins with reduced
510 contamination.

511 For the long-read datasets, the DBSCAN algorithm from a Python library *scipy*
512 (v1.10.0) was used to perform the clustering from scratch (the previous clusters, made
513 by TaxVAMB/VAMB, were not used). Same as in SemiBin2, DBSCAN was run with
514 ϵ value equals to 0.01, 0.05, 0.1, 0.15, 0.2, 0.25, 0.3, 0.35, 0.4, 0.45, 0.5, and 0.55.
515 From all the resulting bins, the best one by the F1-score was recursively selected, and
516 its contigs were removed from all the remaining bins, after which the selection of the
517 best bin repeats. It repeats until no more bins fulfill the criteria for minimal quality
518 (completeness more than 90%, contamination less than 5%). One change that was
519 made in the TaxVAMB long-read reclustering was that it performed the described
520 procedure per a set of contigs assigned to the same genus by the Taxometer refinements
521 of the provided taxonomic annotations.

522 CAMI2 benchmarks

523 For short-read benchmarking, we used five CAMI2 datasets: Airways (10 samples),
524 Oral (10 samples), Skin (10 samples), Gastrointestinal (10 samples), Urogenital (9
525 samples), assemblies for which were sample-specific. We benchmarked the follow-
526 ing binners on the synthetic CAMI2 toy human microbiome dataset: Metabat¹⁰,
527 MetaDecoder⁸³, COMEBin²², SemiBin2⁸², AVAMB²⁸, and VAMB¹¹. We used taxo-
528 nomic labels from four taxonomic classifiers as an input to TaxVAMB: MMSeqs2³³,
529 Metabuli⁷⁹, Kraken2⁸¹ and Centrifuge⁸⁰. AVAMB, VAMB and TaxVAMB bins were
530 postprocessed with reclustering using single-copy genes. We used the number of high
531 quality bins and assemblies estimated via BinBencher (v0.3.0)⁵⁴ as a metric. For the
532 MAG taxonomic annotation experiment we used CheckM2 (v1.0.2)⁴⁰.

533 We benchmarked using BinBencher (v.0.3.0)⁵⁴ against a reference computed from
534 the CAMI2 ground truth. The metrics used were number of near-complete (defined as
535 recall ≥ 0.9 , precision ≥ 0.95) assemblies or genomes. As defined in the BinBencher
536 paper, for genomes, the recall was counted relative to the full length of the genome from
537 which the reads were simulated from, whereas when counting assemblies, the recall was
538 relative to the assembled part of those genomes, i.e. the part of the genomes covered
539 by a contig which was used as input to the binner. The number of near-complete
540 genomes reflect the MAG quality relative to the underlying biological organism and
541 thus depends more on limitations of the dataset, whereas the assembly metric may
542 better reflect the methodological gains from using a different algorithm.

543 Long-read benchmarks

544 For long-read benchmarking we used a human gut microbiome dataset with 4 samples
545 and a dataset from anaerobic digester sludge with 3 samples⁸⁴, both sequenced using
546 Pacific Biosciences HiFi technology. We assembled each sample using metaMDBG
547 (v.b55df39)⁷⁴, mapped reads using minimap2 (v.2.24)⁷⁵ with the '-ax map-hifi'

548 setting, and from there proceeded as with the short reads. For evaluating the qual-
549 ity (completeness and contamination) of the resulting MAGs we used CheckM2
550 (v.1.0.2)⁴⁰. The Metadecoder binner failed on the human gut dataset experiment with
551 an internal code error and its results are thus not displayed on the figure.

552 Multisample scaling

553 For the experiment that evaluates the number of bins given a different number of
554 samples, we used a short-read human gut dataset with 1,000 samples from Almeida
555 et al.⁵⁷, as well as our own wheat phyllosphere dataset with 211 samples. For each
556 dataset, we split all the samples into three sets of chunks: 1) chunks of 100 samples; 2)
557 chunks of 10 samples; 3) chunks of 1 sample. Each chunk was treated as an independent
558 dataset. We then summed the resulting number of near-complete bins within each set
559 of chunks.

560 Wheat phyllosphere dataset: sample collection and processing

561 Twenty-four field plots of *Triticum aestivum* were sampled by collecting composite
562 samples of 30 flag leaves nine times between June 7th and July 14th 2022, at a field trial
563 in Ringsted, Denmark. The experimental design included three wheat cultivars, four
564 replicates, and two treatments, which were unsprayed and sprayed with a fungicide.
565 The samples were washed in 100 ml wash solution (0.9% NaCl + 0.05% Tween80),
566 vigorously shaken for 2 minutes, sonicated for 2 minutes and then vigorously shaken
567 again for 2 minutes, filtered (10 µm), centrifuged (4000 x g, 15 min) and the pellet
568 resuspended in 1 ml 1x PBS and stored at -20°C until DNA extraction using the
569 FastDNA™ SPIN Kit (MP Biomedicals, CA, USA) for Soil according to instructions
570 eluding in 100 µl DES. DNA libraries were build using the Illumina Nextera XT kit
571 (Illumina, CA, USA), but samples with <0,1 ng/µl DNA were built with a 1-fold
572 diluted ATM, 20 PCR cycles and a higher ratio of AMPure XP beads (Beckman

573 Coulter, IN, USA)⁸⁵. Libraries were sequenced using Illumina paired-end (2 x 150bp)
574 technology (NovaSeq 6000 S4 v1.5).

575 Wheat phyllosphere dataset: data analysis

576 Raw sequencing reads were filtered using fastp (v.023.2)⁸⁶ with the following option:
577 '--trim_tail 1 --cut_tail --trim_poly_g --dedup --length_required 80'.
578 Quality control of the filtered reads was assessed using MultiQC (v.1.12)⁸⁷. To
579 remove reads originating from wheat or potential human contamination, the reads
580 were mapped to the reference sequences GCF_018294505.1, MG958554.1 and
581 GCF_000001405.40 (GRCh38.p14). Mapping was performed using Bowtie⁸⁸ (v.2.5.3).
582 Paired reads where both mates were unmapped were extracted using Samtools
583 (v.1.18)⁷³ with the 'fastq -f 13' option.

584 Metagenomic assemblies were generated for each sample using SPAdes (v.3.15.4)⁸⁹
585 with the '--meta -k 21,33,55,77,99' option. Assembly statistics were computed
586 using QUAST (v.5.2.0)⁹⁰.

587 MAGs are assigned the taxonomy using GTDBtk (v.2.4.0) configured with the
588 GTDB database v220. Taxonomic trees are built using *ggtree*⁹¹ (v.3.19), *tidytree*
589 (v.0.4.6), *treeio* R (v.4.4.1) libraries. Mann-Whitney U test is performed on *Pantoea*
590 *agglomerans* abundances by splitting the samples into two groups: 143 samples from
591 the earlier days (2022-06-07, 2022-06-10, 2022-06-14, 2022-06-17, 2022-06-21) and 103
592 samples from the later days (2022-06-28, 2022-07-04, 2022-07-07, 2022-07-14) using
593 *scipy* (v1.10.0).

594 Acknowledgements

595 S.K., M.N., S.R., N.S.O., L.R., L.M.F.J., P.E.D., A.G., K.N.N., S.C. and L.H. were
596 supported by the Novo Nordisk Foundation (grant NNF19SA0059348). P.P., J.N.N.,
597 S.K., K.N.N. and S.R. were supported by the Novo Nordisk Foundation (grant

598 NNF23SA0084103). S.K., P.P., J.N.N., K.N.N. and S.R. were supported by the Novo
599 Nordisk Foundation (grant NNF14CC0001). P.P., J.N.N. and S.R. were supported by
600 the Novo Nordisk Foundation (grant NNF20OC0062223). S.R. was supported by the
601 Novo Nordisk Foundation (grant NNF21SA0072102). We thank Chayan Roy, Sif Chris-
602 tine Lykke Hougaard and Xianfu Liu for contributing to the wheat phyllosphere data
603 collection.

604 Author Contributions

605 S.K., M.N., J.N.N. and S.R. designed the experiments. P.P., J.N.N. and K.N.N. pre-
606 processed the datasets. S.K. wrote the software and performed the analysis. M.N.,
607 P.P., J.N.N. and S.R. provided guidance and input for the analysis. S.C. and J.C.W.
608 selected the trial fields and developed sampling protocols for the wheat phyllosphere
609 dataset. N.S.O., L.R., L.M.F.J., P.E.D., A.G., K.N.N., S.C. and L.H. performed sam-
610 ple collection, sample processing, DNA extractions and library building for the wheat
611 phyllosphere dataset. S.K. wrote the manuscript with contributions from all coauthors.
612 All authors read and approved the final version of the manuscript.

613 Data availability

614 The CAMI2 datasets were downloaded from <https://data.camichallenge.org/participate> from "2nd CAMI Toy Human Microbiome Project Dataset"
615 (5 human microbiome datasets), "2nd CAMI Challenge Marine Dataset" (Marine),
616 "2nd CAMI Challenge Rhizosphere challenge" (Rhizosphere). The long-read human
617 gut dataset is available at <https://downloads.pacbcloud.com/public/dataset/Sequel-IIe-202104/metagenomics/>. The long-read sludge dataset is available at the ENA as
618 part of the study PRJEB39861. The 1000 samples short-read human gut dataset was
619 first published by Almeida et al. The de novo assemblies of the Almeida dataset were
620 obtained through personal communication with A. Almeida and R. D. Finn, and the
621

623 reads downloaded from ENA ERP108418 as specified in their publication. The phyl-
624 losphere short-read dataset is available at the ENA using the accession ERP165292.
625 The near complete and medium quality MAGs from the phyllosphere are available for
626 download at Zenodo⁹² by the link <https://zenodo.org/records/13959411>. The MAGs
627 will be uploaded to ENA upon acceptance as they might change during revisions.

628 **Code availability**

629 All code can be found on GitHub at <https://github.com/RasmussenLab/vamb>
630 and is freely available under the permissive MIT license. The code
631 for making the figures is in the separate repository on Github
632 https://github.com/sgalkina/TaxVAMB_paper_figures.

633 **Competing interests**

634 J.N.N. is the author of the VAMB binning tool, which has been developed using
635 a prototype of BinBencher, which is used to calculate some of the benchmarking
636 metrics in this paper. S.R. is the founder and owner of the Danish company BioAI and
637 have performed consulting for Sidera Bio ApS. Other authors declare no competing
638 interests.

639 **References**

- 640 [1] Grünberger, F., Ferreira-Cerca, S. & Grohmann, D. Nanopore sequencing of RNA
641 and cDNA molecules in escherichia coli. *RNA* **28**, 400–417 (2022).
- 642 [2] Bickhart, D. M. *et al.* Generating lineage-resolved, complete metagenome-
643 assembled genomes from complex microbial communities. *Nat. Biotechnol.* **40**,
644 711–719 (2022).

- 645 [3] Feng, X., Cheng, H., Portik, D. & Li, H. Metagenome assembly of high-fidelity
646 long reads with hifiasm-meta. *Nat. Methods* **19**, 671–674 (2022).
- 647 [4] Sereika, M. *et al.* Oxford nanopore r10.4 long-read sequencing enables the gen-
648 eration of near-finished bacterial genomes from pure cultures and metagenomes
649 without short-read or reference polishing. *Nat. Methods* **19**, 823–826 (2022).
- 650 [5] Quince, C., Walker, A. W., Simpson, J. T., Loman, N. J. & Segata, N. Shotgun
651 metagenomics, from sampling to analysis. *Nat. Biotechnol.* **35**, 833–844 (2017).
- 652 [6] Albertsen, M. *et al.* Genome sequences of rare, uncultured bacteria obtained
653 by differential coverage binning of multiple metagenomes. *Nat. Biotechnol.* **31**,
654 533–538 (2013).
- 655 [7] Alneberg, J. *et al.* Binning metagenomic contigs by coverage and composition.
656 *Nat. Methods* **11**, 1144–1146 (2014).
- 657 [8] Imelfort, M. *et al.* GroopM: an automated tool for the recovery of population
658 genomes from related metagenomes. *PeerJ* **2**, e603 (2014).
- 659 [9] Wu, Y.-W., Simmons, B. A. & Singer, S. W. MaxBin 2.0: an automated
660 binning algorithm to recover genomes from multiple metagenomic datasets.
661 *Bioinformatics* **32**, 605–607 (2016).
- 662 [10] Kang, D. D. *et al.* MetaBAT 2: an adaptive binning algorithm for robust and
663 efficient genome reconstruction from metagenome assemblies. *PeerJ* **7**, e7359
664 (2019).
- 665 [11] Nissen, J. N. *et al.* Improved metagenome binning and assembly using deep
666 variational autoencoders. *Nat. Biotechnol.* **39**, 555–560 (2021).
- 667 [12] Teeling, H., Meyerdierks, A., Bauer, M., Amann, R. & Glöckner, F. O. Appli-
668 cation of tetranucleotide frequencies for the assignment of genomic fragments.
669 *Environ. Microbiol.* **6**, 938–947 (2004).

- 670 [13] Mallawaarachchi, V., Wickramarachchi, A. & Lin, Y. GraphBin: refined binning
671 of metagenomic contigs using assembly graphs. *Bioinformatics* **36**, 3307–3313
672 (2020).
- 673 [14] Zhang, Z. & Zhang, L. METAMVGL: a multi-view graph-based metagenomic
674 contig binning algorithm by integrating assembly and paired-end graphs. *BMC
675 Bioinformatics* **22**, 378 (2021).
- 676 [15] Lamurias, A., Sereika, M., Albertsen, M., Hose, K. & Nielsen, T. D. Metagenomic
677 binning with assembly graph embeddings. *Bioinformatics* **38**, 4481–4487 (2022).
- 678 [16] Lamurias, A., Tibo, A., Hose, K., Albertsen, M. & Nielsen, T. D. Krause,
679 A. *et al.* (eds) *Metagenomic binning using connectivity-constrained variational
680 autoencoders.* (eds Krause, A. *et al.*) *Proceedings of the 40th International
681 Conference on Machine Learning*, Vol. 202 of *Proceedings of Machine Learning
682 Research*, 18471–18481 (PMLR, 2023). URL <https://proceedings.mlr.press/v202/lamurias23a.html>.
- 684 [17] Yu, G., Jiang, Y., Wang, J., Zhang, H. & Luo, H. BMC3C: binning metage-
685 nomic contigs using codon usage, sequence composition and read coverage.
686 *Bioinformatics* **34**, 4172–4179 (2018).
- 687 [18] Lin, H.-H. & Liao, Y.-C. Accurate binning of metagenomic contigs via automated
688 clustering sequences using information of genomic signatures and marker genes.
689 *Sci. Rep.* **6**, 24175 (2016).
- 690 [19] Sieber, C. M. K. *et al.* Recovery of genomes from metagenomes via a dereplication,
691 aggregation and scoring strategy. *Nat Microbiol* **3**, 836–843 (2018).
- 692 [20] Pan, S., Zhu, C., Zhao, X.-M. & Coelho, L. P. A deep siamese neural network
693 improves metagenome-assembled genomes in microbiome datasets across different
694 environments. *Nat. Commun.* **13**, 2326 (2022).

- 695 [21] Pan, S., Zhao, X.-M. & Coelho, L. P. SemiBin2: self-supervised contrastive learn-
696 ing leads to better MAGs for short- and long-read sequencing. *Bioinformatics*
697 **39**, i21–i29 (2023).
- 698 [22] Wang, Z. *et al.* Effective binning of metagenomic contigs using contrastive multi-
699 view representation learning. *Nat. Commun.* **15**, 585 (2024).
- 700 [23] Krause, L. *et al.* Phylogenetic classification of short environmental DNA
701 fragments. *Nucleic Acids Res.* **36**, 2230–2239 (2008).
- 702 [24] Huson, D. H., Mitra, S., Ruscheweyh, H.-J., Weber, N. & Schuster, S. C. Inte-
703 grative analysis of environmental sequences using MEGAN4. *Genome Res.* **21**,
704 1552–1560 (2011).
- 705 [25] Wang, Z., Wang, Z., Lu, Y. Y., Sun, F. & Zhu, S. SolidBin: improving metagenome
706 binning with semi-supervised normalized cut. *Bioinformatics* **35**, 4229–4238
707 (2019).
- 708 [26] Uritskiy, G. V., DiRuggiero, J. & Taylor, J. MetaWRAP-a flexible pipeline for
709 genome-resolved metagenomic data analysis. *Microbiome* **6**, 158 (2018).
- 710 [27] Murovec, B., Deutsch, L. & Stres, B. Computational framework for High-Quality
711 production and Large-Scale evolutionary analysis of metagenome assembled
712 genomes. *Mol. Biol. Evol.* **37**, 593–598 (2020).
- 713 [28] Líndez, P. P., Johansen, J., Sigurdsson, A. I., Nissen, J. N. & Rasmussen, S.
714 Adversarial and variational autoencoders improve metagenomic binning (2023).
- 715 [29] Wickramarachchi, A. & Lin, Y. Carbone, A. & El-Kebir, M. (eds) *LRBinner:*
716 *Binning Long Reads in Metagenomics Datasets.* (eds Carbone, A. & El-Kebir,
717 M.) *21st International Workshop on Algorithms in Bioinformatics (WABI 2021)*,
718 Vol. 201 of *Leibniz International Proceedings in Informatics (LIPIcs)*, 11:1–11:18
719 (Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl, Germany, 2021).
720 URL <https://drops.dagstuhl.de/opus/volltexte/2021/14364>.

- 721 [30] Strous, M., Kraft, B., Bisdorf, R. & Tegetmeyer, H. E. The binning of metagenomic contigs for microbial physiology of mixed cultures. *Front. Microbiol.* **3**, 410 (2012).
- 724 [31] Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation initiation site identification. *BMC Bioinformatics* **11**, 119 (2010).
- 726 [32] Wood, D. E., Lu, J. & Langmead, B. Improved metagenomic analysis with kraken 2. *Genome Biology* **20** (2019). URL <https://doi.org/10.1186/s13059-019-1891-0>.
- 728 [33] Mirdita, M., Steinegger, M., Breitwieser, F., Söding, J. & Karin, E. L. Fast and sensitive taxonomic assignment to metagenomic contigs. *Bioinformatics* **37**, 3029–3031 (2021). URL <https://doi.org/10.1093/bioinformatics/btab184>.
- 731 [34] Lu, J., Breitwieser, F. P., Thielen, P. & Salzberg, S. L. Bracken: estimating species abundance in metagenomics data. *PeerJ Computer Science* **3**, e104 (2017). URL <https://doi.org/10.7717/peerj-cs.104>.
- 734 [35] Kim, D., Song, L., Breitwieser, F. P. & Salzberg, S. L. Centrifuge: rapid and sensitive classification of metagenomic sequences. *Genome Research* **26**, 1721–1729 (2016). URL <https://doi.org/10.1101/gr.210641.116>.
- 737 [36] Blanco-Miguez, A. *et al.* Extending and improving metagenomic taxonomic profiling with uncharacterized species with MetaPhlAn 4 (2022). URL <https://doi.org/10.1101/2022.08.22.504593>.
- 740 [37] Milanese, A. *et al.* Microbial abundance, activity and population genomic profiling with mOTUs2. *Nature Communications* **10** (2019). URL <https://doi.org/10.1038/s41467-019-08844-4>.
- 743 [38] Portik, D. M., Brown, C. T. & Pierce-Ward, N. T. Evaluation of taxonomic classification and profiling methods for long-read shotgun metagenomic sequencing datasets. *BMC Bioinformatics* **23** (2022). URL <https://doi.org/10.1186/s12859-022-05103-0>.

- 747 [39] Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W.
748 CheckM: assessing the quality of microbial genomes recovered from isolates, single
749 cells, and metagenomes. *Genome Res.* **25**, 1043–1055 (2015).
- 750 [40] Chklovski, A., Parks, D. H., Woodcroft, B. J. & Tyson, G. W. CheckM2: a rapid,
751 scalable and accurate tool for assessing microbial genome quality using machine
752 learning. *Nat. Methods* **20**, 1203–1212 (2023).
- 753 [41] Strathern, M. ‘improving ratings’: audit in the british university system. *European*
754 *Review* **5**, 305–321 (1997). URL [https://doi.org/10.1002/\(sici\)1234-981x\(199707\)5:3<305::aid-euro184>3.0.co;2-4](https://doi.org/10.1002/(sici)1234-981x(199707)5:3<305::aid-euro184>3.0.co;2-4)
- 755
- 756 [42] Kutuzova, S., Nielsen, M., Piera, P., Nissen, J. N. & Rasmussen, S. Taxometer:
757 Improving taxonomic classification of metagenomics contigs. *Nat. Commun.* **15**,
758 8357 (2024).
- 759 [43] Palumbo, E., Daunhawer, I. & Vogt, J. E. MMVAE+: ENHANCING THE GEN-
760 ERATIVE QUALITY OF MULTIMODAL VAES WITHOUT COMPROMISES
761 .
- 762 [44] Senellart, A., Chadebec, C. & Allassonnière, S. Improving multimodal joint
763 variational autoencoders through normalizing flows and correlation analysis
764 (2023).
- 765 [45] Hwang, H., Kim, G.-H., Hong, S. & Kim, K.-E. Multi-View representation learn-
766 ing via total correlation objective. *Adv. Neural Inf. Process. Syst.* **34**, 12194–12207
767 (2021).
- 768 [46] Sutter, T. M., Daunhawer, I. & Vogt, J. E. Generalized multimodal ELBO (2021).
- 769 [47] Shi, Y., N, Siddharth, Paige, B. & Torr, P. Variational Mixture-of-Experts autoen-
770 coders for Multi-Modal deep generative models. *Adv. Neural Inf. Process. Syst.*
771 **32** (2019).

- 772 [48] Wu, M. & Goodman, N. Multimodal generative models for scalable Weakly-
773 Supervised learning. *Adv. Neural Inf. Process. Syst.* **31** (2018).
- 774 [49] Suzuki, M., Nakayama, K. & Matsuo, Y. Joint multimodal learning with deep
775 generative models (2016).
- 776 [50] Wu, M. & Goodman, N. Multimodal Generative Models for Compositional
777 Representation Learning. *arXiv e-prints* arXiv:1912.05075 (2019).
- 778 [51] Kutuzova, S., Krause, O., McCloskey, D., Nielsen, M. & Igel, C. Multimodal
779 variational autoencoders for Semi-Supervised learning: In defense of Product-of-
780 Experts. *arXiv preprint arXiv:2101. 07240* (2021).
- 781 [52] Bromley, J., Guyon, I. & LeCun, Y. Signature verification using a siamese time
782 delay neural network. *Advances in neural information processing systems (NIPS)*
783 737–744.
- 784 [53] Valmadre, J. Oh, A. H., Agarwal, A., Belgrave, D. & Cho, K. (eds) *Hierarchical*
785 *classification at multiple operating points*. (eds Oh, A. H., Agarwal, A., Belgrave,
786 D. & Cho, K.) *Advances in Neural Information Processing Systems* (2022). URL
787 <https://openreview.net/forum?id=mNtFhoNRr4i>.
- 788 [54] Nissen, J. N., Lindéz, P. P. & Rasmussen, S. BinBencher: Fast, flexible and
789 meaningful benchmarking suite for metagenomic binning (2024).
- 790 [55] Chaumeil, P.-A., Mussig, A. J., Hugenholtz, P. & Parks, D. H. GTDB-tk v2: mem-
791 ory friendly classification with the genome taxonomy database. *Bioinformatics*
792 **38**, 5315–5316 (2022). URL <https://doi.org/10.1093/bioinformatics/btac672>.
- 793 [56] Mattock, J. & Watson, M. A comparison of single-coverage and multi-coverage
794 metagenomic binning reveals extensive hidden contamination. *Nat. Methods* **20**,
795 1170–1173 (2023).

- 796 [57] Almeida, A. *et al.* A unified catalog of 204,938 reference genomes from the human
797 gut microbiome. *Nat. Biotechnol.* **39**, 105–114 (2021).
- 798 [58] Ibrahim, E. *et al.* Biocontrol efficacy of endophyte *pseudomonas poae* to alleviate
799 fusarium seedling blight by refining the morpho-physiological attributes of wheat.
800 *Plants* **12** (2023).
- 801 [59] Li, X. *et al.* Exploration of phyllosphere microbiomes in wheat varieties with
802 differing aphid resistance. *Environ. Microbiome* **18**, 78 (2023).
- 803 [60] Mikiciński, A., Sobczewski, P., Puławska, J. & Maciorowski, R. Control of fire
804 blight (*erwinia amylovora*) by a novel strain 49M of *pseudomonas graminis* from
805 the phyllosphere of apple (*malus spp.*). *Eur. J. Plant Pathol.* **145**, 265–276 (2016).
- 806 [61] Robinson, R. K. (ed.) *Encyclopedia of food microbiology* (Academic Press, San
807 Diego, CA, 1999).
- 808 [62] Harada, H., Oyaizu, H., Kosako, Y. & Ishikawa, H. *Erwinia aphidicola*, a new
809 species isolated from pea aphid, *acyrthosiphon pisum*. *J. Gen. Appl. Microbiol.*
810 **43**, 349–354 (1997).
- 811 [63] Dougherty, P. E. *et al.* Widespread and largely unknown prophage activity, diver-
812 sity, and function in two genera of wheat phyllosphere bacteria. *ISME J.* **17**,
813 2415–2425 (2023).
- 814 [64] Manni, M., Berkeley, M. R., Seppey, M., Simão, F. A. & Zdobnov, E. M. BUSCO
815 update: Novel and streamlined workflows along with broader and deeper phylo-
816 genetic coverage for scoring of eukaryotic, prokaryotic, and viral genomes. *Mol.*
817 *Biol. Evol.* **38**, 4647–4654 (2021).
- 818 [65] Steinberg, G. Cell biology of *zymoseptoria tritici*: Pathogen cell organization and
819 wheat infection. *Fungal Genet. Biol.* **79**, 17–23 (2015).

- 820 [66] Mylonas, I., Stavrakoudis, D., Katsantonis, D. & Korpetis, E. in *Chapter 1 -*
821 *better farming practices to combat climate change* (eds Ozturk, M. & Gul, A.)
822 *Climate Change and Food Security with Emphasis on Wheat* 1–29 (Academic
823 Press, 2020).
- 824 [67] Parks, D. H. *et al.* GTDB: an ongoing census of bacterial and archaeal diversity
825 through a phylogenetically consistent, rank normalized and complete genome-
826 based taxonomy. *Nucleic Acids Res.* **50**, D785–D794 (2022).
- 827 [68] Sayers, E. W. *et al.* Database resources of the national center for biotechnology
828 information. *Nucleic Acids Res.* **50**, D20–D26 (2022).
- 829 [69] Subramanian, I., Verma, S., Kumar, S., Jere, A. & Anamika, K. Multi-omics
830 Data Integration, Interpretation, and Its Application (2020).
- 831 [70] Abedalrhman Alkhateeb, L. R. (ed.) *Machine Learning Methods for Multi-Omics*
832 *Data Integration* (Springer International Publishing, 2024).
- 833 [71] Allesøe, R. L. *et al.* Discovery of drug-omics associations in type 2 diabetes with
834 generative deep-learning models. *Nat. Biotechnol.* **41**, 399–408 (2023).
- 835 [72] Li, H. Aligning sequence reads, clone sequences and assembly contigs with bwa-
836 mem. *arXiv: Genomics* (2013). URL <https://api.semanticscholar.org/CorpusID:14669139>.
- 838 [73] Li, H. *et al.* The sequence Alignment/Map format and SAMtools. *Bioinformatics*
839 **25**, 2078–2079 (2009).
- 840 [74] Benoit, G. *et al.* Efficient High-Quality Metagenome Assembly from Long Accu-
841 rate Reads using Minimizer-space de Bruijn Graphs (2023). URL <https://www.biorxiv.org/content/10.1101/2023.07.07.548136v1>. Pages: 2023.07.07.548136
842 Section: New Results.

- 844 [75] Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*
845 **34**, 3094–3100 (2018). URL <https://doi.org/10.1093/bioinformatics/bty191>.
- 846 [76] Camargo, A. apcamargo/pycoverm: Simple Python interface to CoverM's fast
847 coverage estimation functions (2023). URL <https://github.com/apcamargo/pycoverm/tree/main>.
- 849 [77] Defazio, A. & Mishchenko, K. Learning-rate-free learning by d-adaptation. *The*
850 *40th International Conference on Machine Learning (ICML 2023)* (2023).
- 851 [78] Paszke, A. *et al.* PyTorch: An imperative style, high-performance deep learning
852 library. *33rd Conference on Neural Information Processing Systems (NeurIPS*
853 *2019)* (2019).
- 854 [79] Kim, J. & Steinegger, M. Metabuli: sensitive and specific metagenomic classifi-
855 cation via joint analysis of amino-acid and DNA (2023). 2023.
- 856 [80] Kim, D., Song, L., Breitwieser, F. P. & Salzberg, S. L. Centrifuge: rapid and
857 sensitive classification of metagenomic sequences. *Genome Res.* **26**, 1721–1729
858 (2016).
- 859 [81] Wood, D. E., Lu, J. & Langmead, B. Improved metagenomic analysis with kraken
860 2. *Genome Biol.* **20**, 257 (2019).
- 861 [82] Pan, S., Zhao, X.-M. & Coelho, L. P. Semibin2: self-supervised contrastive learn-
862 ing leads to better mags for short- and long-read sequencing. *bioRxiv* (2023).
863 URL <https://www.biorxiv.org/content/early/2023/01/09/2023.01.09.523201>.
- 864 [83] Liu, C.-C. *et al.* MetaDecoder: a novel method for clustering metagenomic contigs.
865 *Microbiome* **10**, 46 (2022).
- 866 [84] Quince, C. *et al.* STRONG: metagenomics strain resolution on assembly
867 graphs. *Genome Biology* **22**, 214 (2021). URL <https://doi.org/10.1186/s13059-021-02419-7>.

- 869 [85] Rinke, C. *et al.* Validation of picogram- and femtogram-input DNA libraries for
870 microscale metagenomics. *PeerJ* **4**, e2486 (2016).
- 871 [86] Chen, S. Ultrafast one-pass fastq data preprocessing, quality control, and
872 deduplication using fastp **1**, e107 (2023). URL <https://doi.org/10.1002/imt2.107>.
- 873 [87] Ewels, P., Magnusson, M., Lundin, S. & Käller, M. MultiQC: summarize analysis
874 results for multiple tools and samples in a single report. *Bioinformatics* **32**,
875 3047–3048 (2016). URL <https://doi.org/10.1093/bioinformatics/btw354>.
- 876 [88] Langmead, B. & Salzberg, S. L. Bowtie 2: fast and sensitive read alignment.
877 *Nature methods* **9**, 357–359 (2012).
- 878 [89] Prjibelski, A., Antipov, D., Meleshko, D., Lapidus, A. & Korobeynikov, A. Using
879 spades de novo assembler. *Current Protocols in Bioinformatics* **70**, e102 (2020).
880 URL <https://doi.org/10.1002/cpbi.102>.
- 881 [90] Mikheenko, A., Prjibelski, A., Saveliev, V., Antipov, D. & Gurevich, A. Versatile
882 genome assembly evaluation with quast-lg. *Bioinformatics* **34**, i142–i150 (2018).
883 URL <https://doi.org/10.1093/bioinformatics/bty266>.
- 884 [91] Xu, S. *et al.* Ggtree: A serialized data object for visualization of a phylogenetic
885 tree and annotation data. *Imeta* **1**, e56 (2022).
- 886 [92] Kutuzova, S. *et al.* Wheat phyllosphere metagenome assembled genomes collected
887 in Ringsted, Denmark (2024). URL <https://doi.org/10.5281/zenodo.13959411>.