

1 Metagenomic Binning Pipelines - the State of the Art

2

3 **Contents**

4	1 Abstract	1
5	2 Background	2
6	3 Overview of recent methods for metagenomic binning	3
7	3.1 Innovations in recent binning strategies	3
8	3.1.1 Metagenome Assembled Genomes	4
9	3.1.2 MSPs, binning co-abundant genes	4
10	3.1.3 Metagenomic Species Pan-genomes	4
11	3.1.4 Binning microbial genomes with deep learning	5
12	3.2 Binning for solving new biological challenges/ for viral genome	5
13	4 Choosing the most appropriate binning algorithm (Classification by output)	6
14	4.1 Identify start point variables	6
15	4.2 Identify endpoint	6
16	4.3 Tools that are complementary	6
17	5 Conclusion	6
18	References	6

19 **1 Abstract**

20 New generations of sequencing platforms coupled with numerous bioinformatics tools have led to
21 rapid technological progress in metagenomics to investigate complex microorganism communities.
22 Nevertheless, a combination of different bioinformatic tools remains necessary to draw conclusions

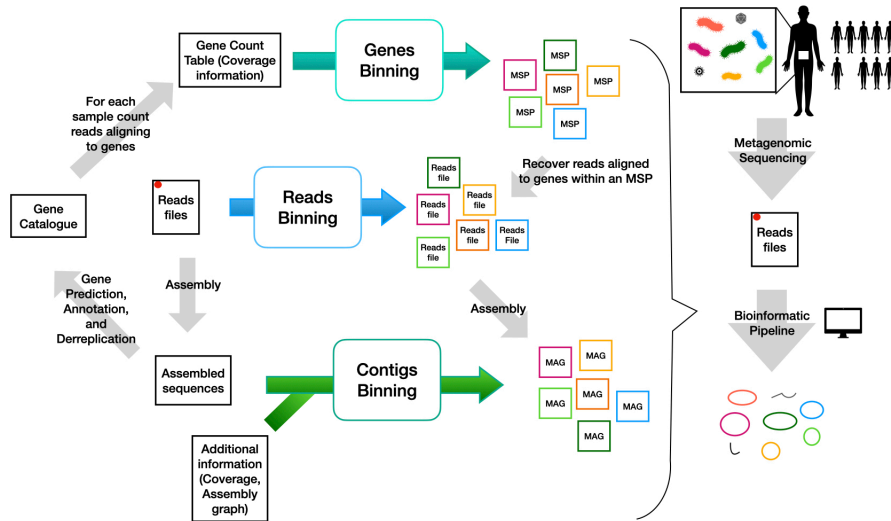


Figure 1: Summary of binning principles and techniques.

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. 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

41 data. As the cost of sequencing is rapidly falling, this burden has been significantly lessened.
42 Whole Genome Shotgun sequences does not require cultivation. At the time of writing, shotgun
43 metagenomic sequencing costs on average three times as much as 16S sequencing in comparison.
44 The objectives of this review is for the reader to be better informed about the latest algorithms
45 (since 2017) for binning metagenomic samples. The second part of this review is for the reader
46 to be informed about distinguishing factors between the methods. The last part is for the reader
47 to make an informed decision based on those factors for their needs. This review will be broken
48 down into the following sections:

49 **3 Overview of recent methods for metagenomic binning**

50 **3.1 Innovations in recent binning strategies**

51 A metagenomic sample is comprised of many organisms and the standard procedure is to retrieve
52 the sequences from the mixture of organisms. The final goal of binning is to reconstruct the
53 sequences from each organism present in the original sample. Currently we can distinguish from 3
54 different strategies in binning algorithms, read binning, contig binning, and gene binning. Among
55 the binning tools developed in recent years a subset of them are dedicated to cluster reads (read-
56 binning) (MetaBBC-LR, BioBloom Tools, CLAME, LVQ-KKN, Meta VW, HirBin, MEGAN-LR).
57 The main purpose of read-binning tools is to pre-process reads into clusters for a posterior targeted
58 assembly, here we find reference-free and non-reference-free tools, and tools designed for short-
59 read or long-read sequencing technologies. The majority of binning tools we can find are oriented
60 toward clustering contigs (contig-binning) into bins, which may represent the genome from a
61 single biological entity/organism. Contig-binning tools normally rely on coverage information
62 and sequence composition. Binning contigs have played a central role in software development
63 in the field, a review on the benchmarking binning algorithms was done by Yue et al., 2020.
64 Progress in contig-binning algorithms can be seen in the proposals to integrate new sources of
65 information (for example, from scaffold-graphs(Binnacle), paired-end reads(COCACOLA), or 3D
66 contact information(MetaTOR)) and state of the art algorithms in machine learning (CoCoNet,
67 Variational Autoencoders for Metagenomic Binning (VAMB)).

68 **3.1.1 Metagenome Assembled Genomes**

69 A Metagenome-Assembled Genome (MAG) is a single-taxon assembly based on one or more binned
70 metagenomes that has been asserted to be a close representation to an actual individual genome
71 (that could match an already existing isolate or represent a novel isolate).

72 **3.1.2 MSPs, binning co-abundant genes**

73 Binning of co-abundant genes represents an alternative proposal to reconstruct species/biological
74 entities from a set of metagenomic samples. Co-abundant gene binning methods assume each gene
75 coming from a shared chromosome will display proportional abundances across samples, if you have
76 enough samples from a common environment you can identify the sets of genes from a common
77 organism of origin (MLGs Chameleon-clust 2012, CAGs and MGSs Canopy 2014, Markovclust-
78 MGCs Karlsson 2013, MSPs MSPminner 2018). To the extent of our knowledge, in the past few
79 years MSPminer was the only Software developed exploiting this approach. MSPminer introduced
80 a robust proportionality measure detecting co abundant but no necessarily co-occurring. This tools
81 groups co-abundant genes into Metagenomic Species Pan-genomes or Metagenomic Species Pan-
82 genomess (MSPs) and classify genes within an MSP as core, accessory and shared. The factors
83 that impact directly on MSP quality include the sample composition, the sequencing depth, the
84 previous bioinformatic steps to build the reference gene dataset and to map the reads. A high
85 number of samples with varying phenotypes improve the quality of MSPs. MSPs can be employed
86 for taxonomic profiles of new samples from similar ecosystems, to compare strains between samples
87 building a presence/absence table of accessory genes and for biomarker discovery. By binning
88 contigs carrying genes from the same MSP it is also possible to build a MAG.

89 **3.1.3 Metagenomic Species Pan-genomes**

90 Microbial pan-genomes are gene repertoires composed of core genes present in all strains and
91 accessory genes present in only some of them (Medini et al., 2005). In a shotgun metagenomic
92 sequencing context, we define as shared the genes detected in some samples where the species is
93 not present. A strain found in a sample is an instance of the species pan-genome: it is made of all
94 the species (shared) core genes and of a subset of (shared) accessory genes. Core genes are suitable
95 for taxonomic profiling at species-level while accessory genes can be used to compare strains across
96 samples. Genes tagged as shared should be used carefully as they contain false positives counts
97 or are subject to horizontal transfer. Core genes are suitable for taxonomic profiling at species-

level while accessory genes can be used to compare strains across samples. Genes tagged as shared should be used carefully as they contain false positives counts or are subject to horizontal transfer.

3.1.4 Binning microbial genomes with deep learning

The integration of deep learning techniques into the field of metagenomics has revolutionised the field of metagenomics. The VAMB pipeline was developed to take advantage of variational autoencoders; a generative machine learning model that uses a deep variational autoencoders (Nissen et al., n.d.)... COCONET (Arisdakessian, Nigro, Steward, Poisson, & Belcaid, 2021)...

3.2 Binning for solving new biological challenges/ for viral genome

Most binning algorithms are designed for prokaryotic organisms leaving viruses out of the design. Viruses are a serious threat to human health CoCoNet uses deep learning to model co-occurrence of contigs from the same viral genome. The method uses a neural network which returns the probability for a pair of contigs coming from the same genome, these probabilities are employed to construct bins representing the species present in the sample. 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 produces better results. Pipelines for Endosymbiont organisms binning also have been recently developed. Extracting endosymbiont sequences from their host poses a similar problem as a metagenomic sample. Seqdex is a tool written in R which tries to separate endosymbiont from host sequences, they proposed they could use specific features in endosymbiotic systems to better solve this problem. This tool combines partial taxonomic annotations obtained through homology searches and sequence composition to predict the contig's organism of origin from host and its endosymbionts and helps the user to understand how effective is the classification.

125 4 Choosing the most appropriate binning algorithm (Clas- 126 sification by output)

127 A review on the benchmarking binning algorithms was done by Yue et al., 2020. Resource man-
128 agement is an important factor in the choice of binning algorithm. The tradeoff between number
129 of Central Processing Units (CPUs), memory, and time are important considerations. Newer ad-
130 vances in pipeline technologies have ameliorated these costs. An analysis pipeline is defined as
131 a program that combines several programs in a defined order to complete a complex analysis.
132 Improperly developed, validated, and/or monitored pipelines may generate inaccurate results.

133 4.1 Identify start point variables

134 4.2 Identify endpoint

135 4.3 Tools that are complementary

136 5 Conclusion

137 Until now binning methods perform poorly in samples that contain similar strains. Also do not
138 perform great assigning 16S sequences to bins maybe due to high copy number within a genome.
139 Binning has been focused mainly in prokaryotic organisms. Binning of organisms outside proka-
140 riyotes need more development, lately some advances have been observed in viral genomes (cite
141 viral catalogue and viral binning organisms) but the huge diversity in viral genomes still poses a
142 challenge for current methodologies. eukaryotic microscopic organisms does not appear in the
143 current picture. The continuously increasing number of sequences available require more effi-
144 cient/faster algorithms and new strategies to reconstruct single organisms from environmental sam-
145 ples. Development of Machine learning algorithms have started in the field and we expect to see
146 more development soon New and open areas of research in which the application of metagenomic
147 pipelines are relevant The increased impact of machine learning in analysis Short section - just for
148 past-present-future completeness Future developments for metagenomic analysis

149 References

150 Arisdakessian, C. G., Nigro, O. D., Steward, G. F., Poisson, G., & Belcaid, M. (2021). Coconet:

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
hyBRCC-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 endosymbionts	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 patterns from previou...	10.7717/peerj.7359	31388474
MetaBMF	2019	Scalable binning algorithm for large-scale meta...	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 extraction of metagenomic bins	NaN	10.1093/bioinformatics/btaz441	30977806
Autmeta	2019	Signal processing method for aligning free met...	Handles eukaryotic contamination	10.1093/bioinformatics/btaz441	30838416
MLBP MrGBP (Algorithm)	2019	Aligning metagenomic reads from individual	Alternative description of sequences designed f...	10.1038/s41586-018-38197-9	30770850
CLAME	2018	Signal processing method for aligning free met...	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/44x348	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	28289664
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

151 an efficient deep learning tool for viral metagenome binning. *Bioinformatics*.
152 Nissen, J. N., Johansen, J., Allesøe, R. L., Sønderby, C. K., Armenteros, J. J. A., Grønbech,
153 C. H., . . . others (n.d.). Improved metagenome binning and assembly using deep variational
154 autoencoders. *Nature Biotechnology*, 1–6.
155 Yue, Y., Huang, H., Qi, Z., Dou, H.-M., Liu, X.-Y., Han, T.-F., . . . Tu, J. (2020). Evaluating
156 metagenomics tools for genome binning with real metagenomic datasets and caml datasets.
157 *BMC bioinformatics*, 21(1), 1–15.