

1 Advances in Metagenomic Binning for the reconstruction of  
2 microbial species

3 Jose Fernando Garcia Guevara<sup>1\*</sup>, Theo Portlock<sup>1\*</sup>, Adil Mardinoglu<sup>1,2</sup>, Mathias  
4 Uhlén<sup>1</sup>, and Saeed Shoaie<sup>1,2</sup>

5 <sup>1</sup>Science for Life Laboratory, Royal Institute of Technology (KTH), Stockholm,  
6 Sweden.

7 <sup>2</sup>Centre for HostMicrobiome Interactions, Faculty of Dentistry, Oral &  
8 Craniofacial Sciences, King's College London, London, UK.

9 **Contents**

10	<b>1 Abstract</b>	<b>2</b>
11	<b>2 Background</b>	<b>2</b>
12	<b>3 Overview of recent methods for metagenomic binning</b>	<b>3</b>
13	3.1 Progress in recent binning strategies . . . . .	3
14	3.1.1 Binning co-abundant genes . . . . .	4
15	3.1.2 Binning microbial genomes with deep learning . . . . .	5
16	3.2 Binning of viral genomes . . . . .	5
17	3.3 Binning Pipelines . . . . .	6
18	<b>4 Suggestions on choosing a binning algorithm</b>	<b>7</b>
19	<b>5 Conclusion</b>	<b>7</b>
20	<b>References</b>	<b>9</b>

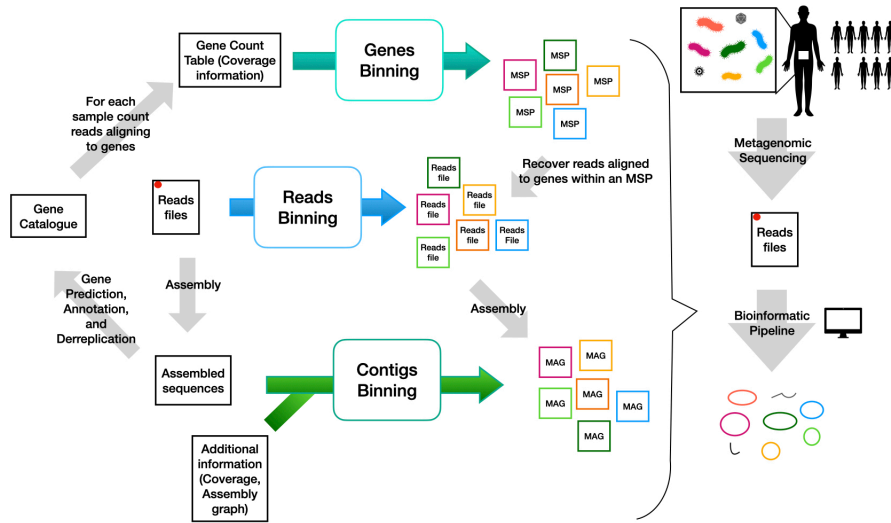


Figure 1: Summary of binning principles and techniques.

## 1 Abstract

New generations of sequencing platforms coupled with numerous bioinformatics tools have led to rapid technological progress in metagenomics to investigate complex microorganism communities. Nevertheless, a combination of different bioinformatic tools remains necessary to draw conclusions out of microbiota studies. As sequencing costs have dropped at a rate above 'Moore's law', a greater number of large data sets are being produced than ever before. Newer algorithms that take advantage of the size of these datasets are continually being developed. Binning algorithms are defined as the grouping of assembled metagenomic contigs by their genome of origin (Figure 1). Selecting the most appropriate binning algorithm can be a daunting task and is influenced by many factors. This review serves as a guide to direct the researcher to the binning algorithm that best suits their needs.

## 2 Background

The explosion in popularity and success in the field of metagenomics over the last 25 years can be largely attributed to the advances in computing technologies. An example of the outcomes of this can be found in the Human Microbiome Project; a project that has been greatly improved the understanding of the microbial flora involved in human health and disease (Turnbaugh et al., 2007). These advances have brought with them greater demands for storage, CPU time, and consequently

38 more efficient algorithms. The main function of binning tools is to reconstruct species/biological  
39 entities from metagenomic samples. Compared to amplicon, shotgun metagenome can provide  
40 functional gene profiles directly and reach a much higher resolution of taxonomic annotation.  
41 However, due to the high demands on computational resources, cost, and expertise necessary  
42 to perform this analysis, researchers have historically been limited in their capacity to collect  
43 and analyse sequencing data. As the cost of sequencing is rapidly falling, this burden has been  
44 significantly lessened. Whole Genome Shotgun sequences does not require cultivation. At the  
45 time of writing, shotgun metagenomic sequencing costs on average three times as much as 16S  
46 sequencing in comparison. Here we will briefly recapitulate recent binning algorithms and highlight  
47 some of the developments in the field, among them, the use of new algorithms and strategies  
48 employed to achieve the goal of identifying the organisms composing microbiome communities.  
49 We hope this overview could aid the reader to choose a binning algorithm or a combination of  
50 them based on their specific needs.

## 51 **3 Overview of recent methods for metagenomic binning**

### 52 **3.1 Progress in recent binning strategies**

53 A metagenomic sample is comprised of many organisms and the goal of binning is to reconstruct  
54 the sequences from each organism present in the original sample. The majority of binning tools  
55 are oriented toward clustering contigs (contig-binning) into bins, which may represent the genome  
56 from a single biological entity/organism. A Metagenome-Assembled Genome (MAG) is a single-  
57 taxon assembly based on one or more binned metagenomes that has been asserted to be a close  
58 representation to an individual genome that could match an already existing isolate or represent  
59 a novel isolate. Current contig-binning tools are commonly reference free (i.e. they do not de-  
60 pend on reference sequences to perform clustering) and rely on coverage information and sequence  
61 composition. Progress in contig-binning algorithms can be seen in the proposals to integrate  
62 new sources of information (for example, from scaffold-graphs (Binnacle), paired-end reads (CO-  
63 CACOLA), or 3D contact information (MetaTOR)) and state of the art algorithms in machine  
64 learning (CoCoNet, VAMB). We also notice the development of Bin refinement tools (DAS-tool,  
65 Binning Refiner) that rely on the outputs from multiple contig-binning algorithms and combine  
66 them to produce better results (Sieber et al., 2018). Binning of contigs have played a central role  
67 in software development in the field, a review on the benchmarking of binning algorithms was done

by Yue et al., 2020. Beside contig-binning tools we can also distinguish read-binning tools and co-abundant-gene-binning tools. The main purpose of read-binning tools is to pre-process reads into clusters for a posterior targeted assembly. Here we find reference-free and non-reference-free tools, and tools designed for short-read or long-read sequencing technologies. Among the binning tools developed in recent years, a subset of them are dedicated to cluster reads (read-binning) (MetaBBC-LR, BioBloom Tools, CLAME, LVQ-KNN, Meta VW, HirBin, MEGAN-LR) (Wickramarachchi, Mallawaarachchi, Rajan, & Lin, 2020; Chu et al., 2014; Benavides, Isaza, Niño-García, Alzate, & Cabarcas, 2018; Belka, Fischer, Pohlmann, Beer, & Höper, 2018; Vervier, Mahé, Tournoud, Veyrieras, & Vert, 2016; Österlund, Jonsson, & Kristiansson, 2017; Huson et al., 2018).

### 3.1.1 Binning co-abundant genes

Binning of co-abundant genes represents an alternative proposal to reconstruct species/biological entities from a set of metagenomic samples. Co-abundant gene binning methods assumes that each gene coming from a shared chromosome will display proportional abundances across samples. Therefore, if there are enough samples from a similar environment you can identify the sets of genes from a common organism of origin (MLGs Chameleon-clust 2012, CAGs and MGSs Canopy 2014, Markovclust-MGCs Karlsson 2013, MSPs MSPminner 2018) (Karypis, Han, & Kumar, 1999; Plaza Oñate et al., 2019). The MSPminer software was developed to exploit this approach. MSPminer introduced a robust proportionality measure to detect co-abundance but no necessarily co-occurrence. This tool groups co-abundant genes into Metagenomic Species Pan-genomess (MSPs) and classifies genes within an MSP as core, accessory and shared. Core genes are present in all strains, accessory are present only in some, and the shared category applies for those genes which may be present in more than one MSP due to horizontal transfer (Tettelin et al., 2005). Factors that impact directly on MSP quality include sample composition, sequencing depth, and previous bioinformatic steps to build the reference gene dataset and map the reads. MSPs can be used for taxonomic profiling of new samples from similar ecosystems at the species level, and also to compare strains between samples by building a presence/absence table of accessory genes and for biomarker discovery. By binning contigs carrying genes from the same MSP it is also possible to build a MAG.

### 3.1.2 Binning microbial genomes with deep learning

The integration of deep learning techniques has revolutionised the field of metagenomics. Deep learning approaches have benefitted from the rapid acceleration in GPU efficiency over the past few years. The Software VAMB and CoCoNet constitute two such examples that employ deep learning for binning (Nissen et al., n.d.; Arisdakessian, Nigro, Steward, Poisson, & Belcaid, 2021).

The main novelty of VAMB is the application of the Deep Learning technique known as Variational Auto Encoders (VAE). In this case, variational autoencoders learn how to integrate two data types, coabundance and k-mer composition. The resulting latent representation is able to cluster better than either of the inputs alone. In principle this technology is not limited by only two input data types. VAMB also applies a "multisplit" approach whereby each cluster should correspond to an organism representation across samples and each bin in a cluster to a per-sample representation of the genome of that organism.

The CoCoNet software uses deep learning and clustering to bin contigs into clusters representing species present in the samples. The algorithm consists of two phases. During the first phase, a neural network is trained to estimate the probability that two contigs arise from the same genome given their composition and coverage information. The second use a heuristic to bin the contigs using the probabilities inferred in the first stage. An interesting feature in CoCoNet is it was trained on viral genomes. In the following section we discuss more about binning on viral genomes.

## 3.2 Binning of viral genomes

Most binning algorithms are designed for prokaryotic organisms leaving viruses out of the software scope. In recent years the virome and its importance in health and disease has recognised. CoCoNet uses deep learning to model co-occurrence of contigs from the same viral genome. The network was optimized for diverse viral metagenomes, the network learns to model coverage variability within samples, a critical feature in viral metagenomes where DNA amplification methods are needed to increase input genetic material. VirBin clusters contigs for genome reconstruction of viral strains, different strains within viral species may show different biological properties such as transmissibility or virulence. Composition based features are usually not enough to separate haplotypes, VirBin receives contigs as inputs and outputs the estimated number of haplotypes via contig alignment and returns the contigs for each haplotype based on relative abundance distribution, when the contigs are long enough VirBin produce better results. Newer strategies has been

128 proposed and employed to reconstruct viral genomes from metagenomic samples, in a recent work  
129 (Nayfach et al., 2021) a new compendium of 189680 DNA viruses from the human gut microbiome  
130 was produced. In this work they use viral informative features including presence of viral protein  
131 families (Paez-Espino et al., 2016), absence of non-viral families (El-Gebali et al., 2019), gene  
132 strand switch rate (Roux 2015), and the score produced from the VirFinder (Ren, Ahlgren, Lu,  
133 Fuhrman, & Sun, 2017) software.

### 134 **3.3 Binning Pipelines**

135 Other advances in binning can be found in the integration of existing tools and software into  
136 bioinformatic pipelines. These innovations allow the automatic complete processing of read sam-  
137 ples into bins or the addition of extra processing steps to address specific biological questions or  
138 problems related to the sample of origin. MetaWRAP is a modular pipeline ready to perform  
139 common tasks in metagenomic analysis, starting from read quality checks up to bin creation,  
140 refinement, reassembly quantification, taxonomic annotation and functional annotation. MAGO  
141 pipeline integrates metagenome assembly, binning, bin improvement, bin quality check, bin func-  
142 tional annotation, and bin taxonomic annotation. SqueezeMeta also integrates external software  
143 to perform the complete analysis of metagenomic samples from sequences reading to MAG con-  
144 struction and annotation (Tamames & Puente-Sánchez, 2019) nf-mag supports both short and  
145 long reads, performs quality and adapter trimming, quality check, performs assembly, binning,  
146 checks bin quality and assigns taxonomy (Ewels et al., 2020). Autometa was developed to deal  
147 with non-model Eukariotic host contamination and complex single metagenomes, the application  
148 integrate sequence homology, nucleotide composition, coverage and single-copy marker genes to  
149 separate microbial genomes from non model host genomes (Miller et al., 2019). Seqdex is a tool  
150 written in R which separates endosymbionts from their host sequences (Chiodi et al., 2019). Their  
151 approach uses specific features in endosymbiotic systems to better solve this problem. This tool  
152 combines partial taxonomic annotations obtained trough homology searches and sequence compo-  
153 sition to predict the contig’s organism of origin from host and its endosymbionts and helps the user  
154 to understand how effective is the classification. Reproducibility, scalability, and ease of use from  
155 people with little computational experience are attractive features that pipelines for metagenomic  
156 analysis provide.

## 157 4 Suggestions on choosing a binning algorithm

158 A number of aspects should be considered when performing binning analysis on metagenomic  
159 samples. Computational resources available, sequencing technology, number of samples, and the  
160 sample's source are important factors to consider. Some tools employ more resources than others,  
161 and some perform better under specific circumstances (as reviewed by Yue et al., 2020). If you are  
162 dealing with a large number of samples, a gene-binning strategy could be taken into consideration.  
163 Tools such as CoMet were built around single sample binning (Herath, Tang, Tandon, Ackland,  
164 & Halgamuge, 2017). Long read sequence technology is gaining momentum and some tools also  
165 integrate the characteristic features generated with this technology. The environment under study  
166 also play an important role for binning. Sometimes there exists host organisms whose genome  
167 sequences would be removed before starting the analysis. The environment also has a profound ef-  
168 fect on the sample's diversity with samples that have greater diversity requiring greater sequencing  
169 depth making binning more difficult. It is also difficult for binning tools to discern between similar  
170 strains within the same sample. It is also worth mentioning that there is no mutual exclusivity  
171 between the currently available tools and it is possible to benefit from the relative advantages each  
172 has to offer and merge the results depending of the aim of the study. Besides binning, other types  
173 of metagenomic analysis can be performed on microbiomes. Recent reviews provide an overview of  
174 the complete process and practical guides to apply available software (Breitwieser, Lu, & Salzberg,  
175 2019).

## 176 5 Conclusion

177 Popularity and successes of metagenomic binning have accelerated in the last ten years. Current  
178 limitations that still remain include the difficulty in classifying similar strains within samples.  
179 They additionally do not perform well assigning 16S sequences to bins likely due to the high copy  
180 number of these sequences within a genome. As binning has been focused mainly in prokariotic  
181 organisms, binning of organisms outside prokariotes need more development. Although there  
182 have been signifiante advances in the characterisation of viral genomes as of late (Nayfach et al.,  
183 2021), the huge diversity in viral genomes still poses a challenge for current methodologies. The  
184 continuously increasing number of sequences available require more efficient/faster algorithms and  
185 new strategies to reconstruct single organisms from environmental samples. However, with the  
186 breakneck pace of technological advancements in computing resources, this requirement is sure to

Table 1: Comparison of binning algorithms

Software/Algorithm	Year	Description/purpose	Comment/Highlight	Doi	PubmedID
CoCoNet	2021	Deep learning tool for Viral Metagenome Binning	Reconstructs viral genomes	10.1093/bioinformatics/btab213	33822891
Binnacle	2021	Using scaffolds to improve Metagenomic bin qua...	Incorporates scaffold information	10.3389/fmicb.2021.638561	33717033
VAMB	2021	Metagenome binning using deep variational auto...	Autoencoder algorithm, fast processing	10.1038/s41587-020-00777-4	33398153
phyloFlash	2020	ssrRNA profiling and MAG assembly	Incorporates ssrRNA profiling info into MAG as...	10.1128/mSystems.00920-20	33109753
MetaBCC-LR	2020	Metagenomic binning for Long-Reads	Suitable for Long Reads sequencing technology	10.1093/bioinformatics/btaa441	32657364
BioBloom Tools	2020	Reads binning for targeted assembly, alignment...	Data preparation for targeted assembly, using ...	10.1073/pnas.1903436117	32641514
GraphBin	2020	Refined binning of metagenomic contigs using a...	Incorporates assembly graphs information	10.1093/bioinformatics/btaa180	32167528
MetaSPSim	2020	Simulating metagenomic stable isotope probing ...	Augment binning resolution with extra experime...	10.1186/s12859-020-3372-6	32000676
MetaCon	2019	Unsupervised binning k-mers and coverage, focu...	Focus different lengths contigs	10.1186/s12859-019-2904-4	31757198
VirBin	2019	Binning viral haplotypes from assembled contigs	Viral haplotypes MAGs	10.1186/s12859-019-3138-1	31684876
MAGO (*only tool pipeline)	2019	Framework for Production and analysis of MAGs	Identification of endosymbiont	10.1093/molbev/msz237	31633780
SeqDex	2019	Genome separation of Endosymbionts from mixed ...	Incorporates 3D contact information	10.3389/fgene.2019.00853	31608107
MetaTOR	2019	High quality MAGs from mammalian guts using me...	Eliminates manual parameter tuning from previo...	10.3389/fgene.2019.00753	31481973
MetaBAT (v2)	2019	Adaptive binning algorithm for genome recon...	Employs sample X contigs cf mapped read counts	10.7717/peerj.7359	31388474
MetaBMF	2019	Scalable binning algorithm for large scale met...	Haplotypes for polyploid genomes	10.1093/bioinformatics/btz577	31347687
PolyCRACKER	2019	Method for partitioning polyploid sub genomes ...	NaN	10.1186/s12864-019-5828-5	31299888
SolidBin	2019	Improving metagenome binning with semi-supervi...	NaN	10.1093/bioinformatics/btz253	30977806
Autometa	2019	Improvement of microbial genomes from individua...	Handles eukaryotic contamination	10.1093/nar/gkz148	30838416
MLBP MrGBP (Algorithm)	2019	Signal processing method for alignment free me...	Alternative description of sequences designed ...	10.1038/s41598-018-38197-9	30770850
CLAME	2018	Alignement based algorithm allowed description o...	Alignement based for reads	10.1186/s12864-018-5191-y	30537931
3D BH SNE (Algorithm)	2018	Fuzzy binning of metagenomic sequence fragments	Horizontal gene transfer and regions of uncert...	10.1109/EMBC.2018.8512529	30440633
LVQ-KNN	2018	Composition based RNA or DNA binning of short ...	Classify into DNA or RNA sequence	10.1016/j.virusres.2018.10.002	30291874
MSPminer	2018	Abundance based reconstitution of microbial pa...	Pan genome reconstitution	10.1093/bioinformatics/bty830	30252023
MetaWRAP*	2018	Flexible pipeline for genome resolved metageno...	Hybrid bin extraction algorithm	10.1186/s40168-018-0541-1	30219103
MetaVW	2018	Large scale Machine Learning Sequence classifi...	Machine learning for reads based on Khmer profile	10.1007/978-1-4939-8561-6_2	30030800
Opal (algorithm*)	2018	Metagenomic binning through low density binning	Improvement at higher taxonomic levels, discov...	10.1093/bioinformatics/bty611	30010790
BMC3C	2018	Binning contigs using codon usage sequence com...	Add codon usage information	10.1093/bioinformatics/bty519	29947757
AMBER tool	2018	Assessment of Metagenome Binner	NaN	10.1093/gigascience/giy069	29893851
DAS Tool	2018	Derreplication aggregation and scoring strategy	Combines several binning algorithm results	10.1038/s41564-018-0171-1	29807988
MEGAN-LR	2018	Long Read/ contigs taxonomic binning	Alignement of long reads against reference seque...	10.1186/s13062-018-0208-7	29678199
CoMet	2017	Binning workflow using contain coverage and co...	Single sample, include gc content and 4mer fr...	10.1186/s12859-017-1967-3	29297295
MetaGen	2017	reference-free learning with multiple metageno...	Requires multiple samples	10.1186/s13059-017-1323-y	28974263
d2sBin add onn	2017	Improved formula for calculate oligonucleotide...	Math formula to calculate oligo sequence dissi...	10.1186/s12859-017-1835-1	28931373
BusyBee Web	2017	Bootstrapped supervises binning and annotation	2d interactive scatterplots supervised binning	10.1093/nar/gkx348	28472498
ICoVer	2017	High resolution identification of differential...	Interactive visualisation tool	10.1186/s12859-017-1653-5	28464793
HirBin*	2017	Unsupervised clustering tool for verification...	Supervised annotation, unsupervised clustering...	10.1186/s12864-017-3686-6	28431529
BinSanity	2017	Improve genome bins through the combination of...	Reduce bias for high/low abundance	10.7717/peerj.3035	28289564
Binning_refinner	2017	Improve genome bins through the combination of...	Combination of different binning algorithms	10.1093/bioinformatics/btx086	28186226
IFCM add on	2016	Improved binning using Fuzzy C-Means Method	Add estimated distribution of real genome lengths	10.1109/TCBB.2016.2576452	27295684
COCACOLA	2016	binning contigs using composition, read covera...	Adds paired end read and coalignment information	10.1093/bioinformatics/btw290	27256312
GroopM (v2)	2014	Tool for automatic recovery of population geno...	Adds differential coverage to complement compo...	10.7717/peerj.603	25289188



187 be met and will pave the way for greater insights into the microbial world. With the integration  
 188 of Machine learning algorithms into binning, we expect to see significant developments in the near  
 189 future.

## 190 References

- 191 Arisdakessian, C. G., Nigro, O. D., Steward, G. F., Poisson, G., & Belcaid, M. (2021). Coconet:  
 192 an efficient deep learning tool for viral metagenome binning. *Bioinformatics*.
- 193 Belka, A., Fischer, M., Pohlmann, A., Beer, M., & Höper, D. (2018). Lvq-knn: Composition-based  
 194 dna/rna binning of short nucleotide sequences utilizing a prototype-based k-nearest neighbor  
 195 approach. *Virus research*, 258, 55–63.
- 196 Benavides, A., Isaza, J. P., Niño-García, J. P., Alzate, J. F., & Cabarcas, F. (2018). Clame: a new  
 197 alignment-based binning algorithm allows the genomic description of a novel xanthomon-  
 198 adaceae from the colombian andes. *BMC genomics*, 19(8), 9–30.
- 199 Breitwieser, F. P., Lu, J., & Salzberg, S. L. (2019). A review of methods and databases for  
 200 metagenomic classification and assembly. *Briefings in bioinformatics*, 20(4), 1125–1136.
- 201 Chiodi, A., Comandatore, F., Sassera, D., Petroni, G., Bandi, C., & Brilli, M. (2019). Seqdex: a  
 202 sequence deconvolution tool for genome separation of endosymbionts from mixed sequencing  
 203 samples. *Frontiers in genetics*, 10, 853.
- 204 Chu, J., Sadeghi, S., Raymond, A., Jackman, S. D., Nip, K. M., Mar, R., ... Birol, I. (2014).  
 205 Biobloom tools: fast, accurate and memory-efficient host species sequence screening using  
 206 bloom filters. *Bioinformatics*, 30(23), 3402–3404.
- 207 El-Gebali, S., Mistry, J., Bateman, A., Eddy, S. R., Luciani, A., Potter, S. C., ... others (2019).  
 208 The pfam protein families database in 2019. *Nucleic acids research*, 47(D1), D427–D432.
- 209 Ewels, P. A., Peltzer, A., Fillinger, S., Patel, H., Alneberg, J., Wilm, A., ... Nahnsen, S. (2020).  
 210 The nf-core framework for community-curated bioinformatics pipelines. *Nature biotechnol-  
 211 ogy*, 38(3), 276–278.
- 212 Herath, D., Tang, S.-L., Tandon, K., Ackland, D., & Halgamuge, S. K. (2017). Comet: a workflow  
 213 using contig coverage and composition for binning a metagenomic sample with high precision.  
 214 *BMC bioinformatics*, 18(16), 161–172.
- 215 Huson, D. H., Albrecht, B., Bağcı, C., Bessarab, I., Gorska, A., Jolic, D., & Williams, R. B.  
 216 (2018). Megan-lr: new algorithms allow accurate binning and easy interactive exploration  
 217 of metagenomic long reads and contigs. *Biology direct*, 13(1), 1–17.

218 Karypis, G., Han, E.-H., & Kumar, V. (1999). Chameleon: Hierarchical clustering using dynamic  
219 modeling. *Computer*, 32(8), 68–75.

220 Miller, I. J., Rees, E. R., Ross, J., Miller, I., Baxa, J., Lopera, J., ... Kwan, J. C. (2019). Au-  
221 tometa: automated extraction of microbial genomes from individual shotgun metagenomes.  
222 *Nucleic acids research*, 47(10), e57–e57.

223 Nayfach, S., Páez-Espino, D., Call, L., Low, S. J., Sberro, H., Ivanova, N. N., ... others (2021).  
224 Metagenomic compendium of 189,680 dna viruses from the human gut microbiome. *Nature*  
225 *Microbiology*, 1–11.

226 Nissen, J. N., Johansen, J., Allesøe, R. L., Sønderby, C. K., Armenteros, J. J. A., Grønbech,  
227 C. H., ... others (n.d.). Improved metagenome binning and assembly using deep variational  
228 autoencoders. *Nature Biotechnology*, 1–6.

229 Österlund, T., Jonsson, V., & Kristiansson, E. (2017). Hirbin: high-resolution identification of  
230 differentially abundant functions in metagenomes. *BMC genomics*, 18(1), 1–11.

231 Paez-Espino, D., Eloie-Fadrosh, E. A., Pavlopoulos, G. A., Thomas, A. D., Huntemann, M.,  
232 Mikhailova, N., ... Kyrpides, N. C. (2016). Uncovering earth’s virome. *Nature*, 536(7617),  
233 425–430.

234 Plaza Oñate, F., Le Chatelier, E., Almeida, M., Cervino, A. C., Gauthier, F., Magoulès, F., ...  
235 Pichaud, M. (2019). Mspminer: abundance-based reconstitution of microbial pan-genomes  
236 from shotgun metagenomic data. *Bioinformatics*, 35(9), 1544–1552.

237 Ren, J., Ahlgren, N. A., Lu, Y. Y., Fuhrman, J. A., & Sun, F. (2017). Virfinder: a novel k-mer  
238 based tool for identifying viral sequences from assembled metagenomic data. *Microbiome*,  
239 5(1), 1–20.

240 Sieber, C. M., Probst, A. J., Sharrar, A., Thomas, B. C., Hess, M., Tringe, S. G., & Banfield,  
241 J. F. (2018). Recovery of genomes from metagenomes via a dereplication, aggregation and  
242 scoring strategy. *Nature microbiology*, 3(7), 836–843.

243 Tamames, J., & Puente-Sánchez, F. (2019). Squeezemeta, a highly portable, fully automatic  
244 metagenomic analysis pipeline. *Frontiers in microbiology*, 9, 3349.

245 Tettelin, H., Massignani, V., Cieslewicz, M. J., Donati, C., Medini, D., Ward, N. L., ... others  
246 (2005). Genome analysis of multiple pathogenic isolates of streptococcus agalactiae: impli-  
247 cations for the microbial “pan-genome”. *Proceedings of the National Academy of Sciences*,  
248 102(39), 13950–13955.

249 Turnbaugh, P. J., Ley, R. E., Hamady, M., Fraser-Liggett, C. M., Knight, R., & Gordon, J. I.

250 (2007). The human microbiome project. *Nature*, 449(7164), 804–810.  
 251 Vervier, K., Mahé, P., Tournoud, M., Veyrieras, J.-B., & Vert, J.-P. (2016). Large-scale machine  
 252 learning for metagenomics sequence classification. *Bioinformatics*, 32(7), 1023–1032.  
 253 Wickramarachchi, A., Mallawaarachchi, V., Rajan, V., & Lin, Y. (2020). Metabcc-lr: meta  
 254 genomics binning by coverage and composition for long reads. *Bioinformatics*,  
 255 36(Supplement\_1), i3–i11.  
 256 Yue, Y., Huang, H., Qi, Z., Dou, H.-M., Liu, X.-Y., Han, T.-F., ... Tu, J. (2020). Evaluating  
 257 metagenomics tools for genome binning with real metagenomic datasets and caml datasets.  
 258 *BMC bioinformatics*, 21(1), 1–15.