

Metagenomic Binning Pipelines - the State of the Art

Abstract

New generations of sequencing platforms coupled with numerous bioinformatics tools have led to rapid technological progress in metagenomics and metatranscriptomics 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. 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.

Index

- *Abstract*
- *Background/Introduction*
 - *Binning problem definition (recover biological entities from metagenomic sequencing)*
 - *problem relevance (Explosion in metagenomics, reduction in sequencing cost, increased computer capacity)*
 - *Review objectives (Brief summary on popular tools, innovations overview of recent tools)*
- *Popular/Previous Binning software*
 - *Proposed solutions (bin contigs into bins(MAG if good quality) based on their kmer composition and abundance/coabundance)*

- 24 – *Tools available (Cite recent benchmark)*
- 25 • *Overview of recent metagenomic binning tools*
- 26 – *Inovations in binning tools*
- 27 * *Inovations in proposed solutions/ strategy innovations* Read binning gene-abundance
- 28 binning (CAG, MGS, MSPi) Integrate new experimental data
- 29 * *Software/algorithms innovations* machine-learning/deep-learning implementation
- 30 – *Inovations in specific biological questions* Viral genomes and viral strains; Endosym-
- 31 bionts
- 32 • *Choosing a binning algorithm*
- 33 – *Identifiy start point variables*
- 34 * *Sample origin (Host contamination, diversity)*
- 35 * *Number of samples (some tools require many samples to perform well)*
- 36 * *Sequencing technology (Most tools employ illumina, LongReads are increasing)*
- 37 * *Computational resources available*
- 38 – *Identify endpoint*
- 39 * *organism of interest viral(ref viral catalogue), bacteria, all*
- 40 – *Tools are complementary MSP/Metabat*
- 41 • *Conclusions*
- 42 – *Current limitations and future directions* Do not perform well on multiple strains, on
- 43 the same sample

44 Figure. Binning software historical citations barplot Figure. Decision tree, overview of

45 metagenomic binning Table. List of binning software since 2017

46 Background

47 The explosion in popularity and success in the field of metagenomics over the last 25 years can

48 be largely attributed to the advances in computing tecnologies. An example of the outcomes of

49 this can be found in the Human Microbiome Project; a project that has been greatly imprved

50 the understanding of the microbila flora involved in human health and disease. These advances

51 have brought with them greater demands for storage, CPU time, and consequently more efficient
 52 algorithms. The main function of binning tools is to reconstruct species/biological entities from
 53 metagenomic samples. Compared to amplicon, shotgun metagenome can provide functional gene
 54 profiles directly and reach a much higher resolution of taxonomic annotation. However, due to the
 55 high demands on computational resources, cost, and expertise necessary to perform this analysis,
 56 researchers have historically been limited in their capacity to collect and analyse sequencing data.
 57 As the cost of sequencing is rapidly falling, this burden has been significantly lessened At the time
 58 of writing, shotgun metagenomic sequeng costs on average three times as much as 16S sequencing
 59 comparitavely. The objectives of this review is for the reader to be better informed about the
 60 latest algorithms (since 2017) for binning metagenomic samples. The second part of this review
 61 is for the reader to be informed about distinguishing factors between the methods. The last part
 62 is for the reader to make an informed decision based on those factors for their needs. This review
 63 will be broken down into the following sections:

- 64 • *List of the binning algorithms*
- 65 • *Classify binning algorithms based on their objectives, guideline for algorithm choice, subsec-*
 66 *tion msp mag*
- 67 • *Current limitations and Future directions*

68 Recent methods for metagenomic binning

69 A metagenomic sample is comformed of many organisms and the standard procedure is to retrieve
 70 the sequences from the mixture of organisms. The final goal of binning is to reconstruct the se-
 71 quences from each organism present in the original sample. Among the binning tools developed
 72 in recent years we can distinguish a subset dedicated to cluster reads (read-binning) (MetaBBC-
 73 LR, BioBloom Tools, CLAME, LVQ-KKN, Meta VW, HirBin, MEGAN-LR). The main purpose
 74 of read-binning tools is to preprocess reads into clusters for a posterior targeted assembly, here
 75 we find reference-free and non-reference-free tools, and tools designed for short-read or long-read
 76 sequencing technologies. The majority of binning tools we can find are oriented toward cluster-
 77 ing contigs (contig-binning) into bins, which may represent the genome from a single biological
 78 entity/organism. Contig-binning tools normally rely on coverage information and sequence compo-
 79 sition. Progress in contig-binning algorithms can be seen in the proposals to integrate new sources
 80 of information (for example, from scaffold-graphs(Binnacle), paired-end reads(COCACOLA), or

81 3D contact information(MetaTOR)) and state of the art algorithms in machine learning (CoCoNet,
82 Variational Autoencoders for Metagenomic Binning (VAMB)).

83 **Metagenome Assembled Genomes**

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

87 **Binning microbial genomes with deep learning**

88 The integration of deep learning techniques into the field of metagenomics has revolutionised the
89 field of metagenomics. The VAMB pipeline was developed to take advantage of variational au-
90 toencoders; a generative machine learning model that uses a combination. Improved metagenome
91 binning and assembly using deep variational autoencoders. Nature biotechnology the VAMB
92 pipeline (?, ?).

93 **Binning for viral genomes**

94 New insights from uncultivated genomes of the global human gut microbiome (?, ?).

95 **Chosing the most appropriate binning algorithm (Classifica-** 96 **tion by output)**

97 A review on the benchmarking binning algorithms was done by ?, ?. Resource management is
98 an important factor in the choice of binning algorithm. The tradeoff between number of Central
99 Processing Units (CPUs), memory, and time are important considerations. Newer advances in
100 pipeline technologies have ameliorated these costs. Alignment based or alignment free. An analysis
101 pipeline is defined as a program that combines several programs in a defined order to complete
102 a complex analysis. Improperly developed, validated, and/or monitored pipelines may generate
103 inaccurate results.

104 **MSPs, binning co-abundant genes**

105 Binning of co-abundant genes represents an alternative proposal to reconstruct species/biological
106 entities from a set of metagenomic samples. Co-abundant gene binning methods assume each

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
BiBin	2020	Refined binning of metagenomic contigs using as...	Data preparation for targeted assembly, using s...	10.1093/bioinformatics/btaz17	32641514
MetaBin Tools	2020	Metagenomic binning of contigs using as...	Incorporates assembly graphs	10.1093/bioinformatics/btaz180	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/mbe/mbz237	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 guts 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 meta...	Employs sample X mappings of mapped read counts	10.1093/bioinformatics/btaz577	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/btaz253	30977806
Autmeta	2019	Signal processing method for aligning free met...	Handles eukaryotic contamination	10.1093/bat/bkz148	30838416
MLBP MrGBP (Algorithm)	2019	Alternative description of sequences designed f...	Alternative description of sequences designed f...	10.1038/s41598-018-38197-9	30770850
CLAME	2018	Aligner based algorithm 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	30447633
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/btaz830	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/btaz611	30010790
BM3C	2018	Binning contigs using codon usage sequence comp...	Add codon usage information	10.1093/bioinformatics/btaz611	29947757
Assessment of Metagenome Binners	2018	Assessment of Metagenome Binners	NaN	10.1093/gigascience/gix069	29893851
AMBER tool	2018	Derreplication aggregation and scoring strategy	Combines several binning algorithm results	10.1038/s41564-018-0171-1	29807988
DAS Tool	2018	Long Read/ contigs taxonomic binning	Alignment of long reads against reference sequences	10.1186/s13062-018-0208-7	29678199
MEGAN-LR	2018	Binning workflow using contain coverage and com...	Single sample, include gc content and 4mer fre...	10.1186/s12859-017-1967-3	29297295
CoMet	2017	Metagenomic binning and association of plasmids...	Plasmid binning at strain level using methylati...	10.1038/nbt.4037	29227468
?	2017	Reference-free learning with multiple metagenom...	Requires multiple samples	10.1186/s13059-017-1323-y	28974263
MetaGen	2017	Improved formula for calculate oligonucleotide...	Math formula to calculate oligo sequence dissim...	10.1186/s12859-017-1835-1	28931373
d2sBin add onn	2017	Bootstrapped supervised binning and annotation	2d interactive scatterplots supervised binning	10.1093/nar/bkx348	28472498
BuscBee Web	2017	Interactive visualisation tool for verification...	Interactive visualisation tool	10.1186/s12859-017-1653-5"	28464793
ICoVer	2017	High resolution identification of differential...	Supervised annotation, unsupervised clustering ...	10.1186/s12864-017-3686-6	28431529
HiBin*	2017	Unsupervised clustering using coverage and affi...	Reduce bias for high/low abundance	10.7717/peerj.3035	28289564
BinSanity	2017	Improve genome bins through the combination of ...	Combination of different binning algorithms	10.1093/bioinformatics/btaz086	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/btaz290	27256312
GroupM (2)	2014	Tool for automatic recovery of population genom...	Adds differential coverage to complement compos...	10.7717/peerj.603	25289188

gene coming from a shared chromosome will display proportional abundances across samples, if you have enough samples from a common 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). `##### HEAD`

To the extent of our knowledge, in the past few years MSPminer is the only available Software exploiting this approach. MSPminer introduced a robust proportionality measure detecting co-abundant but not necessarily co-occurring. This tool groups co-abundant genes into Metagenomic Species Pan-genomes or MSPs and classifies genes within an MSP as core, accessory and shared.

The factors that impact directly on MSPs quality include the sample composition, the sequencing depth, the previous bioinformatic steps to build the reference gene dataset and to map the reads. A high number of samples with varying phenotypes improve the quality of MSPs.

MSPs can be employed for taxonomic profiles of new samples from similar ecosystems, to compare strains between samples 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.

Co-abundant gene binning methods perform better in large sample datasets

0.1 Metagenomic Species Pan-genomes

Microbial pan-genomes are gene repertoires composed of core genes present in all strains and accessory genes present in only some of them (Medini et al., 2005). In a shotgun metagenomic sequencing context, we define as shared the genes detected in some samples where the species is not present.

A strain found in a sample is an instance of the species pan-genome: it is made of all the species (shared) core genes and of a subset of (shared) accessory genes. 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 positive counts or are subject to horizontal transfer.

1 Conclusion

- *New and open areas of research in which the application of metagenomic pipelines are relevant*

- 136 • *HMP and other*
- 137 • *The increased impact of machine learning in analysis*
- 138 • *Short section - just for past-present-future completeness*
- 139 • *Future developments for metagenomic analysis*

140 **1.1 Weaknesses and future developments**

141 Until now binning methods perform poorly in samples containing similar strains.

142 ===== To the extent of our knowledge, in the past few years MSPminer is the only avail-
 143 able Software exploiting this approach. MSPminer introduced a robust proportionality measure
 144 detecting co abundant but no necessarily co occurring. This tools groups co-abundant genes into
 145 Metagenomic Species Pan-genomes or Metagenomic Species Pan-genomess (MSPs) and classify
 146 genes within an MSP as core, accesory and shared. The factors that impact directly on MSP
 147 quality include the sample composition, the sequencing depth, the previos bioinforamtic steps to
 148 build the reference gene dataset and to map the reads. A high number of samples with varying
 149 phenotypes improve the quality of MSPs. MSPs can be employed for taxonomic profiles of new
 150 samples from similar ecosystems, to compare strains between samples building a presence/absence
 151 table of accesory genes and for biomarker discovery. By binning contigs carrying genes from the
 152 same MSP it is also possible to build a MAG.

153 **Metagenomic Species Pan-genomes**

154 Microbial pan-genomes are gene repertoires composed of core genes present in all strains and
 155 accessory genes present in only some of them (?, ?). In a shotgun metagenomic sequencing context,
 156 we define as shared the genes detected in some samples where the species is not present. A strain
 157 found in a sample is an instance of the species pan-genome: it is made of all the species (shared)
 158 core genes and of a subset of (shared) accesory genes. Core genes are suitable for taxonomic
 159 profiling at species-level while accesory genes can be used to compare strains across samples.
 160 Genes tagged as shared should be used carefully as they contain false positives counts or are
 161 subject to horizontal transfer.

162 Conclusion

163 New and open areas of research in which the application of metagenomic pipelines are relevant
164 The increased impact of machine learning in analysis Short section - just for past-present-future
165 completeness Future developments for metagenomic analysis
166 `0b81c509cb38a218555a68fc24a137c2e6560df2`

167 References

- 168 Nayfach, S., Shi, Z. J., Seshadri, R., Pollard, K. S., & Kyrpides, N. C. (2019). New insights from
169 uncultivated genomes of the global human gut microbiome. *Nature*, 568(7753), 505–510.
- 170 Nissen, J. N., Johansen, J., Allesøe, R. L., Sønderby, C. K., Armenteros, J. J. A., Grønbech,
171 C. H., ... others (n.d.). Improved metagenome binning and assembly using deep variational
172 autoencoders. *Nature Biotechnology*, 1–6.
- 173 Tettelin, H., Massignani, V., Cieslewicz, M. J., Donati, C., Medini, D., Ward, N. L., ... others
174 (2005). Genome analysis of multiple pathogenic isolates of streptococcus agalactiae: impli-
175 cations for the microbial “pan-genome”. *Proceedings of the National Academy of Sciences*,
176 102(39), 13950–13955.
- 177 Yue, Y., Huang, H., Qi, Z., Dou, H.-M., Liu, X.-Y., Han, T.-F., ... Tu, J. (2020). Evaluating
178 metagenomics tools for genome binning with real metagenomic datasets and caml datasets.
179 *BMC bioinformatics*, 21(1), 1–15.