

# Metagenome-Assembled Genome Binning Methods Disproportionately Fail for Plasmids and Genomic Islands

This manuscript ([permalink](#)) was automatically generated from [fmaguire/mag\\_sim\\_paper@79a8542](#) on November 5, 2019.

## Authors

---

Please note the current author order is chronological and does not reflect the final order.

- **Finlay Maguire\***

 [0000-0002-1203-9514](#) ·  [fmaguire](#) ·  [fmaguire](#)

Faculty of Computer Science, Dalhousie University · Funded by ['Genome Canada', 'Donald Hill Family Fellowship']

- **Baofeng Jia\***

 [XXXX-XXXX-XXXX-XXXX](#) ·  [imasianxd](#)

Department of Biochemistry and Molecular Biology, Simon Fraser University

- **Kristen Gray**

 [XXXX-XXXX-XXXX-XXXX](#)

Department of Biochemistry and Molecular Biology, Simon Fraser University

- **Venus Lau**

 [XXXX-XXXX-XXXX-XXXX](#)

Department of Biochemistry and Molecular Biology, Simon Fraser University

- **Robert G. Beiko**

Faculty of Computer Science, Dalhousie University

- **Fiona S.L. Brinkman**

 [XXXX-XXXX-XXXX-XXXX](#)

Department of Biochemistry and Molecular Biology, Simon Fraser University

## Abstract

---

### final numbers to add

Metagenomic methods, in which all the DNA in sample is simultaneously sequenced, is an increasingly popular method in the life sciences. They have a major advantage over genomic or phenotypic methods as they do not require time-intensive and bias-inducing culturing steps. This means a much greater diversity can be profiled with minimal *a priori* assumptions. Due to this strength, metagenomics is emerging as a key tool in public health microbiology for surveillance of virulence and antimicrobial resistance (AMR) genes. The most important sequences for surveillance purposes are those associated with mobile genetic elements such as plasmids and genomic islands (GIs). Unfortunately, metagenomic data, even when assembled, results in complex mixed set of DNA fragments rather than nicely resolved individual genomes. Recently, methods have been developed that attempt to group these fragments into bins likely to have been derived from the same underlying genome. These bins are commonly known as metagenome-assembled genomes (MAGs). MAG based approaches have been used to great effect in revealing huge amounts of previously uncharacterised microbial diversity. These methods perform this grouping using aspects of the sequence composition and the relative abundance of that sequence in the dataset. Unfortunately, plasmids are often represented at different copy numbers than the corresponding chromosomes. Additionally, both plasmids and genomic islands often feature significantly different sequence composition than the rest of the source genome as a whole. Due to this we hypothesise these types of sequences will be highly under-represented in MAG based approaches.

To evaluate this we generated a simulated metagenomic dataset comprised of genomes with large numbers of plasmids and considerable proportion of chromosomal DNA consisting of GIs at varying relative abundances. MAGs were then recovered from this data using a variety of different established MAG pipelines and parameterisations and correct binning of plasmid and GI sequences calculated relative to the genomes as a whole. We show that regardless of the MAG approach used, plasmid and GI dominated sequences will systematically be left unbinned or incorrectly binned. This indicates the importance of read based approaches for thorough evaluation of resistome complements in metagenomic data.

## Introduction

---

Metagenomics, the untargeted sequencing of all DNA within a sample, has become the dominant approach for characterising viral and microbial communities over the last 17 years [1,2]. By targeting all genomic contents, these methods allow researchers to profile the functional potential and the taxonomic composition of a sample. This is in contrast to barcoding based approaches such as 16S or 18S rRNA sequencing which only provide taxonomic information [3] (although you can attempt to predict functional potential from taxonomic data [4,5]). One of many areas where metagenomics has been very useful is in the analysis of antimicrobial resistance (AMR). Using these approaches has been instrumental in developing our understanding of the distribution and evolutionary history of AMR genes [6,7,8]. It has also formed a very useful tool for pathogen tracking in public health outbreak analyses [9]

While 3rd generation long-read technology has begun to be adopted in metagenomics analyses [10,11] the majority of analyses still involve high-throughput 2nd generation sequencing. These 2nd generation platforms such as Illumina's MiSeq provide high numbers (10s-100s of millions) of relatively short reads (150-250bp) randomly sampled from the underlying DNA in the sample. This sampling is therefore in proportion to the relative abundance of different organisms (i.e. more abundant organisms will be more represented in the reads). There are 2 main approaches to the analysis of 2nd generation metagenomic data: read homology and metagenome assembly. Read-based approaches involve using reference databases and BLAST based sequence similarity search tools (e.g. DIAMOND [12]), read mapping (e.g. Bowtie 2 [13]), Hidden Markov Models (e.g. HMMER3 [14]) or k-mer hashing (e.g. CLARK [15]). These read-based approaches allow analysis of all reads with detectable similarity to the genes you are interested even if the organism is relatively under-represented in the dataset. However, read-based methods are reliant on quality of the reference database (i.e. you don't detect things you don't already know about) and does not provide any information about the genomic organisation of the genes.

In order to get more data about the relative genomic context and organisation of your genes of interest it is possible (although computationally demanding) to assemble the short reads into longer fragments of DNA (contigs). This approach has been used successfully in even very early metagenomic analysis papers [16]. There are a variety of specialised de Bruijn graph assemblers developed to handle the particular challenges of this type of assembly (such as metaSPAdes [17], IDBA-UD [18], and megahit [19]) each with a range of different strengths and weaknesses [20]. While metagenomic assembly does provide longer stretches of DNA incorporating information about multiple genes without further analysis it still leaves you with a large collection of DNA fragments with no obvious groupings.

An increasingly common way to deal with this is to attempt to group these assembled contigs into bins all derived from the same underlying genome in the sample. These resulting bins are known as metagenome assembled genomes (MAGs). This binning is typically performed by grouping all the contigs with similar abundance and similar sequence composition into the same bin. A range of tools have been released to perform this binning including CONCOCT [21], MetaBAT 2 [22], and MaxBin 2 [23]. There is also the metabinning tool DAS Tool [24] which combines predictions from multiple binning tools together. These MAG approaches have been used to great effect in unveiling huge amounts of previously uncharacterised genomic diversity [25,26,27].

Unfortunately, only a relatively small proportion of reads are successfully assembled and binned in large complex metagenome datasets e.g. 24.2-36.4% of reads from permafrost [28] and soil metagenomes [29]. Additionally, a large number of detected genomes are not reconstructed at all with only ~23% of all detected genomes recovered in the soil metagenome [29]. There have been attempts to benchmark and compare these tools such as the Critical Assessment of Metagenome Interpretation (CAMI) challenge's (<https://data.cami-challenge.org/>) Assessment of Metagenome BinnERs (AMBER) [30] however these only investigate the completeness and purity of recovered MAGs relative to true genomes in the sample. They don't attempt to assess whether there are specific components of the underlying genomes that are disproportionately lost. Two such genomic elements of great importance to the study and lateral gene transfer of AMR are genomic islands and plasmid sequences.

Genomic islands (GIs) are clusters of genes known or predicted to have been acquired through lateral gene transfer (LGT) events. These include integrons, transposons, integrative and conjugative elements (ICEs) and prophages (integrated phages) [31,32]. They have been shown to disproportionately encode virulence factors [33] and are a major mechanism of LGT of AMR genes [34,35]. However, these GIs often have different nucleotide composition compared to the rest of the genome [31], which is exploited by tools such as SIGI-HMM [36] and IslandPath-DIMOB [37] to detect GIs. Additionally, GIs may exist as multiple copies within a genome [38] leading to potential difficulties in correctly assembling these regions in metagenome assemblies as well as likely biases in the calculation of coverage statistics. Similarly, plasmids are a major source of the dissemination and translocation of AMR genes throughout microbial ecosystems [34,39]. They also exist at variable copy number [40,41] and with markedly different sequence composition to the genome they are associated with [42,43].

As MAG binning is performed on the basis of sequence composition and coverage this suggests that these sequences are liable to being incorrectly binned or lost in the process of recovering MAGs. Due to the importance of these genomic components in the function and spread of pathogenic traits such as AMR and virulence it is vital that we assess the impact of assembly and binning on the representation of these elements. This is particularly important with the increasing popularity of MAG approaches within microbial and public health research. Therefore, to address this issue we performed an analysis of GI and plasmid recovery accuracy across a range of assembly and binning approaches using a simulated medium complexity metagenome comprised of GI and plasmid rich taxa.

## Materials and Methods

---

All analyses presented in this paper can be reproduced and inspected using the Jupyter (version XX) [citation] notebook (.ipynb) within the associated github repository ([https://github.com/fmaguire/MAG\\_gi\\_plasmid\\_analysis](https://github.com/fmaguire/MAG_gi_plasmid_analysis)). The specific code version used for this paper is also archived under DOI:XXXXXX.

### Metagenome Simulation

All genomes were selected from the set of completed RefSeq genomes as of April 2019. Genomic islands for these genomes were previously predicted using IslandPath-DIMOB [37] and collated into the IslandViewer database (<http://www.pathogenomics.sfu.ca/islandviewer>) [44]. Plasmid sequences and numbers were recovered for each genome using the linked genbank Project IDs.

30 genomes were arbitrarily chosen to exemplify the following criteria: - 10 with high numbers of plasmids - 10 with a very high proportion (>10%) of chromosomes corresponding to composition detected GIs. - 10 with a very low proportion (<1%) of chromosomes corresponding to composition detected GIs.

The data used to select the taxa is listed in Supplemental Table 1 and the details of the selected subset taxa are listed in Supplemental Table 2 with their NCBI accessions. The sequences themselves are in `data/sequences/sequences.tar.bz2`

In accordance to the recommendation in the CAMI challenge [45] the genomes were randomly assigned a relative abundance following a log-normal distribution ( $\mu = 1$ ,  $\sigma = 2$ ). Plasmid copy number estimates could not be accurately found for all organisms, therefore, plasmids were randomly assigned a copy number regime: low (1-20), medium (20-100), or high (500-1000) at a 2:1:1 rate. Within each regime the exact copy number was selected using an appropriately scaled gamma distribution ( $\mu = 4$ ,  $\sigma = 1$ ) or the minimum edge of the regime. Finally, the effective plasmid relative abundance was determined by multiplying the plasmid copy number with the genome relative abundance. The full set of randomly assigned relative abundances and copy numbers can be found in Supplemental Table 3.

Sequences were then concatenated into a single fasta file with the appropriate relative abundance. MiSeq v3 250bp paired-end reads with a mean fragment length of 1000bp (standard deviation of 50bp) were then simulated using `art_illumina` (v2016.06.05) [46] at a fold coverage of 2.9 resulting in a simulated metagenome of 31,174,411 read pairs.

The selection of relative abundance and metagenome simulation itself was performed using the `data_simulation/simulate_metagenome.py` script.

### Metagenome Assembled Genome Recovery

Reads were trimmed using `sickle` (v1.33) [47] resulting in 25,682,644 surviving read pairs. The trimmed reads were then assembled using 3 different metagenomic assemblers: metaSPAdes (v3.13.0) [17], IDBA-UD (v1.1.3) [18], and megahit (v1.1.3) [19]. The resulting assemblies were summarised using metaQUAST (v5.0.2) [48]. The assemblies were indexed and reads mapped back using Bowtie 2 (v2.3.4.3) [13].

Samtools (v1.9) were then used to sort the read mappings and the read coverage calculated using the MetaBAT 2 accessory script (jgi\_summarize\_bam\_contig\_depths). The 3 metagenome assemblies were then separately binned using CONCOCT (v0.4.2) [21], MetaBAT 2 (v2.13) [22], and MaxBin 2 (v2.2.6) [23]. As per the specified manual instructions CONCOCT used a slightly different approach to estimate read coverage. The supplied accessory scripts were also used to cut contigs into 10 kilobase fragments (cut\_up\_fasta.py) and read coverage calculated for the fragments (concoct\_coverage\_table.py). These fragment coverages were then used to bin the 10kb fragments before the clustered fragments were merged (merge\_cutup\_clustering.py) to create the final CONCOCT MAG bins (extra\_fasta\_bins.py). Finally, for each metagenome assembly the predicted bins from these 3 binners were combined using DAS Tool (v1.1.1) [24]. This resulted in 12 separate sets of MAGs (one set for each assembler and binner pair).

## MAG assessment

### Chromosomal Coverage

The MAG assessment for chromosomal coverage was performed by creating a BLASTN 2.9.0+ [49] database consisting of all the chromosomes of the input reference genomes. Each MAG contig was then used as a query against this database and the coverage of the underlying chromosomes tallied by merging the overlapping aligning regions and summing the total length of aligned MAG contigs. The most represented genome in each MAG was assigned as the “identity” of that MAG for further analyses. Coverages less than 5% were filtered out and the number of different genomes that a MAG contain contigs aligning to were tallied. Finally, the overall proportion of chromosomes that were not present in any MAG were tallied for each binner and assembler.

### Plasmid and GI Coverage

Plasmid and GI coverage were assessed in the same way as one another. Firstly, each a BLASTN database was generated for each set of MAG contigs. Then each MAG database was searched for plasmid and GI sequences. Any plasmid or GI with greater than 50% coverage in a MAG was retained. All plasmids or GIs which could be found in the unbinned contigs or the MAGs was recorded as having been successfully assembled. The subset of these which were found in the binned MAGs was then separately tallied. Finally, we evaluated the proportion of plasmids or GIs which were binned correctly in the bin which was maximally composed of chromosomes from the same source genome. This was determined using the bin “IDs” from the chromosomal coverage analysis.

## Antimicrobial Resistance and Virulence Factors Assessment

### Detection of AMR/VF Genes

For each of the 12 MAGs, and the reference chromosome and plasmids, AMR genes were predicted using Resistance Gene Identifier (RGI v5.0.0; default parameters) and the Comprehensive Antibiotic Resistance Database (CARD v3.0.2) [50]. Virulence factors were predicted using BLASTX against the Virulence Factors Database (VFDB; obtained on Aug 26, 2019) with an e-value cut-off of 0.001. [51]. Each MAG was then assigned to a reference chromosome and plasmid using the above mentioned mapping criteria for downstream analysis.

### AMR/VF Gene Recovery

The ability for MAGs to properly recover AMR and VF genes was assessed using only the megahit-DasTool assembler-binner combination as it was the best performing pair. For each bin, we counted the total number of AMR/VF genes recovered then compared this to the number predicted in their assigned reference chromosome and plasmids to determine MAG’s gene recovery ability. We then mapped the location of reference replicon’s predicted genes to the bins to determined the location of those genes in MAGs.

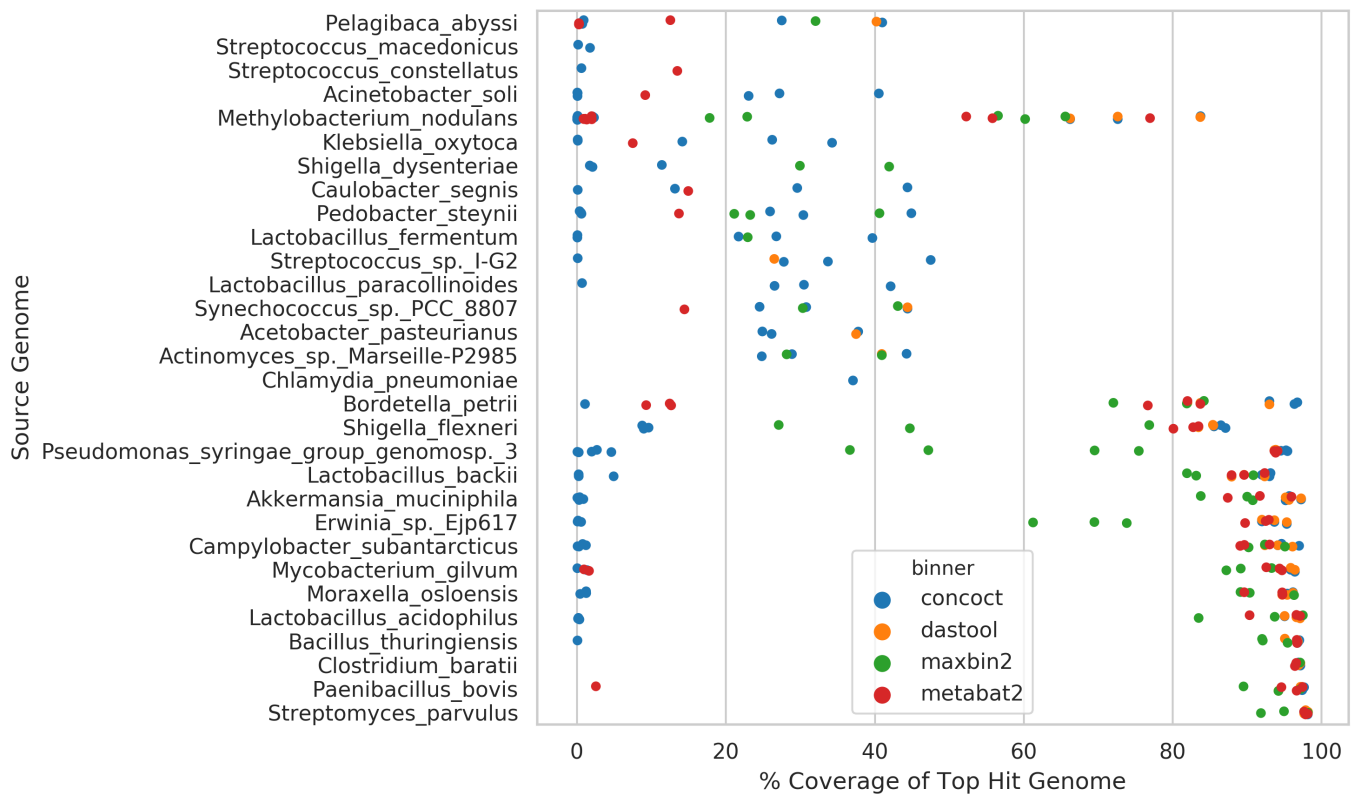
### Protein subcellular localization predictions

The MAG bins from megahit-DasTool assembler-binner combination was inputted into prodigal [52] to predict open reading frames (ORFs) using the default parameter. The list of predicted proteins is inputted into PSORTb v3.0 with default parameters [53].

## Results

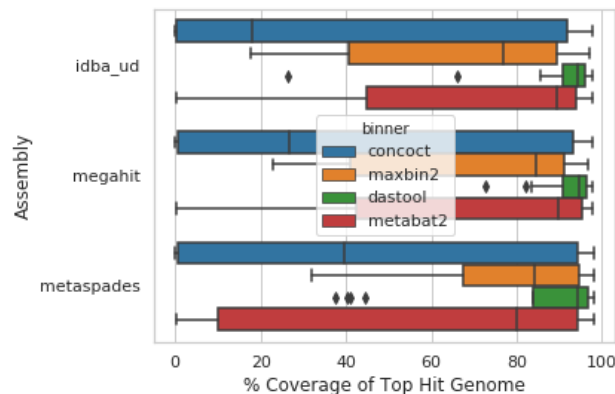
---

The overall ability of MAG methods to recapitulate the original chromosomal source genome results varied widely. We consider the identity of a given MAG bin to be that of the genome that composes the largest proportion of sequence within that bin. In other words if a bin is identifiably 70% species A and 30% species B we consider that to be a bin of species A. Ideally, we wish to generate a single bin for each source genome comprised of the entire genome and no contigs from other genomes. Some genomes are cleanly and accurately binned regardless of the assembler and binning method used (See Fig. (6)). Specifically, greater than 90% of *Streptomyces parvulus* (minimum 91.8%) and *Clostridium baratii* (minimum 96.4%) chromosomes are represented in individual bins across all methods. However, no other genomes were consistently recovered by all methods for more than a 3rd of the chromosomes. The 3 *Streptococcus* genomes were particularly problematic, likely due to their similarity, with the best recovery for each ranging from 1.7% to 47.49%.



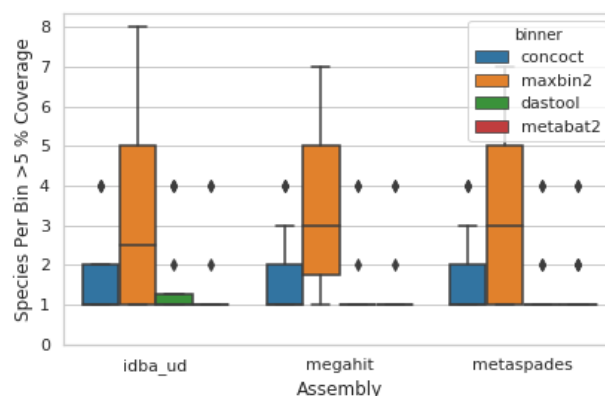
**Figure 1:** Top Species Coverage

In terms of assembler, megahit resulted in the highest median chromosomal coverage across all binners (81.9%) with metaSPAdes performing worst (76.8%) (see Fig. (2)). In terms of binning tool, CONCOCT performed very poorly with a median 26% coverage for top hit per bin, followed by maxbin2 (83.1%), and metabat2 (88.5%). It is perhaps unsurprising that the best performing binner in terms of bin top hit coverage was the metabinner DAS-TOOL that combines predictions from the other 3 binners (94.3% median top hit chromosome coverage per bin, (2)).



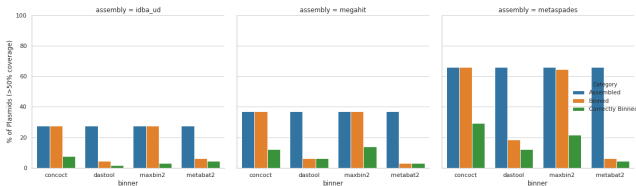
**Figure 2:** Chromosomal coverages of most prevalent genome in each bin across binners and metagenome assemblies

Bin purity, i.e. the number of genomes present in a bin at >5% coverage, was largely equivalent across assemblers (see Fig. (3)), with a very marginally higher purity for IDBA. In terms of binning tool, however, maxbin2 proved an outlier with nearly twice as many bins containing multiple species as the next binner. The remaining binning tools were largely equivalent, producing chimeric bins at approximately the same rates.



**Figure 3:** Distribution of bin purities. Showing the number of genomes with >5% chromosomal coverage across the bins and methods

Regardless of method, a very small proportion of plasmids were correctly binned in the bin that mostly contained chromosomal contigs from the same source genome. Specifically, between 1.5% (IDBA-UD assembly with DAS Tool bins) and 29.2% (metaSPAdes with CONCOCT bins) were correctly binned at over 50% coverage. In terms of metagenome assembly MetaSPAdes was far and away the most successful assembler at assembling plasmids with 66.2% of plasmids identifiable at greater than 50% coverage. IDBA-UD performed worst with 17.1% of plasmids recovered, and megahit recovered 36.9%. If the plasmid was successfully assembled it was placed in a bin by maxbin2 and CONCOCT, although a much smaller fraction correctly binned (typically less than 1/3rd). Interestingly, metatbat2 and DAS tool binners were a lot more conservative in assigning plasmid contigs to bins, however, of those assigned to bins nearly all were correctly binned (see Fig. (4)).

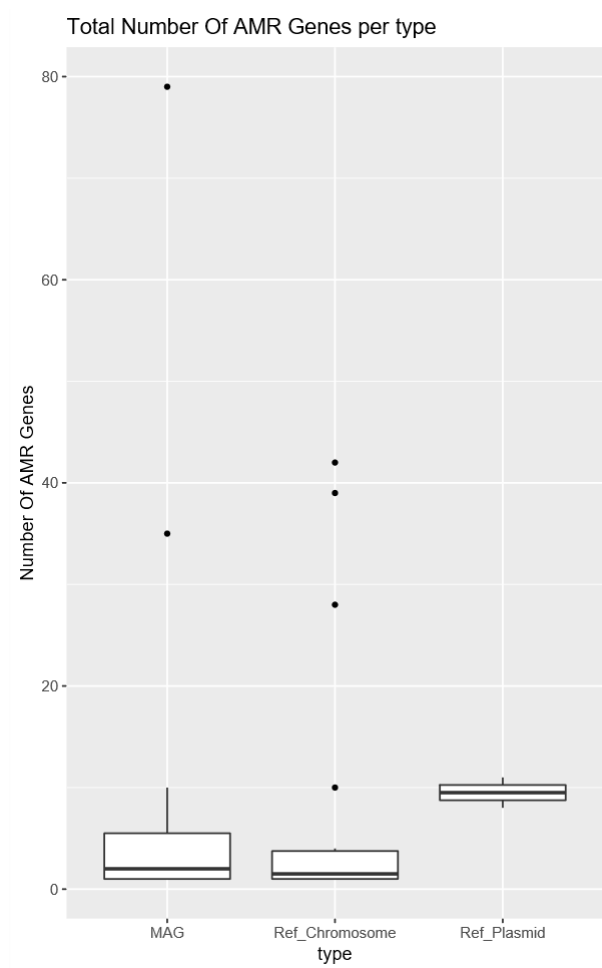
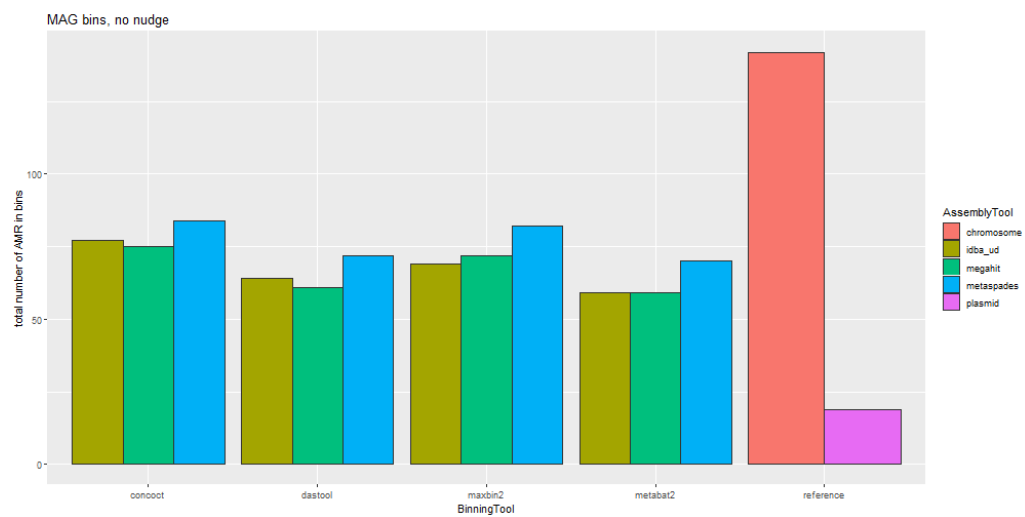


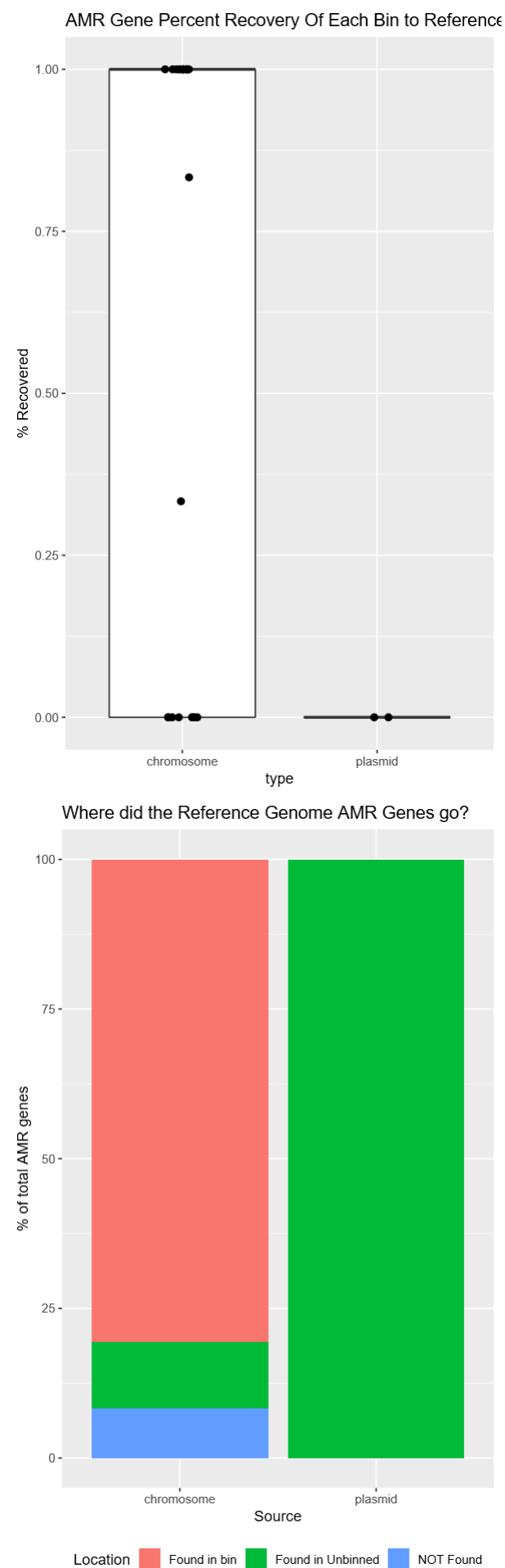
**Figure 4:** Plasmid Coverage

GIs displayed a similar pattern of assembly and correct binning performance as plasmids (see Fig. ([???])). These sequences were assembled uniformly badly (37.8-44.1%) with metaSPAdes outperforming the other two assembly approaches. For CONCOCT and maxbin2 binning tools all GIs that were assembled were assigned to a bin although the proportion of binned GIs that were correctly binned was lower than for DAS Tool and metatbat2. DAS Tool, metatbat2 and CONCOCT didn't display the same precipitous drop-off between those assembled and those correctly binned as was observed for plasmids. In terms of overall correct binning with the chromosomes from the same genome the metaSPAdes assembly with CONCOCT (44.1%) and maxbin2 (43.3%) binners performed best.



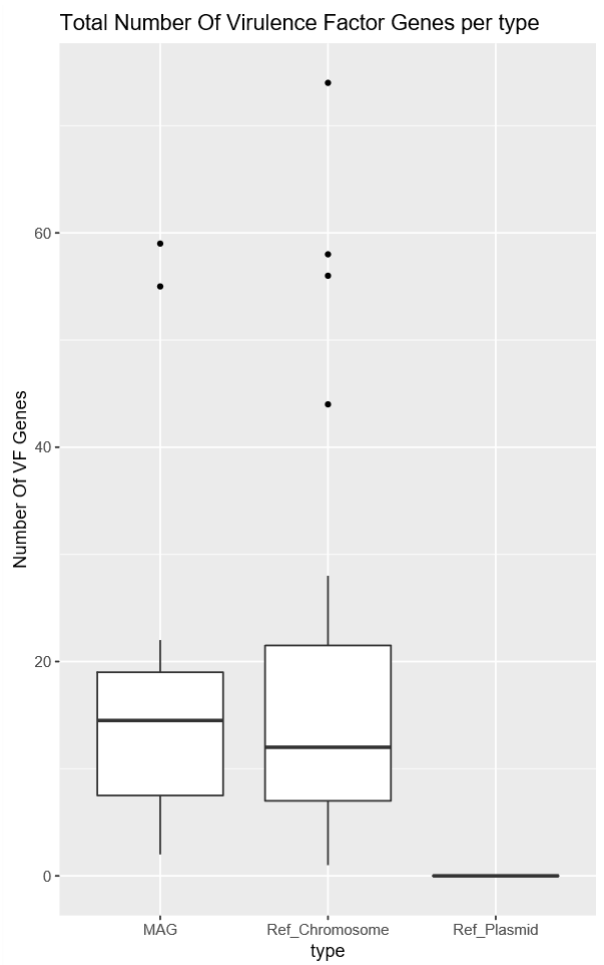
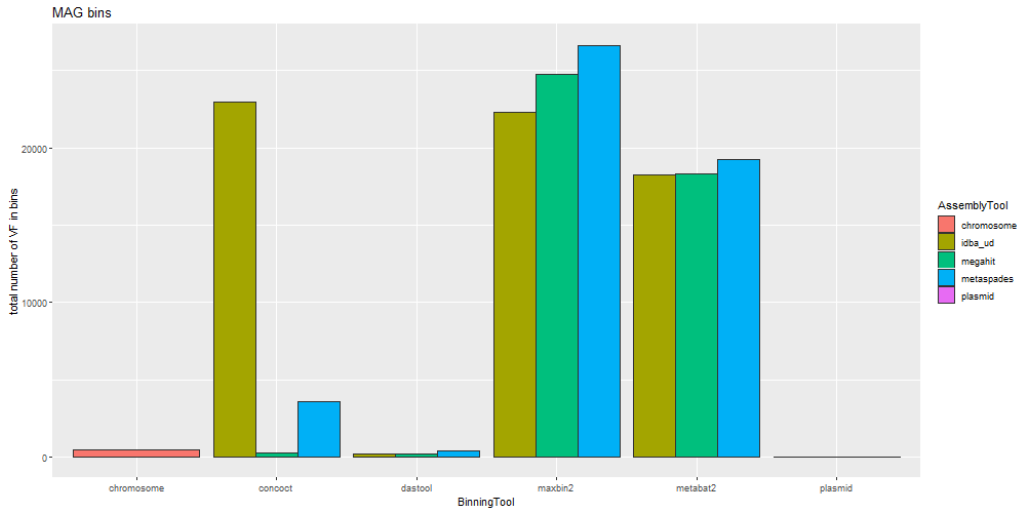
Switching over to gene content, we first explored the ability to find open reading frames (ORFs) within MAGs. Overall, the total number of predicted ORFs in MAGs followed a similar trend to the chromosomal coverage and purity ([???]). Of the 4 binning tools, CONCOCT performed the worst, finding <30% of the number of ORFs in our reference genomes. Metatbat2 performed second worst at ~80%. DASTool recovered a similar number to our reference and Maxbin2 seemed to predicted 7-46% more genes. The Assembler method did not significantly impact the number of genes predicted with the exception of Maxbin2 in which idba\_ud was the closest to reference and metaspades predicted 46% more ORFs.



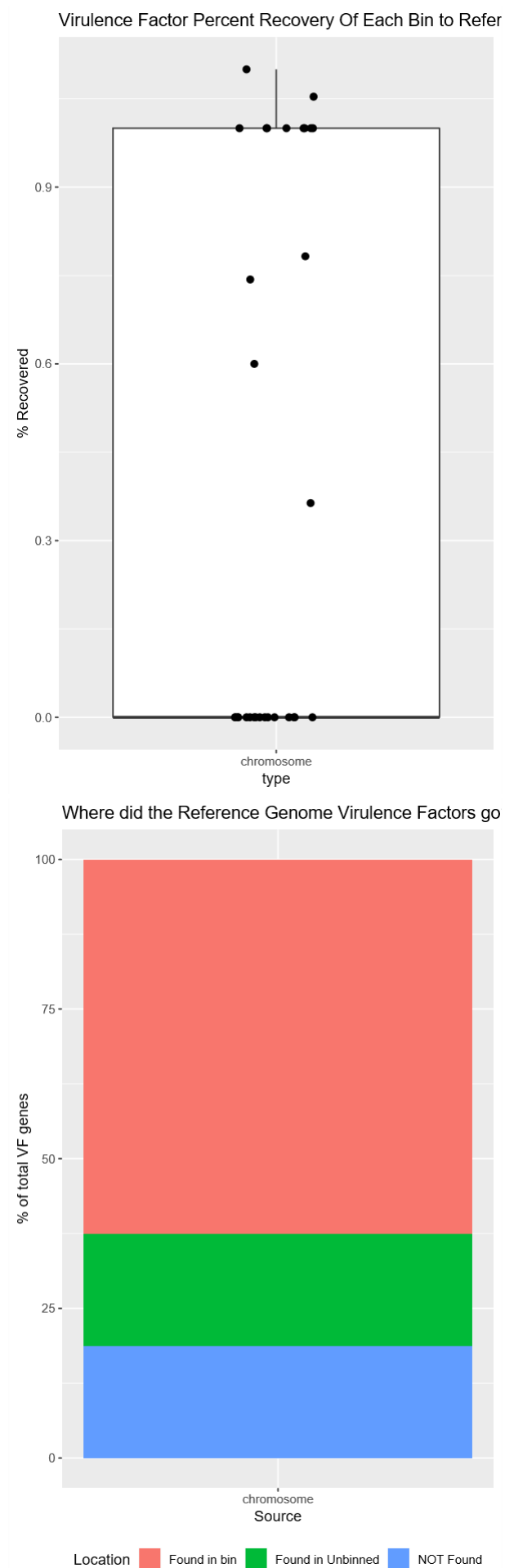


With respect to AMR genes, in total, MAGs were only able to recover between 40-53% of the AMR genes predicted in our reference genomes across all assembler-binner pairs (Fig. ([???])). We then took the best assembler-binner pair (MegaHit-DasTools) and examined the AMR genes recovered in detail. We noticed that, for majority of the bins (85%), MAGs were able to correctly recover either 100% or 0% of the AMR genes (Median value 100%) that are

contained in the reference chromosome assigned to that bin. However, MAGs were not able to correctly recover any of the AMR genes that were present on plasmids (Fig. ([???]), Fig. ([???])). Lastly, we asked the question of where reference replicon AMR genes went in the MAGs. For chromosome, majority (81%) of the AMR genes was found in a bin of the MAG. A small portion (12%) was left unbinned and 7% were not found in MAGs at all. On the other hand, for plasmid born AMR genes, all of the recovered genes (n=20) were identified in the unbinned fraction of our MAG (Fig. ([???])).

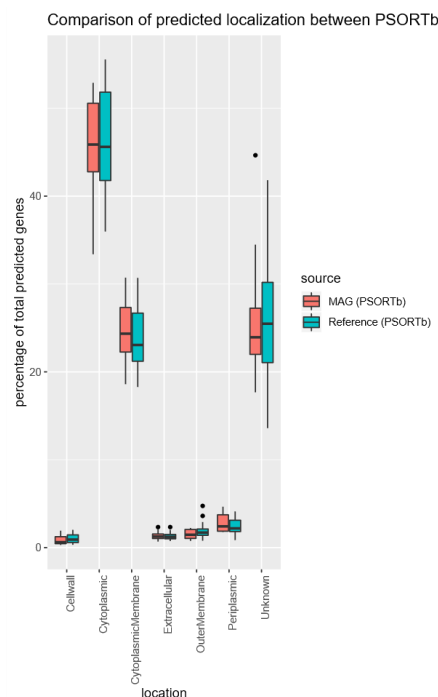






Aside from AMR genes, we also examined virulence factors in our dataset. The number of VF varied dramatically across assembler-binner combinations. Maxbin2 and metabat2 binned MAGs, regardless of assembly method, predicted over 20,000 virulence genes, roughly 50x more than our reference genome. The DasTool-megahit MAG predicted just under 250 VFs, which represent 61% of the number of VFs in our reference replicons (Fig. ([???])).

Furthermore, each bin's percent recovery of predicted VF compared to the number in their assigned reference chromosome varied much more compared to AMR genes. The %recovered value ranges from 0% to 112% (Median 0). (Fig. ([???]), Fig. ([???])). Finally, we examined the location of reference chromosomal VF genes in our MAG bins. 63% of VF genes were found in a MAG bin, 19% were found in the unbinned portion and 18% were not found at all. There were no VF predicted on reference plasmids (Fig. ([???])).



**Figure 5:** Distribution of Predicted Protein Subcellular Localization

Lastly, we looked at the ability for MAGs to predict subcellular localization of proteins using PSORTb. Overall, the localization distribution of predicted proteins were very similar in MAGs compared to the reference genome (Fig. (5)).

## Discussion

In this paper, we evaluated the ability and accuracy of metagenome-assembled genomes (MAGs) to correctly recover mobile genetic elements (i.e. genomic islands and plasmids) from metagenomic samples across different tools used to assemble and bin MAGs.

Overall, the best assembler-binner pair was megahit-DASTOOL in term of both chromosomal coverage (94.3%) and bin purity (1). Looking at genomes with the lowest coverage, the 3 *Streptococcus* genomes were particularly problematic, likely due to their similarity, with the best recovery for each ranging from 1.7% to 47.49%. This suggest that MAGs might not be able to distinguish between closely related species (COMMENT: Point 1, MAGs cannot distinguish closely related species). While CONCOCT performed significantly worse compared to the other bidders, we did notice that CONCOCT seems to display a trend of generating lots of small partial bins. Perhaps CONCOCT bins might be able to distinguish between closely related species to a higher resolution (COMMENT: Small partial bins... what does it mean overall. is this assumption correct? Would it be able to distinguish closely related species?)

While the overall recovery of chromosomes was okay, we were interested in MAG's ability to correctly bin mobile genetic elements due to their importance in the functions and spread of pathogenic traits such as AMR and virulence. In term of plasmids, a very small proportion of plasmids were correctly binned regardless of the method (<33% at best). Similarly, the same trend exists for genomic islands (<43.3%). This poor result is not unexpected as genomic islands and plasmids have divergent composition features relative to the chromosomes. Furthermore, the difference between the percentages suggest that binning plasmids are harder than GIs. This difference might be due to the problem of plasmid assembly. Therefore, the binning efficiency might improve if we use an assembler targeted at assembling plasmids [54].

Looking at predicted gene content, our best assembler-binner pair produced a similar number of predicted ORFs as our reference genomes. (interestingly we still missed a bunch of AMR genes. perhaps theses predicted ORFs are fragmented? idk need ideas)

Due to the importance of mobile genetic elements to disseminate clinically relevant antimicrobial resistance genes and virulence factors, we explored whether or not MAGs can be used to provide useful lateral gene transfer insights.

With respect to AMR genes, MAGs were able to recover roughly half of all AMR genes present in our reference genome. The correct bins were assigned for majority of the chromosomally located AMR genes (81%). The accuracy of chromosomal AMR genes were as expected given the accuracy of MAGs to recover chromosomes as discussed previously. However, while MAGs were able to detect all of plasmid-born AMR genes, none of these were placed in any of the bins. We specifically included a few high threat AMR genes in our dataset: namely KPC and OXA, which are plasmid borne carbapenemases of increasing prevalence in the clinics that are rendering our last resort antibiotics useless. These genes were successfully detected from the metagenomics assembly, but they were not assigned to a bin. This could mean a limited ability for MAGs to be used in the public health research to pinpoint the lateral transfer of AMR genes and to conduct epidemiological analysis (COMMENT: does this make sense?).

Virulence factors had shown a similar trend as AMR genes, recovering ~60% of virulence factors present in the reference genome. Interestingly, while the detection of virulence factors is better than AMR genes, the binning accuracy was worse, with more being present in the unbinned fraction. Previous studies has found that VFs are disproportionally present on GIs[33], which might be the reason to why the binning accuracy was worse compared to AMR genes.

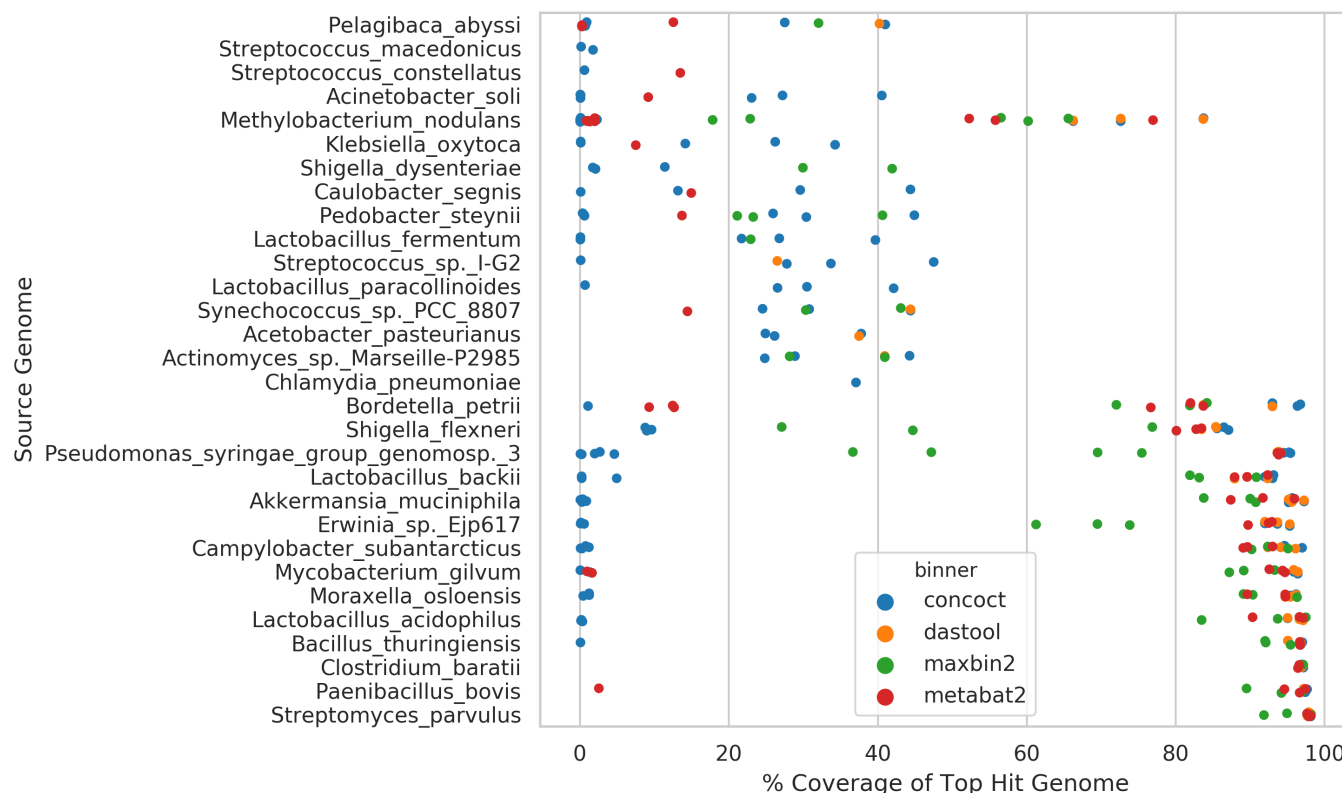
Lastly, previous works have shown that AMR genes that are on mobile genetic elements disproportionally encode secrete proteins. Given that the recovery of plasmid-borne genes were not great, we asked if MAGs would affect the ability to predict the subcellular localization of proteins. We found

that the proportion of predicted localizations were very similar between MAGs and our reference genomes, suggesting that there is not a significant penalty to use MAGs as input for protein localization predictions.

## Conclusions

Using a simulated medium complexity metagenome, this study had shown that MAGs provides a great tool to study a bacterial species' chromosomal elements but presented difficulties in the recovery of mobile genetic elements from metagenomic samples. These mobile genetic elements are liable to being incorrectly binned or lost in this process. Due to the importance of these mobile genomic components in the function and spread of pathogenic traits such as AMR and virulence, it is vital that we utilize a combination of MAGs and other methods (e.g. read-based methods) in public health metagenomic researches. This would allow both the detection of the sample microbial diversity and the thorough evaluation of resistome in metagenomic data to provide meaningful epidemiological information.

## Supplementals



**Figure 6:** Top Species Coverage

### 1. Genomic analysis of uncultured marine viral communities

M. Breitbart, P. Salamon, B. Andresen, J. M. Mahaffy, A. M. Segall, D. Mead, F. Azam, F. Rohwer  
*Proceedings of the National Academy of Sciences* (2002-10-16) <https://doi.org/br7jq3>  
 DOI: [10.1073/pnas.202488399](https://doi.org/10.1073/pnas.202488399) · PMID: [12384570](https://pubmed.ncbi.nlm.nih.gov/12384570/) · PMCID: [PMC137870](https://pubmed.ncbi.nlm.nih.gov/PMC137870/)

### 2. Shotgun metagenomics, from sampling to analysis

Christopher Quince, Alan W Walker, Jared T Simpson, Nicholas J Loman, Nicola Segata  
*Nature Biotechnology* (2017-09) <https://doi.org/gbv6nf>  
 DOI: [10.1038/nbt.3935](https://doi.org/10.1038/nbt.3935) · PMID: [28898207](https://pubmed.ncbi.nlm.nih.gov/28898207/)

### 3. Analysis of a marine picoplankton community by 16S rRNA gene cloning and sequencing.

TM Schmidt, EF DeLong, NR Pace  
*Journal of bacteriology* (1991-07) <https://www.ncbi.nlm.nih.gov/pubmed/2066334>  
 DOI: [10.1128/jb.173.14.4371-4378.1991](https://doi.org/10.1128/jb.173.14.4371-4378.1991) · PMID: [2066334](https://pubmed.ncbi.nlm.nih.gov/2066334/) · PMCID: [PMC208098](https://pubmed.ncbi.nlm.nih.gov/PMC208098/)

### 4. Predictive functional profiling of microbial communities using 16S rRNA marker gene sequences

Morgan GI Langille, Jesse Zaneveld, J Gregory Caporaso, Daniel McDonald, Dan Knights, Joshua A Reyes, Jose C Clemente, Deron E Burkepille, Rebecca L Vega Thurber, Rob Knight, ... Curtis Huttenhower  
*Nature Biotechnology* (2013-08-25) <https://doi.org/f49xzcd>  
 DOI: [10.1038/nbt.2676](https://doi.org/10.1038/nbt.2676) · PMID: [23975157](https://pubmed.ncbi.nlm.nih.gov/23975157/) · PMCID: [PMC3819121](https://pubmed.ncbi.nlm.nih.gov/PMC3819121/)

### 5. PICRUSt2: An improved and extensible approach for metagenome inference

Gavin M. Douglas, Vincent J. Maffei, Jesse Zaneveld, Svetlana N. Yurgel, James R. Brown, Christopher M. Taylor, Curtis Huttenhower, Morgan G. I. Langille  
*Cold Spring Harbor Laboratory* (2019-06-15) <https://doi.org/gf5ffb>  
 DOI: [10.1101/672295](https://doi.org/10.1101/672295)

### 6. A Systematic Analysis of Biosynthetic Gene Clusters in the Human Microbiome Reveals a Common Family of Antibiotics

Mohamed S. Donia, Peter Cimermancic, Christopher J. Schulze, Laura C. Wieland Brown, John Martin, Makedonka Mitreva, Jon Clardy, Roger G. Linington, Michael A. Fischbach

Cell (2014-09) <https://doi.org/f6k3fg>  
DOI: [10.1016/j.cell.2014.08.032](https://doi.org/10.1016/j.cell.2014.08.032) · PMID: [25215495](https://pubmed.ncbi.nlm.nih.gov/25215495/) · PMCID: [PMC4164201](https://pubmed.ncbi.nlm.nih.gov/PMC4164201/)

#### 7. Expanding the soil antibiotic resistome: exploring environmental diversity

Vanessa M D'Costa, Emma Griffiths, Gerard D Wright  
*Current Opinion in Microbiology* (2007-10) <https://doi.org/cfbpjj>  
DOI: [10.1016/j.mib.2007.08.009](https://doi.org/10.1016/j.mib.2007.08.009) · PMID: [17951101](https://pubmed.ncbi.nlm.nih.gov/17951101/)

#### 8. Antibiotic resistance is ancient

Vanessa M. D'Costa, Christine E. King, Lindsay Kalan, Mariya Morar, Wilson W. L. Sung, Carsten Schwarz, Duane Froese, Grant Zazula, Fabrice Calmels, Regis Debruyne, ... Gerard D. Wright  
*Nature* (2011-08-31) <https://doi.org/b3wbvx>  
DOI: [10.1038/nature10388](https://doi.org/10.1038/nature10388) · PMID: [21881561](https://pubmed.ncbi.nlm.nih.gov/21881561/)

#### 9. A Culture-Independent Sequence-Based Metagenomics Approach to the Investigation of an Outbreak of Shiga-Toxigenic *Escherichia coli* O104:H4

Nicholas J. Loman, Chrystala Constantinidou, Martin Christner, Holger Rohde, Jacqueline Z.-M. Chan, Joshua Quick, Jacqueline C. Weir, Christopher Quince, Geoffrey P. Smith, Jason R. Betley, ... Mark J. Pallen  
*JAMA* (2013-04-10) <https://doi.org/f5rqft>  
DOI: [10.1001/jama.2013.3231](https://doi.org/10.1001/jama.2013.3231) · PMID: [23571589](https://pubmed.ncbi.nlm.nih.gov/23571589/)

#### 10. Ultra-deep, long-read nanopore sequencing of mock microbial community standards

Samuel M Nicholls, Joshua C Quick, Shuiquan Tang, Nicholas J Loman  
*GigaScience* (2019-05-01) <https://doi.org/gf39g3>  
DOI: [10.1093/gigascience/giz043](https://doi.org/10.1093/gigascience/giz043) · PMID: [31089679](https://pubmed.ncbi.nlm.nih.gov/31089679/) · PMCID: [PMC6520541](https://pubmed.ncbi.nlm.nih.gov/PMC6520541/)

#### 11. Long-read based de novo assembly of low-complexity metagenome samples results in finished genomes and reveals insights into strain diversity and an active phage system

Vincent Somerville, Stefanie Lutz, Michael Schmid, Daniel Frei, Aline Moser, Stefan Irmeler, Jürg E. Frey, Christian H. Ahrens  
*BMC Microbiology* (2019-06-25) <https://doi.org/gf5ffc>  
DOI: [10.1186/s12866-019-1500-0](https://doi.org/10.1186/s12866-019-1500-0) · PMID: [31238873](https://pubmed.ncbi.nlm.nih.gov/31238873/) · PMCID: [PMC6593500](https://pubmed.ncbi.nlm.nih.gov/PMC6593500/)

#### 12. Fast and sensitive protein alignment using DIAMOND

Benjamin Buchfink, Chao Xie, Daniel H Huson  
*Nature Methods* (2014-11-17) <https://doi.org/gftzcs>  
DOI: [10.1038/nmeth.3176](https://doi.org/10.1038/nmeth.3176) · PMID: [25402007](https://pubmed.ncbi.nlm.nih.gov/25402007/)

#### 13. Fast gapped-read alignment with Bowtie 2

Ben Langmead, Steven L Salzberg  
*Nature Methods* (2012-03-04) <https://doi.org/gd2xzn>  
DOI: [10.1038/nmeth.1923](https://doi.org/10.1038/nmeth.1923) · PMID: [22388286](https://pubmed.ncbi.nlm.nih.gov/22388286/) · PMCID: [PMC3322381](https://pubmed.ncbi.nlm.nih.gov/PMC3322381/)

#### 14. nhmmer: DNA homology search with profile HMMs

T. J. Wheeler, S. R. Eddy  
*Bioinformatics* (2013-07-09) <https://doi.org/f5xm9x>  
DOI: [10.1093/bioinformatics/btt403](https://doi.org/10.1093/bioinformatics/btt403) · PMID: [23842809](https://pubmed.ncbi.nlm.nih.gov/23842809/) · PMCID: [PMC3777106](https://pubmed.ncbi.nlm.nih.gov/PMC3777106/)

#### 15. CLARK: fast and accurate classification of metagenomic and genomic sequences using discriminative k-mers

Rachid Ounit, Steve Wanamaker, Timothy J Close, Stefano Lonardi  
*BMC Genomics* (2015-03-25) <https://doi.org/gb3h2t>  
DOI: [10.1186/s12864-015-1419-2](https://doi.org/10.1186/s12864-015-1419-2) · PMID: [25879410](https://pubmed.ncbi.nlm.nih.gov/25879410/) · PMCID: [PMC4428112](https://pubmed.ncbi.nlm.nih.gov/PMC4428112/)

#### 16. Community structure and metabolism through reconstruction of microbial genomes from the environment

Gene W. Tyson, Jarrod Chapman, Philip Hugenholtz, Eric E. Allen, Rachna J. Ram, Paul M. Richardson, Victor V. Solovyev, Edward M. Rubin, Daniel S. Rokhsar, Jillian F. Banfield  
*Nature* (2004-02-01) <https://doi.org/b85j5j>  
DOI: [10.1038/nature02340](https://doi.org/10.1038/nature02340) · PMID: [14961025](https://pubmed.ncbi.nlm.nih.gov/14961025/)

#### 17. metaSPAdes: a new versatile metagenomic assembler

Sergey Nurk, Dmitry Meleshko, Anton Korobeynikov, Pavel A. Pevzner  
*Genome Research* (2017-03-15) <https://doi.org/f97jky>  
DOI: [10.1101/gr.213959.116](https://doi.org/10.1101/gr.213959.116) · PMID: [28298430](https://pubmed.ncbi.nlm.nih.gov/28298430/) · PMCID: [PMC5411777](https://pubmed.ncbi.nlm.nih.gov/PMC5411777/)

#### 18. IDBA-UD: a de novo assembler for single-cell and metagenomic sequencing data with highly uneven depth

Y. Peng, H. C. M. Leung, S. M. Yiu, F. Y. L. Chin  
*Bioinformatics* (2012-04-11) <https://doi.org/f3z7hv>  
DOI: [10.1093/bioinformatics/bts174](https://doi.org/10.1093/bioinformatics/bts174) · PMID: [22495754](https://pubmed.ncbi.nlm.nih.gov/22495754/)

#### 19. MEGAHIT: an ultra-fast single-node solution for large and complex metagenomics assembly via succinct de Bruijn graph

Dinghua Li, Chi-Man Liu, Ruibang Luo, Kunihiro Sadakane, Tak-Wah Lam  
*Bioinformatics* (2015-01-20) <https://doi.org/f7fb5z>  
DOI: [10.1093/bioinformatics/btv033](https://doi.org/10.1093/bioinformatics/btv033) · PMID: [25609793](https://pubmed.ncbi.nlm.nih.gov/25609793/)

#### 20. Assembling metagenomes, one community at a time

Andries Johannes van der Walt, Marc Warwick van Goethem, Jean-Baptiste Ramond, Thulani Peter Makhalanyane, Oleg Reva, Don Arthur Cowan  
*BMC Genomics* (2017-07-10) <https://doi.org/gf5fhs>  
DOI: [10.1186/s12864-017-3918-9](https://doi.org/10.1186/s12864-017-3918-9) · PMID: [28693474](https://pubmed.ncbi.nlm.nih.gov/28693474/) · PMCID: [PMC5502489](https://pubmed.ncbi.nlm.nih.gov/PMC5502489/)

#### 21. COCACOLA: binning metagenomic contigs using sequence COMposition, read CoverAge, CO-alignment and paired-end read LinkAge

Yang Young Lu, Ting Chen, Jed A. Fuhrman, Fengzhu Sun

Bioinformatics (2016-06-02) <https://doi.org/f9x7sc>  
DOI: [10.1093/bioinformatics/btw290](https://doi.org/10.1093/bioinformatics/btw290) · PMID: [27256312](https://pubmed.ncbi.nlm.nih.gov/27256312/)

**22. MetaBAT 2: an adaptive binning algorithm for robust and efficient genome reconstruction from metagenome assemblies**

Dongwan Kang, Feng Li, Edward S Kirton, Ashleigh Thomas, Rob S Egan, Hong An, Zhong Wang  
*PeerJ* (2019-02-06) <https://doi.org/gf5fhv>  
DOI: [10.7287/peerj.preprints.27522v1](https://doi.org/10.7287/peerj.preprints.27522v1)

**23. MaxBin 2.0: an automated binning algorithm to recover genomes from multiple metagenomic datasets**

Yu-Wei Wu, Blake A. Simmons, Steven W. Singer  
*Bioinformatics* (2015-10-29) <https://doi.org/f8c9n2>  
DOI: [10.1093/bioinformatics/btv638](https://doi.org/10.1093/bioinformatics/btv638) · PMID: [26515820](https://pubmed.ncbi.nlm.nih.gov/26515820/)

**24. Recovery of genomes from metagenomes via a dereplication, aggregation and scoring strategy**

Christian M. K. Sieber, Alexander J. Probst, Allison Sharrar, Brian C. Thomas, Matthias Hess, Susannah G. Tringe, Jillian F. Banfield  
*Nature Microbiology* (2018-05-28) <https://doi.org/gfwwfg>  
DOI: [10.1038/s41564-018-0171-1](https://doi.org/10.1038/s41564-018-0171-1) · PMID: [29807988](https://pubmed.ncbi.nlm.nih.gov/29807988/) · PMCID: [PMC6786971](https://pubmed.ncbi.nlm.nih.gov/PMC6786971/)

**25. Unusual biology across a group comprising more than 15% of domain Bacteria**

Christopher T. Brown, Laura A. Hug, Brian C. Thomas, Itai Sharon, Cindy J. Castelle, Andrea Singh, Michael J. Wilkins, Kelly C. Wrighton, Kenneth H. Williams, Jillian F. Banfield  
*Nature* (2015-06-15) <https://doi.org/f7h5xj>  
DOI: [10.1038/nature14486](https://doi.org/10.1038/nature14486) · PMID: [26083755](https://pubmed.ncbi.nlm.nih.gov/26083755/)

**26. Recovery of nearly 8,000 metagenome-assembled genomes substantially expands the tree of life**

Donovan H. Parks, Christian Rinke, Maria Chuvochina, Pierre-Alain Chaumeil, Ben J. Woodcroft, Paul N. Evans, Philip Hugenholtz, Gene W. Tyson  
*Nature Microbiology* (2017-09-11) <https://doi.org/cczd>  
DOI: [10.1038/s41564-017-0012-7](https://doi.org/10.1038/s41564-017-0012-7) · PMID: [28894102](https://pubmed.ncbi.nlm.nih.gov/28894102/)

**27. The genomic and proteomic landscape of the rumen microbiome revealed by comprehensive genome-resolved metagenomics**

Robert D. Stewart, Marc D. Auffret, Amanda Warr, Alan W. Walker, Rainer Roehe, Mick Watson  
*Cold Spring Harbor Laboratory* (2018-12-08) <https://doi.org/gf5fhr>  
DOI: [10.1101/489443](https://doi.org/10.1101/489443)

**28. Genome-centric view of carbon processing in thawing permafrost**

Ben J. Woodcroft, Caitlin M. Singleton, Joel A. Boyd, Paul N. Evans, Joanne B. Emerson, Ahmed A. F. Zayed, Robert D. Hoelzle, Timothy O. Lamberton, Carmody K. McCalley, Suzanne B. Hodgkins, ... Gene W. Tyson  
*Nature* (2018-07-16) <https://doi.org/gdth6p>  
DOI: [10.1038/s41586-018-0338-1](https://doi.org/10.1038/s41586-018-0338-1) · PMID: [30013118](https://pubmed.ncbi.nlm.nih.gov/30013118/)

**29. Mediterranean grassland soil C-N compound turnover is dependent on rainfall and depth, and is mediated by genomically divergent microorganisms**

Spencer Diamond, Peter F. Andeer, Zhou Li, Alexander Crits-Christoph, David Burstein, Karthik Anantharaman, Katherine R. Lane, Brian C. Thomas, Chongle Pan, Trent R. Northen, Jillian F. Banfield  
*Nature Microbiology* (2019-05-20) <https://doi.org/gf5fcx>  
DOI: [10.1038/s41564-019-0449-y](https://doi.org/10.1038/s41564-019-0449-y) · PMID: [31110364](https://pubmed.ncbi.nlm.nih.gov/31110364/) · PMCID: [PMC6784897](https://pubmed.ncbi.nlm.nih.gov/PMC6784897/)

**30. AMBER: Assessment of Metagenome BinnERS**

Fernando Meyer, Peter Hofmann, Peter Belmann, Ruben Garrido-Oter, Adrian Fritz, Alexander Sczyrba, Alice C McHardy  
*GigaScience* (2018-06-01) <https://doi.org/gdptz9>  
DOI: [10.1093/gigascience/giy069](https://doi.org/10.1093/gigascience/giy069) · PMID: [29893851](https://pubmed.ncbi.nlm.nih.gov/29893851/) · PMCID: [PMC6022608](https://pubmed.ncbi.nlm.nih.gov/PMC6022608/)

**31. Detecting genomic islands using bioinformatics approaches**

Morgan G. I. Langille, William W. L. Hsiao, Fiona S. L. Brinkman  
*Nature Reviews Microbiology* (2010-05) <https://doi.org/d6ss55>  
DOI: [10.1038/nrmicro2350](https://doi.org/10.1038/nrmicro2350) · PMID: [20395967](https://pubmed.ncbi.nlm.nih.gov/20395967/)

**32. Horizontal gene transfer: building the web of life**

Shannon M. Soucy, Jinling Huang, Johann Peter Gogarten  
*Nature Reviews Genetics* (2015-07-17) <https://doi.org/f7j3d9>  
DOI: [10.1038/nrg3962](https://doi.org/10.1038/nrg3962) · PMID: [26184597](https://pubmed.ncbi.nlm.nih.gov/26184597/)

**33. The Association of Virulence Factors with Genomic Islands**

Shannan J. Ho Sui, Amber Fedynak, William W. L. Hsiao, Morgan G. I. Langille, Fiona S. L. Brinkman  
*PLoS ONE* (2009-12-01) <https://doi.org/c7hsyv>  
DOI: [10.1371/journal.pone.0008094](https://doi.org/10.1371/journal.pone.0008094) · PMID: [19956607](https://pubmed.ncbi.nlm.nih.gov/19956607/) · PMCID: [PMC2779486](https://pubmed.ncbi.nlm.nih.gov/PMC2779486/)

**34. Dissemination of Antimicrobial Resistance in Microbial Ecosystems through Horizontal Gene Transfer**

Christian J. H. von Wintersdorff, John Penders, Julius M. van Niekerk, Nathan D. Mills, Snehal Majumder, Lieke B. van Alphen, Paul H. M. Savelkoul, Petra F. G. Wolffs  
*Frontiers in Microbiology* (2016-02-19) <https://doi.org/gf5fht>  
DOI: [10.3389/fmicb.2016.00173](https://doi.org/10.3389/fmicb.2016.00173) · PMID: [26925045](https://pubmed.ncbi.nlm.nih.gov/26925045/) · PMCID: [PMC4759269](https://pubmed.ncbi.nlm.nih.gov/PMC4759269/)

**35. Transfer of antibiotic-resistance genes via phage-related mobile elements**

Maryury Brown-Jaque, William Calero-Cáceres, Maite Muniesa  
*Plasmid* (2015-05) <https://doi.org/f7dvxy>  
DOI: [10.1016/j.plasmid.2015.01.001](https://doi.org/10.1016/j.plasmid.2015.01.001) · PMID: [25597519](https://pubmed.ncbi.nlm.nih.gov/25597519/)

36.  
Rainer Merkl

*BMC Bioinformatics* (2004) <https://doi.org/bt5x8h>  
DOI: [10.1186/1471-2105-5-22](https://doi.org/10.1186/1471-2105-5-22) · PMID: [15113412](https://pubmed.ncbi.nlm.nih.gov/15113412/) · PMCID: [PMC394314](https://pubmed.ncbi.nlm.nih.gov/PMC394314/)

**37. Improved genomic island predictions with IslandPath-DIMOB**

Claire Bertelli, Fiona SL Brinkman  
*Bioinformatics* (2018-02-23) <https://doi.org/gdphgs>  
DOI: [10.1093/bioinformatics/bty095](https://doi.org/10.1093/bioinformatics/bty095) · PMID: [29905770](https://pubmed.ncbi.nlm.nih.gov/29905770/) · PMCID: [PMC6022643](https://pubmed.ncbi.nlm.nih.gov/PMC6022643/)

**38. Microbial genomic island discovery, visualization and analysis**

Claire Bertelli, Keith E Tilley, Fiona SL Brinkman  
*Briefings in Bioinformatics* (2018-06-03) <https://doi.org/gdnhfv>  
DOI: [10.1093/bib/bby042](https://doi.org/10.1093/bib/bby042) · PMID: [29868902](https://pubmed.ncbi.nlm.nih.gov/29868902/)

**39. Understanding the mechanisms and drivers of antimicrobial resistance.**

Alison H Holmes, Luke SP Moore, Arnfinn Sundsfjord, Martin Steinbakk, Sadie Regmi, Abhilasha Karkey, Philippe J Guerin, Laura JV Piddock  
*Lancet (London, England)* (2015-11-18) <https://www.ncbi.nlm.nih.gov/pubmed/26603922>  
DOI: [10.1016/s0140-6736\(15\)00473-0](https://doi.org/10.1016/s0140-6736(15)00473-0) · PMID: [26603922](https://pubmed.ncbi.nlm.nih.gov/26603922/)

**40. Multicopy plasmids potentiate the evolution of antibiotic resistance in bacteria**

Alvaro San Millan, Jose Antonio Escudero, Danna R. Gifford, Didier Mazel, R. Craig MacLean  
*Nature Ecology & Evolution* (2016-11-07) <https://doi.org/bs76>  
DOI: [10.1038/s41559-016-0010](https://doi.org/10.1038/s41559-016-0010) · PMID: [28812563](https://pubmed.ncbi.nlm.nih.gov/28812563/)

**41. Small-Plasmid-Mediated Antibiotic Resistance Is Enhanced by Increases in Plasmid Copy Number and Bacterial Fitness**

Alvaro San Millan, Alfonso Santos-Lopez, Rafael Ortega-Huedo, Cristina Bernabe-Balas, Sean P. Kennedy, Bruno Gonzalez-Zorn  
*Antimicrobial Agents and Chemotherapy* (2015-03-30) <https://doi.org/f7k8bk>  
DOI: [10.1128/aac.00235-15](https://doi.org/10.1128/aac.00235-15) · PMID: [25824216](https://pubmed.ncbi.nlm.nih.gov/25824216/) · PMCID: [PMC4432117](https://pubmed.ncbi.nlm.nih.gov/PMC4432117/)

**42. cBar: a computer program to distinguish plasmid-derived from chromosome-derived sequence fragments in metagenomics data**

Fengfeng Zhou, Ying Xu  
*Bioinformatics* (2010-08-02) <https://doi.org/cn7486>  
DOI: [10.1093/bioinformatics/btq299](https://doi.org/10.1093/bioinformatics/btq299) · PMID: [20538725](https://pubmed.ncbi.nlm.nih.gov/20538725/) · PMCID: [PMC2916713](https://pubmed.ncbi.nlm.nih.gov/PMC2916713/)

**43. Modal Codon Usage: Assessing the Typical Codon Usage of a Genome**

J. J. Davis, G. J. Olsen  
*Molecular Biology and Evolution* (2009-12-17) <https://doi.org/bhsmq5>  
DOI: [10.1093/molbev/msp281](https://doi.org/10.1093/molbev/msp281) · PMID: [20018979](https://pubmed.ncbi.nlm.nih.gov/20018979/) · PMCID: [PMC2839124](https://pubmed.ncbi.nlm.nih.gov/PMC2839124/)

**44. IslandViewer 3: more flexible, interactive genomic island discovery, visualization and analysis: Figure 1.**

Bhavjinder K. Dhillon, Matthew R. Laird, Julie A. Shay, Geoffrey L. Winsor, Raymond Lo, Fazmin Nizam, Sheldon K. Pereira, Nicholas Waglechner, Andrew G. McArthur, Morgan G. I. Langille, Fiona S. L. Brinkman  
*Nucleic Acids Research* (2015-04-27) <https://doi.org/f7n2xs>  
DOI: [10.1093/nar/gkv401](https://doi.org/10.1093/nar/gkv401) · PMID: [25916842](https://pubmed.ncbi.nlm.nih.gov/25916842/) · PMCID: [PMC4489224](https://pubmed.ncbi.nlm.nih.gov/PMC4489224/)

**45. Critical Assessment of Metagenome Interpretation—a benchmark of metagenomics software**

Alexander Sczyrba, Peter Hofmann, Peter Belmann, David Koslicki, Stefan Janssen, Johannes Dröge, Ivan Gregor, Stephan Majda, Jessika Fiedler, Eik Dahms, ... Alice C McHardy  
*Nature Methods* (2017-10-02) <https://doi.org/gbzpspt>  
DOI: [10.1038/nmeth.4458](https://doi.org/10.1038/nmeth.4458) · PMID: [28967888](https://pubmed.ncbi.nlm.nih.gov/28967888/) · PMCID: [PMC5903868](https://pubmed.ncbi.nlm.nih.gov/PMC5903868/)

**46. ART: a next-generation sequencing read simulator**

Weichun Huang, Leping Li, Jason R. Myers, Gabor T. Marth  
*Bioinformatics* (2011-12-23) <https://doi.org/fzf84c>  
DOI: [10.1093/bioinformatics/btr708](https://doi.org/10.1093/bioinformatics/btr708) · PMID: [22199392](https://pubmed.ncbi.nlm.nih.gov/22199392/) · PMCID: [PMC3278762](https://pubmed.ncbi.nlm.nih.gov/PMC3278762/)

**47. Sickle: A sliding-window, adaptive, quality-based trimming tool for FastQ files**

NA Joshi, JN Fass  
*GitHub* (2011) <https://github.com/najoshi/sickle>

**48. MetaQUAST: evaluation of metagenome assemblies**

Alla Mikheenko, Vladislav Saveliev, Alexey Gurevich  
*Bioinformatics* (2015-11-26) <https://doi.org/f8jdjj>  
DOI: [10.1093/bioinformatics/btv697](https://doi.org/10.1093/bioinformatics/btv697) · PMID: [26614127](https://pubmed.ncbi.nlm.nih.gov/26614127/)

**49. BLAST+: architecture and applications**

Christiam Camacho, George Coulouris, Vahram Avagyan, Ning Ma, Jason Papadopoulos, Kevin Bealer, Thomas L Madden  
*BMC Bioinformatics* (2009) <https://doi.org/cnjxgz>  
DOI: [10.1186/1471-2105-10-421](https://doi.org/10.1186/1471-2105-10-421) · PMID: [20003500](https://pubmed.ncbi.nlm.nih.gov/20003500/) · PMCID: [PMC2803857](https://pubmed.ncbi.nlm.nih.gov/PMC2803857/)

**50. CARD 2017: expansion and model-centric curation of the comprehensive antibiotic resistance database**

Baofeng Jia, Amogelang R. Raphenya, Brian Alcock, Nicholas Waglechner, Peiyao Guo, Kara K. Tsang, Briony A. Lago, Biren M. Dave, Sheldon Pereira, Arjun N. Sharma, ... Andrew G. McArthur  
*Nucleic Acids Research* (2016-10-26) <https://doi.org/f9wbjs>  
DOI: [10.1093/nar/gkw1004](https://doi.org/10.1093/nar/gkw1004) · PMID: [27789705](https://pubmed.ncbi.nlm.nih.gov/27789705/) · PMCID: [PMC5210516](https://pubmed.ncbi.nlm.nih.gov/PMC5210516/)

**51. VFDB 2019: a comparative pathogenomic platform with an interactive web interface**

Bo Liu, Dandan Zheng, Qi Jin, Lihong Chen, Jian Yang  
*Nucleic Acids Research* (2018-11-05) <https://doi.org/gf4zfr>  
DOI: [10.1093/nar/gky1080](https://doi.org/10.1093/nar/gky1080) · PMID: [30395255](https://pubmed.ncbi.nlm.nih.gov/30395255/) · PMCID: [PMC6324032](https://pubmed.ncbi.nlm.nih.gov/PMC6324032/)

**52. Prodigal: prokaryotic gene recognition and translation initiation site identification**

Doug Hyatt, Gwo-Liang Chen, Philip F LoCascio, Miriam L Land, Frank W Larimer, Loren J Hauser

*BMC Bioinformatics* (2010-03-08) <https://doi.org/cktxnm>

DOI: [10.1186/1471-2105-11-119](https://doi.org/10.1186/1471-2105-11-119) · PMID: [20211023](https://pubmed.ncbi.nlm.nih.gov/20211023/) · PMCID: [PMC2848648](https://pubmed.ncbi.nlm.nih.gov/PMC2848648/)

**53. PSORTb 3.0: improved protein subcellular localization prediction with refined localization subcategories and predictive capabilities for all prokaryotes**

Nancy Y. Yu, James R. Wagner, Matthew R. Laird, Gabor Melli, Sébastien Rey, Raymond Lo, Phuong Dao, S. Cenk Sahinalp, Martin Ester, Leonard J. Foster, Fiona S. L. Brinkman

*Bioinformatics* (2010-05-13) <https://doi.org/bz3q2w>

DOI: [10.1093/bioinformatics/btq249](https://doi.org/10.1093/bioinformatics/btq249) · PMID: [20472543](https://pubmed.ncbi.nlm.nih.gov/20472543/) · PMCID: [PMC2887053](https://pubmed.ncbi.nlm.nih.gov/PMC2887053/)

**54. On the (im)possibility of reconstructing plasmids from whole-genome short-read sequencing data**

Sergio Arredondo-Alonso, Rob J. Willems, Willem van Schaik, Anita C. Schürch

*Microbial Genomics* (2017-10-01) <https://doi.org/gf6b63>

DOI: [10.1099/mgen.0.000128](https://doi.org/10.1099/mgen.0.000128) · PMID: [29177087](https://pubmed.ncbi.nlm.nih.gov/29177087/) · PMCID: [PMC5695206](https://pubmed.ncbi.nlm.nih.gov/PMC5695206/)