

1 Metagenomic Binning Pipelines - the State of the Art

2

3 **Contents**

4	<b>1 Abstract</b>	<b>1</b>
5	<b>2 Background</b>	<b>2</b>
6	<b>3 Overview of recent methods for metagenomic binning</b>	<b>3</b>
7	3.1 Progress in recent binning strategies . . . . .	3
8	3.1.1 Binning co-abundant genes . . . . .	4
9	3.1.2 Binning microbial genomes with deep learning . . . . .	4
10	3.2 Binning of viral genomes . . . . .	5
11	3.3 Binning Pipelines . . . . .	5
12	<b>4 Suggestions on choosing a binning algorithm</b>	<b>6</b>
13	<b>5 Conclusion</b>	<b>8</b>
14	<b>References</b>	<b>8</b>

15 **1 Abstract**

16 New generations of sequencing platforms coupled with numerous bioinformatics tools have led to  
17 rapid technological progress in metagenomics to investigate complex microorganism communities.  
18 Nevertheless, a combination of different bioinformatic tools remains necessary to draw conclusions  
19 out of microbiota studies. As sequencing costs have dropped at a rate above 'Moore's law', a  
20 greater number of large data sets are being produced than ever before. Newer algorithms that  
21 take advantage of the size of these datasets are continually being developed. Binning algorithms  
22 are defined as the grouping of assembled metagenomic contigs by their genome of origin (Figure 1).

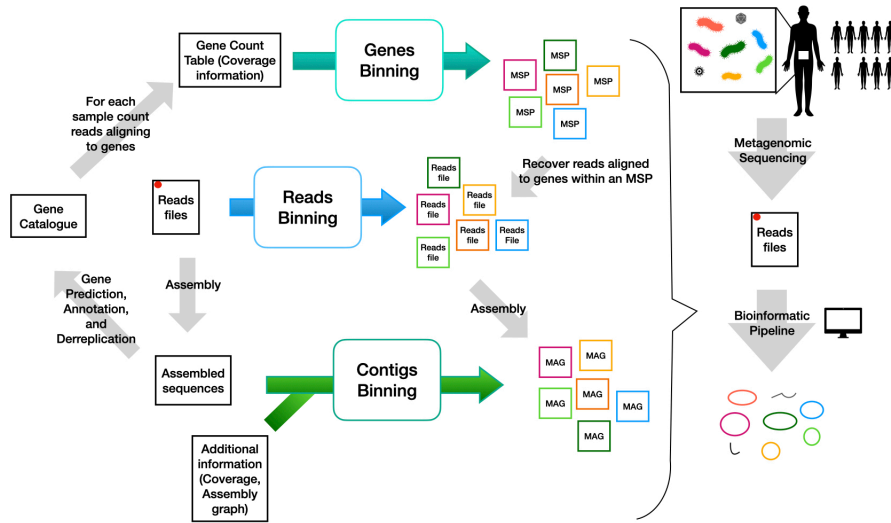


Figure 1: Summary of binning principles and techniques.

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. These advances have brought with them greater demands for storage, CPU time, and consequently more efficient algorithms. The main function of binning tools is to reconstruct species/biological entities from metagenomic samples. Compared to amplicon, shotgun metagenome can provide functional gene profiles directly and reach a much higher resolution of taxonomic annotation. However, due to the high demands on computational resources, cost, and expertise necessary to perform this analysis, researchers have historically been limited in their capacity to collect and analyse sequencing data. As the cost of sequencing is rapidly falling, this burden has been significantly lessened. Whole Genome Shotgun sequences does not require cultivation. At the time of writing, shotgun metagenomic sequencing costs on average three times as much as 16S sequencing in comparison. Here we will briefly recapitulate recent binning algorithms and highlight some of the developments in

41 the field, among them, the use of new algorithms and strategies employed to achieve the goal of  
42 identifying the organisms composing microbiome communities. We hope this overview could aid  
43 the reader to choose a binning algorithm or a combination of them based on their specific needs.

## 44 **3 Overview of recent methods for metagenomic binning**

### 45 **3.1 Progress in recent binning strategies**

46 A metagenomic sample is comprised of many organisms and the goal of binning is to reconstruct  
47 the sequences from each organism present in the original sample. The majority of binning tools  
48 we can find are oriented toward clustering contigs (contig-binning) into bins, which may represent  
49 the genome from a single biological entity/organism. A Metagenome-Assembled Genome (MAG)  
50 is a single-taxon assembly based on one or more binned metagenomes that has been asserted to  
51 be a close representation to an actual individual genome (that could match an already existing  
52 isolate or represent a novel isolate).

53 Current contig-binning tools normally are reference free (i.e do not depend on reference se-  
54 quences to perform clustering) and rely on coverage information and sequence composition. Progress  
55 in contig-binning algorithms can be seen in the proposals to integrate new sources of information  
56 (for example, from scaffold-graphs(Binnacle), paired-end reads(COCACOLA), or 3D contact in-  
57 formation(MetaTOR)) and state of the art algorithms in machine learning (CoCoNet, Variational  
58 Autoencoders for Metagenomic Binning (VAMB)).

59 We also notice the development of Bin refinement tools (DAS-tool, Binning Refiner), this tools  
60 rely on the outputs from multiple contig-binning algorithms and attempt to combine them to  
61 produce better results.

62 Binning of contigs have played a central role in software development in the field, a review on  
63 the benchmarking binning algorithms was done by Yue et al., 2020.

64 Beside contig-binning tools we can also distinguish read-binning tools and co-abundant-gene-  
65 binning tools.

66 The main purpose of read-binning tools is to pre-process reads into clusters for a posterior  
67 targeted assembly, here we find reference-free and non-reference-free tools, and tools designed for  
68 short-read or long-read sequencing technologies. Among the binning tools developed in recent  
69 years a subset of them are dedicated to cluster reads (read-binning) (MetaBBC-LR, BioBloom  
70 Tools, CLAME, LVQ-KKN, Meta VW, HirBin, MEGAN-LR).

### 71 3.1.1 Binning co-abundant genes

72 Binning of co-abundant genes represents an alternative proposal to reconstruct species/biological  
73 entities from a set of metagenomic samples. Co-abundant gene binning methods assume each gene  
74 coming from a shared chromosome will display proportional abundances across samples, if you have  
75 enough samples from a similar environment you can identify the sets of genes from a common  
76 organism of origin (MLGs Chameleon-clust 2012, CAGs and MGSs Canopy 2014, Markovclust-  
77 MGCs Karlsson 2013, MSPs MSPminner 2018) (Karypis, Han, & Kumar, 1999; Plaza Oñate et  
78 al., 2019). In the past few years the MSPminer software was developed exploiting this approach.  
79 MSPminer introduced a robust proportionality measure detecting co-abundant but not necessarily  
80 co-occurring. This tool groups co-abundant genes into Metagenomic Species Pan-genomes or  
81 Metagenomic Species Pan-genomes (MSPs) and classifies genes within an MSP as core, accessory  
82 and shared. Core genes are present in all strains, accessory are present only in some (Medini et al.,  
83 2005), the shared category applies for those genes which may be present in more than one MSP  
84 due to horizontal transfer. The factors that impact directly on MSP quality include the sample  
85 composition, the sequencing depth, the previous bioinformatic steps to build the reference gene  
86 dataset and to map the reads. MSPs can be employed for taxonomic profiles of new samples from  
87 similar ecosystems at the species level, and also to compare strains between samples building a  
88 presence/absence table of accessory genes and for biomarker discovery. By binning contigs carrying  
89 genes from the same MSP it is also possible to build a MAG.

### 90 3.1.2 Binning microbial genomes with deep learning

91 The integration of deep learning techniques into the field of metagenomics has revolutionised the  
92 field of metagenomics. The Software VAMB and CoCoNet constitute two such examples in the  
93 binning area.

94 The main feature VAMB is the application of the Deep Learning technique known as Variational  
95 Auto Encoders (VAE). The variational autoencoders in VAMB learn how to integrate two data  
96 types, coabundance and kmer composition. The resulting latent representation clusters better than  
97 either of the inputs, but in principle is not limited by only two data types, and it would be possible  
98 to incorporate more data as input to the VAE. VAMB also applies a "multisplit" approach where  
99 each cluster should correspond to an organism representation across samples and each bin in a  
100 cluster to a per-sample representation of the genome of that organism. Deep learning approaches  
101 also benefit from the fact GPU technology has advanced rapidly over the past few years.

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

## 108 **3.2 Binning of viral genomes**

109 Most binning algorithms are designed for prokariotic organisms leaving viruses out of the software  
 110 scope. Viruses important for many reasons, thus it was not unexpected binning algorithms focusing  
 111 on sequences of viral origin also have shown some progress.

112 CoCoNet uses deep leaning to model co-occurrence of contigs from the same viral genome. The  
 113 network was optimized for diverse viral metagenomes, the network learns to model coverage vari-  
 114 ability within samples, a critical feature in viral metagenomes where DNA amplification methods  
 115 are needed to increase input genetic material.

116 VirBin clusters contigs for genome reconstruction of viral strains, different strains within viral  
 117 species may show different biological properties such as transmissibility or virulence. Composition  
 118 based features are usually are not enough to separate haplotypes, VirBin receives contigs as inputs  
 119 and outputs the estimated number of haplotypes via contig alignment and returns the contigs for  
 120 each haplotype based on relative abundance distribution, when the contigs are long enough VirBin  
 121 produce better results.

122 Newer strategies has been proposed and employed to reconstruct viral genomes from metage-  
 123 nomic samples, in a recent work (Natfah 2021) a new compendium of 189680 DNA viruses from  
 124 the human gut microbiome was produced. In this work they use viral informative features, among  
 125 them are presence of viral protein families (Paez-Espino 2016), and absence of non-viral fam-  
 126 ilies (El Gebali 2019), gene strand switch rate (Roux 2015) and the score produced from the  
 127 VirFinder(Ren et al 2017) software

## 128 **3.3 Binning Pipelines**

129 Other advances in binning consist in the integration of existing tools and software into bioinfor-  
 130 matic pipelines, which allow the automatic processing from beginning to end of read samples into  
 131 bins or the addition of extra processing steps to address specific biological questions or problems

132 related to the sample of origin.

133 MetaWRAP is a modular pipeline ready to perform common tasks in metagenomic analy-  
134 sis, starting from read quality checks up to bin creation, refinement, reassembly quantification,  
135 taxonomic annotation and functional annotation. MAGO pipeline integrates metagenome assem-  
136 bly, binning, bin improvement, bin quality check, bin functional annotation, and bin taxonomic  
137 annotation. SqueezeMeta also integrates external software to perform the complete analysis of  
138 metagenomic samples from sequences reading to MAG construction and annotation.nf-coreMAG  
139 supports both short and long reads, performs quality and adapter trimming, quality check, per-  
140 forms assembly, binning, checks bin quality and assigns taxonomy.

141 Autometa was developed to deal with non-model Eukariotic host contamination and complex  
142 single metagenomes, the application integrate sequence homology, nucleotide composition, cover-  
143 age and single-copy marker genes to separate microbial genomes from non model host genomes.  
144 Seqdex is a tool written in R which tries to separate endosymbionts from their host sequences,  
145 they propose the use specific features in endosymbiotic systems to better solve this problem. This  
146 tool combines partial taxonomic annotations obtained trough homology searches and sequence  
147 composition to predict the contig’s organism of origin from host and its endosymbionts and helps  
148 the user to understand how effective is the classification.

149 Among pipelines benefits we can mention they ease the reproducibility and scalability of  
150 metagenomic analysis, and they allow people with little computational experience to perform  
151 complete analysis in less time.

## 152 4 Suggestions on choosing a binning algorithm

153 A number of aspects should be considered when performing binning analysis on metagenomic  
154 samples, among them we can mention the computational resources available, the sequencing tech-  
155 nology, the number of samples and the sample’s source. Some tools employ more resources than  
156 others, and some perform better under specific circumstances, a review on the benchmarking  
157 contig-binning algorithms was done by Yue et al., 2020. If you are dealing with a large number  
158 of samples a gene-binning strategy could be taken into consideration, tools like CoMet were built  
159 around single sample binning. Long read sequence technology is gaining momentum and some  
160 tools also integrate the characteristic features generated under this technology. The environment  
161 under study also play an important role for binning, sometimes there exists host organisms whose  
162 genome sequences would be removed before starting the analysis. The environment also has a

Table 1: Comparison of binning algorithms

Software/Algorithm	Year	Description/purpose	Comment/Highlight	Doi	PubmedID
CoNet	2021	Deep learning tool for Viral Metagenome Binning	Reconstructs viral genomes	10.1093/bioinformatics/btab213	33822891
Biatic	2021	Using scaffolds to improve Metagenomic bin quality	Incorporates scaffold information	10.3389/fmicb.2021.685614	33717083
VAMB	2021	Metagenomic binning and MAG assembly	Autocoder algorithm, fast processing	10.1186/s13057-020-00777-4	33398153
phyloFlash	2020	mRNA profiling and MAG assembly	Incorporates asRNA profiling info into MAG as...	10.1093/bioinformatics/btaz441	33109753
hyBRICC-LR	2020	Metagenomic binning for Long-Reads	Suitable for Long Reads sequencing technology	10.1093/bioinformatics/btaz441	32657364
BinBam Tools	2020	Refined binning of metagenomic contigs using as...	Data preparation for targeted assembly, using s...	10.1093/bioinformatics/btaz441	32641514
MetaBin	2020	Metagenomic binning using assembly graphs	Incorporates assembly graphs	10.1093/bioinformatics/btaz441	32167328
MetaSPSim	2020	Simulating metagenomic stable isotope probing d...	Augment binning resolution with extra experimen...	10.1186/s12859-020-3372-6	32000876
MetaCon	2019	Unsupervised binning k-mers and coverage, focus...	Augment binning resolution with extra experimen...	10.1186/s12859-019-2904-4	31757198
VireBin	2019	Binning viral haplotypes from assembled contigs	Viral haplotypes MAGs	10.1186/s12859-019-3138-1	31634576
MAGO (*only tool pipeline)	2019	Framework for Proton and analysis of MAGs	Identification of endosymbiont	10.1093/bioinformatics/btaz441	31633780
SeqDox	2019	Genome separation of Endosymbionts from mixed s...	Identifies endosymbiont	10.3389/fgene.2019.00853	31608107
MetaTOR	2019	High quality MAGs from mammalian gits using met...	Incorporates 2D contact information	10.3389/fgene.2019.00753	31481973
MetaBAT 2	2019	Adaptive binning algorithm for genome recons...	Eliminates misassembly, post binning from previou...	10.7717/peerj.7359	31388474
MetaBMF	2019	Scalable binning algorithm for large-scale metab...	Employs sample X mappings of mapped read counts	10.1093/bioinformatics/btaz441	31347687
PolyCRACKER	2019	Method for partitioning polyploid bacterial genomes b...	Haplotypes for polyploid genomes	10.1186/s12864-019-5828-5	31299888
SolidBin	2019	Improving metagenomic binning with individual ext...	NaN	10.1093/bioinformatics/btaz441	30977806
Autmeta	2019	Inspecting metagenomic binning from individual ext...	Handles eukaryotic contamination	10.1093/bioinformatics/btaz441	30838416
MLBP MrGBP (Algorithm)	2019	Signal processing method for alignat free met...	Alternative description of sequences designed f...	10.1038/s41598-018-38197-9	30770850
CLAME	2018	Alignmet based algorithm allowed description of	Alignment 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 uncerta...	10.1109/EMBC.2018.8512529	30440833
LVQ-KNN	2018	Classification based RNA or DNA binning of short s...	Classify into DNA or RNA sequence	10.1016/j.virusres.2018.10.002	30291874
MSPinner	2018	Abundance based reconstruction of microbial pan...	Pan genome reconstruction	10.1093/bioinformatics/btaz441	30252023
MetaWRAP*	2018	Flexible pipeline for genome resolved metagenom...	Hybrid bin extraction algorithm	10.1186/s40168-018-0541-1	30219103
MetaVW	2018	Large scale Machine Learning Sequence classific...	Machine learning for reads based on Kmer profile	10.1007/978-1-4939-8561-6_2	30030800
BM3C	2018	Metagenomic binning through low density binning	Improvement at higher taxonomic levels, discove...	10.1093/bioinformatics/btaz441	30010790
BM3C	2018	Binning contigs using codon usage sequence comp...	Add codon usage information	10.1093/bioinformatics/btaz441	29947757
AMBER tool	2018	Assessment of Metagenome Binners	NaN	10.1093/bioinformatics/btaz441	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	Alignment of long reads against reference sequences	10.1186/s13062-018-0208-7	29678199
CoMet	2018	Binning workflow using contain coverage and com...	Single sample, include gc content and 4mer fre...	10.1186/s12859-017-1967-3	29297295
?	2017	Metagenomic binning and association of plasmids...	Plasmid binning at strain level using methylati...	10.1038/nbr.4037	29227468
MetaGen	2017	Reference-free learning with multiple metagenom...	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 dissim...	10.1186/s12859-017-1835-1	28931373
BusyBee Web	2017	Bootstrapped supervised binning and annotation	2d interactive scatterplots supervised binning	10.1093/nar/akx348	28472498
ICoVer	2017	Interactive visualisation tool for verification...	Interactive visualisation tool	10.1186/s12859-017-1653-5"	28464793
HiBin*	2017	High resolution identification of differential...	Supervised annotation, unsupervised clustering ...	10.1186/s12864-017-3686-6	28431529
BinSanity	2017	Unsupervised clustering using coverage and affi...	Reduce bias for high/low abundance	10.7717/peerj.3035	28289564
Binning-refiner	2017	Improve genome bins through the combination of ...	Combination of different binning algorithms	10.1093/bioinformatics/btaz441	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 coverage...	Adds paired end read and coalignment information	10.1093/bioinformatics/btaz441	27256312
GroupM (2)	2014	Tool for automatic recovery of population genom...	Adds differential coverage to complement compos...	10.7717/peerj.603	25289188

profound effect on the sample’s diversity, higher diversity need more sequencing depth and are harder to bin. It is also difficult for binning tools to discern between similar strains within the same sample. It is also worth mentioning that most of the available tools are not mutually exclusive and it is possible to benefit from the advantages each one has to offer and merge results depending of the aim of the study. Besides binning, other types of metagenomic analysis can be performed on microbiomes, recent reviews provide an overview of the complete process and practical guides to apply available software.

## 5 Conclusion

Until now binning methods perform poorly in samples that contain similar strains. Also do not perform great assigning 16S sequences to bins maybe due to high copy number of these sequences within a genome. Binning has been focused mainly in prokariotic organisms. Binning of organisms outside prokariotes need more development, lately some advances have been observed in viral genomes (cite viral catalogue and viral binning organisms) but the huge diversity in viral genomes still poses a challenge for current methodologies. Eukariotic microscopic organisms does not appear in the current picture.

The continuously increasing number of sequences available require more efficient/faster algorithms and new strategies to reconstruct single organisms from environmental samples. New sources of experimental information might add up into solving the binning central problem. Development of Machine learning algorithms have started in the field and we expect to see more development soon

## References

- Karypis, G., Han, E.-H., & Kumar, V. (1999). Chameleon: Hierarchical clustering using dynamic modeling. *Computer*, 32(8), 68–75.
- Plaza Oñate, F., Le Chatelier, E., Almeida, M., Cervino, A. C., Gauthier, F., Magoulès, F., ... Pichaud, M. (2019). Mspminer: abundance-based reconstitution of microbial pan-genomes from shotgun metagenomic data. *Bioinformatics*, 35(9), 1544–1552.
- Yue, Y., Huang, H., Qi, Z., Dou, H.-M., Liu, X.-Y., Han, T.-F., ... Tu, J. (2020). Evaluating metagenomics tools for genome binning with real metagenomic datasets and cami datasets. *BMC bioinformatics*, 21(1), 1–15.