

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 Binning co-abundant genes	4
9	3.1.2 Binning microbial genomes with deep learning	4
10	3.2 Binning of viral genomes	4
11	3.3 Binning Pipelines	5
12	4 Choosing a binning algorithm (Classification by output)	6
13	4.1 Identify start point variables	6
14	4.2 Identify endpoint	6
15	5 Conclusion	6
16	References	8

17 **1 Abstract**

18 New generations of sequencing platforms coupled with numerous bioinformatics tools have led to
19 rapid technological progress in metagenomics to investigate complex microorganism communities.
20 Nevertheless, a combination of different bioinformatic tools remains necessary to draw conclusions
21 out of microbiota studies. As sequencing costs have dropped at a rate above 'Moore's law', a
22 greater number of large data sets are being produced than ever before. Newer algorithms that

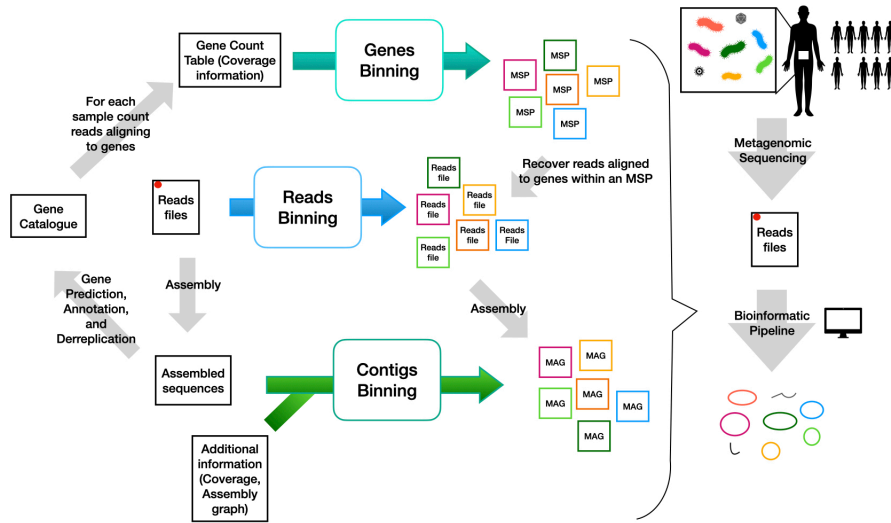


Figure 1: Summary of binning principles and techniques.

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 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 metage-

41 nomic sequencing costs on average three times as much as 16S sequencing in comparison. Here
42 we will briefly recapitulate recent binning algorithms and highlight some of the developments in
43 the field, among them, the use of new algorithms and strategies employed to achieve the goal of
44 identifying the organisms composing microbiome communities. We hope this overview could aid
45 the reader to choose a binning algorithm or a combination of them based on their specific needs.

46 **3 Overview of recent methods for metagenomic binning**

47 **3.1 Innovations in recent binning strategies**

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

55 Current contig-binning tools normally are reference free (i.e do not depend on reference se-
56 quences to perform clustering) and rely on coverage information and sequence composition. Progress
57 in contig-binning algorithms can be seen in the proposals to integrate new sources of information
58 (for example, from scaffold-graphs(Binnacle), paired-end reads(COCACOLA), or 3D contact in-
59 formation(MetaTOR)) and state of the art algorithms in machine learning (CoCoNet, Variational
60 Autoencoders for Metagenomic Binning (VAMB)). We also notice the development of Bin refine-
61 ment tools (DAS-tool, Binning Refiner), this tools rely on the outputs from multiple contig-binning
62 algorithms and attempt to combine them to produce better results.

63 Binning of contigs have played a central role in software development in the field, a review on
64 the benchmarking binning algorithms was done by Yue et al., 2020. Beside contig-binning tools
65 we can also distinguish read-binning tools and co-abundant-gene-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
74 gene coming from a shared chromosome will display proportional abundances across samples, if
75 you have enough samples from a similar environment you can identify the sets of genes from
76 a common organism of origin (MLGs Chameleon-clust 2012, CAGs and MGSs Canopy 2014,
77 Markovclust-MGCs Karlsson 2013, MSPs MSPminner 2018). In the past few years the MSPminer
78 software was developed exploiting this approach. MSPminer introduced a robust proportionality
79 measure detecting co-abundant but not necessarily co-occurring. This tool groups co-abundant
80 genes into Metagenomic Species Pan-genomes or Metagenomic Species Pan-genomess (MSPs) and
81 classify genes within an MSP as core, accessory and shared. Core genes are present in all strains,
82 accessory are present only in some (Medini et al., 2005), the shared category applies for those
83 genes which may be present in more than one MSP due to horizontal transfer. The factors
84 that impact directly on MSP quality include the sample composition, the sequencing depth, the
85 previous bioinformatic steps to build the reference gene dataset and to map the reads. MSPs can
86 be employed for taxonomic profiles of new samples from similar ecosystems at the species level,
87 and also to compare strains between samples building a presence/absence table of accessory genes
88 and for biomarker discovery. By binning contigs carrying genes from the same MSP it is also
89 possible to build a MAG.

90 **3.1.2 Binning microbial genomes with deep learning**

91 An increasingly attractive approach to the field of metagenomic binning is the utilization of
92 deep learning. Components of ML that have been employed for binning include K-NN, The
93 VAMB pipeline was developed to take advantage of variational autoencoders; a generative ma-
94 chine learning model that uses a deep variational autoencoders (Nissen et al., n.d.)... COCONET
95 (Arisdakessian, Nigro, Steward, Poisson, & Belcaid, 2021)...

96 **3.2 Binning of viral genomes**

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

100 CoCoNet uses deep learning to model co-occurrence of contigs from the same viral genome. The

101 method uses a neural network which returns the probability for a pair of contigs coming from the
102 same genome, this probabilities are employed to construct bins representing the species present
103 in the sample. The network was optimized for diverse viral metagenomes, the network learns to
104 model coverage variability within samples, a critical feature in viral metagenomes where DNA
105 amplification methods are needed to increase input genetic material.

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

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

118 **3.3 Binning Pipelines**

119 Other advances in binning consist in the integration of existing tools and software into bioinfor-
120 matic pipelines, which allow the automatic processing from beginning to end of read samples into
121 bins or the addition of extra processing steps to address specific biological questions or problems
122 related to the sample of origin.

123 MetaWRAP is a modular pipeline ready to perform common tasks in metagenomic analy-
124 sis, starting from read quality checks up to bin creation, refinement, reassembly quantification,
125 taxonomic annotation and functional annotation. MAGO pipeline integrates metagenome assem-
126 bly, binning, bin improvement, bin quality check, bin functional annotation, and bin taxonomic
127 annotation. SqueezeMeta also integrates external software to perform the complete analysis of
128 metagenomic samples from sequences reading to MAG construction and annotation.nf-coreMAG
129 supports both short and long reads, performs quality and adapter trimming, quality check, per-
130 forms assembly, binning, checks bin quality and assigns taxonomy. Autometa was developed to
131 deal with non-model Eukariotic host contamination and complex single metagenomes, the applica-

tion integrate sequence homology, nucleotide composition, coverage and single-copy marker genes to separate microbial genomes from non model host genomes. Seqdex is a tool written in R which tries to separate endosymbionts from their host sequences, they propose the 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.

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

4 Choosing a binning algorithm (Classification by output)

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

4.1 Identify start point variables

4.2 Identify endpoint

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

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	Identifies 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 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 binning	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/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

160 samples. New sources of experimental information might add up into solving the binning central
161 problem.

162 Development of Machine learning algorithms have started in the field and we expect to see
163 more development soon

164 Short section - just for past-present-future completeness

165 **References**

166 Arisdakessian, C. G., Nigro, O. D., Steward, G. F., Poisson, G., & Belcaid, M. (2021). Coconet:
167 an efficient deep learning tool for viral metagenome binning. *Bioinformatics*.

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

171 Yue, Y., Huang, H., Qi, Z., Dou, H.-M., Liu, X.-Y., Han, T.-F., . . . Tu, J. (2020). Evaluating
172 metagenomics tools for genome binning with real metagenomic datasets and caml datasets.
173 *BMC bioinformatics*, 21(1), 1–15.