

Biases in genome reconstruction from metagenomic data

William C Nelson^{Corresp., 1}, Benjamin J Tully^{2, 3}, Jennifer M Mobberley⁴

¹ Biological Sciences Division, Pacific Northwest National Laboratory, Richland, Washington, USA

² Department of Biological Sciences, Marine Environmental Biology Section, University of Southern California, Los Angeles, CA, United States

³ Center for Dark Energy Biosphere Investigations, University of Southern California, Los Angeles, California, USA

⁴ Chemical and Biological Signature Science Group, Pacific Northwest National Laboratory, Richland, Washington, USA

Corresponding Author: William C Nelson

Email address: william.nelson@pnnl.gov

Background: Advances in sequencing, assembly, and assortment of contigs into species-specific bins has enabled the reconstruction of genomes from metagenomic data (MAGs). Though a powerful technique, it is difficult to determine whether assembly and binning techniques are accurate when applied to environmental metagenomes due to a lack of complete reference genome sequences against which to check the resulting MAGs.

Methods: We compared MAGs derived from an enrichment culture containing ~20 organisms to complete genome sequences of 10 organisms isolated from the enrichment culture. Factors commonly considered in binning software - nucleotide composition and sequence repetitiveness - were calculated for both the correctly binned and not-binned regions. This direct comparison revealed biases in sequence characteristics and gene content in the not-binned regions. Additionally, the composition of three public data sets representing MAGs reconstructed from the *Tara* Oceans metagenomic data was compared to a set of representative genomes available through NCBI RefSeq to verify that the biases identified were observable in more complex data sets and using three contemporary binning software packages.

Results: Repeat sequences were frequently not binned in the genome reconstruction processes, as were sequence regions with variant nucleotide composition. Genes encoded on the not-binned regions were strongly biased towards ribosomal RNAs, transfer RNAs, mobile element functions and genes of unknown function. Our results support genome reconstruction as a robust process and suggest that reconstructions determined to be >90% complete are likely to effectively represent organismal function, however, population-level genotypic heterogeneity in natural populations, such as uneven distribution of plasmids, can lead to incorrect inferences.

Biases in Genome Reconstruction from Metagenomic Data

William C. Nelson¹, Benjamin J. Tully^{2,3} and Jennifer M. Mobberley⁴

¹ Biological Sciences Division, Pacific Northwest National Laboratory, Richland, WA, USA

² Department of Biological Sciences, Marine and Environmental Biology Section, University of Southern California, Los Angeles, CA, USA

³ Center for Dark Energy Biosphere Investigations, University of Southern California, Los Angeles, CA, USA

⁴ Chemical and Biological Signature Science Group, Pacific Northwest National Laboratory, Richland, WA, USA

Corresponding Author:

William C. Nelson¹

Pacific Northwest National Laboratory, PO Box 999, MSIN J4-18, Richland, WA 99352

william.nelson@pnnl.gov

ABSTRACT

Background: Advances in sequencing, assembly, and assortment of contigs into species-specific bins has enabled the reconstruction of genomes from metagenomic data (MAGs). Though a powerful technique, it is difficult to determine whether assembly and binning techniques are accurate when applied to environmental metagenomes due to a lack of complete reference genome sequences against which to check the resulting MAGs.

Methods: We compared MAGs derived from an enrichment culture containing ~20 organisms to complete genome sequences of 10 organisms isolated from the enrichment culture. Factors commonly considered in binning software - nucleotide composition and sequence repetitiveness - were calculated for both the correctly binned and not-binned regions. This direct comparison revealed biases in sequence characteristics and gene content in the not-binned regions. Additionally, the composition of three public data sets representing MAGs reconstructed from the *Tara* Oceans metagenomic data was compared to a set of representative genomes available through NCBI RefSeq to verify that the biases identified were observable in more complex data sets and using three contemporary binning software packages.

Results: Repeat sequences were frequently not binned in the genome reconstruction processes, as were sequence regions with variant nucleotide composition. Genes encoded on the not-binned regions were strongly biased towards ribosomal RNAs, transfer RNAs, mobile element functions and genes of unknown function. Our results support genome reconstruction as a robust process and suggest that reconstructions determined to be >90% complete are likely to effectively represent organismal function, however, population-level genotypic heterogeneity in natural populations, such as uneven distribution of plasmids, can lead to incorrect inferences.

INTRODUCTION

High-throughput sequencing has revolutionized microbiology by circumventing “the great plate count anomaly” (1) and allowing direct investigation of natural communities in a culture-independent manner (2–8). One goal of metagenomics has always been to obtain organism-specific, complete, genomic information from the complex mixture of sequence data generated from environmental samples. Having a complete genome sequence provides a platform for understanding the range of metabolic roles an organism can play within a community and the interactions it has with other organisms (9–11), and it can provide specific context for interpretation of transcriptomics and proteomics (12,13). Metagenome-assembled genomes (MAGs) are produced by segregating assembled contigs/scaffolds into organism-specific ‘bins’. This process of genome reconstruction has benefitted from continuing advances in sequencing technologies, sequence assembly algorithms, and segregation methods (14). Early success assembling genomes from a simple community (15) has led to more recent studies reconstructing many organisms from complex environments (16–30). The accuracy of these techniques in the context of a complex environmental community is difficult to gauge, however, because most available complete microbial genome sequences that could serve as references are from cultured isolates, and these isolates are rarely present in environmental metagenomes. Techniques that have been developed to evaluate the accuracy of the binning process rely on conserved genes and consistency of nucleotide composition (31–35). These techniques, however, cannot make accurate determinations of how much sequence is missing or the functional potential of missing content. Genome reconstruction techniques have been tested using synthetic communities of cultured organisms (36) and simulated metagenomic datasets. Over time, increasingly sophisticated methods have been developed to simulate metagenomic read data sets, from the

earlier Grinder (37), MetaSim (38), GemSIM (39), BEAR (40), and NeSSM (41), to the more recent CAMISIM (42), which was developed as part of the community effort to address standards in metagenome analysis software development (43). Generally these simulators concern themselves with modeling community structure and sequencing attributes, such as read length and error rates, but are limited to presenting data generated from a reference genomic database, thus cannot model the genetic diversity found in most environments, although CAMISIM addresses this issue by implementing the genome evolution simulator sgEvolver (44). Because genetic variability within natural populations is, as yet, ill-defined (45), it is unlikely that such test data can accurately replicate the type and amount of variability found in natural communities, and the complications this variability causes.

Unicyanobacterial consortia (UCC) were developed as model systems to investigate the mechanisms of metabolic interaction between cyanobacteria and heterotrophs. These systems provide an opportunity to compare MAGs against a matching reference genome set and learn about potential gaps and pitfalls of current reconstruction processes. Two consortia, each containing a single unique cyanobacterial species and sharing an additional 18 heterotrophic species, were derived from a natural mat community (46). The communities have been sequenced, and genome reconstruction has been performed (47), yielding near-complete genome sequences revealing the presence and maintenance of microdiversity, such as might be found within an intact environmental sample. Thus, this system more accurately reflects *in situ* community diversity compared to synthetic communities constructed from isolated organisms. In parallel, isolates of 10 of the member species have also been sequenced (47,48). This paired genomic and metagenomic data set allows direct comparison of MAGs from diverse organisms against ‘ground truth’ genomic data. Previously, we have shown that common aspects of the

genome reconstruction process (assembly from a complex sequence space and segregation of contigs based on read depth profiles and sequence composition) to be both specific and sensitive (47).

We have investigated the nature of genomic regions that under current standard genome reconstruction techniques are not recovered (herein referred to as **not-binned regions**, or **NRs**) to evaluate how these regions differ from recovered regions (**correctly binned regions**, or **CRs**), and to what extent the missing genomic information might impact conclusions drawn from analysis of MAGs. Two common elements of current sequence segregation protocols are analysis of sequence composition and comparison of coverage profiles between samples, so we compared the nucleotide content of NRs vs CRs, examining both %G+C and tetranucleotide content, and the redundancy of sequence information both within the individual genome (*i.e.*, repetitiveness within the genome) and across the entire metagenomic data set (*i.e.*, sequence shared between populations). To determine the impact on downstream functional analyses, the gene content was examined for biases in the cellular roles of genes found within NRs and CRs. To verify that the biases observed extended to more complex metagenomic datasets and across binning algorithms, the *Tara* Oceans metagenome, which has been binned by different groups using MetaBAT (22,49), Anvi'o (31,50), and BinSanity (21,51), was subjected to similar sequence and repeat compositional analysis.

MATERIALS & METHODS

Data and Code Availability.

The UCC MAG and genome data analyzed are available in the GenBank repository as listed in Table 1. The metagenomic data used to construct the UCC MAGs is available from the NCBI

SRA (accessions SRX1063989 and SRX1065184). MAGs reconstructed from the *Tara* Oceans metagenomic data (21,22) are available in the GenBank repository. MAGs from Delmont et al. (50) are available through figshare (doi: 10.6084/m9.figshare.4902923). A list of MAGs and corresponding identifiers are available in Supplemental Table 1. Complete bacterial and archaeal genomes were collected from NCBI RefSeq (52) (accessed Aug 2019) based on assignment as either “reference genome” or “representative genome” for the data column “refseq_category” and “Complete Genome” in the “assembly_level” column. A list of genomes used in the analysis are available in Supplemental Table 2. All analysis scripts are available at http://github.com/wichne/biases_in_genome_reconstruction.

data

data

data

Identification of CR and NR regions.

The UCC scaffolds comprising each MAG were searched against their cognate complete genome sequence using nucmer using the maxmatch option (53). Regions of the genomes that aligned end-to-end to MAG scaffolds at $\geq 99\%$ identity were cataloged as CR regions. All other genome regions were considered NR regions.

Compositional analysis.

For the UCC MAGs and genomes, %G+C calculation and tetranucleotide frequency (TNF) chi-square test were performed using custom Perl scripts (available at

http://github.com/wichne/biases_in_genome_reconstruction). Compositional analysis was

data

restricted to CR or NR regions longer than 1000 bp to ensure sufficient sequence for meaningful results. For TNF, the chi-squared statistic was calculated for each region using the TNF for the whole genome as the expected values, and the mean and standard deviation for the CR and NR pools calculated. For %G+C analysis, the mean %G+C for the CR and NR regions was calculated, and the absolute difference was calculated between each region and the genome

average, and average differences determined for CR and NR pools. To estimate *p-values* for the %G+C and TNF analyses, one thousand random coordinate sets yielding the same number and length of fragments as in each genome's CR or NR set were generated from the genome sequence and evaluated.

For comparison of the UCC data set to the *Tara* Oceans MAGs and RefSeq genome data sets, sequence composition variance (i.e., deviation from the mean) was calculated for the %G+C and tetranucleotide frequency using a custom Python script. The %G+C was calculated for 2kb segments (sliding window of 500bp) for each MAG or genome. A genome-wide variance value was calculated for each MAG or genome based on the segments and plotted as a box plot per source data set. TNF was calculated for 10kb segments (sliding window 5kb) for each MAG or genome. Using the calculation described in Teeling (54), each segment had a Z-score calculated for each tetranucleotide based on the observed-vs-expected frequency of the tetranucleotide in the 10kb segment. A Pearson correlation was then calculated in a pairwise fashion for all segments. Variance of the Pearson correlation values within a MAG or genome was calculated and plotted as a box plot per source data set.

Repetitiveness analysis

To calculate intragenome sequence repetitiveness, we determined the fraction of each genome that was comprised of repeat sequence. Each genome sequence was searched against itself using nucmer v3.0 (53) with the maxmatch option, and the lengths of regions that aligned to another part of the genome/MAG with $\geq 97\%$ identity were summed and divided by the length of the genome/MAG.

To determine the repetitiveness of sequences across the entire metagenomic data set, metagenome reads were searched against genome sequences using Bowtie2 (55). Per-base

coverage was calculated using the samtools (56) depth command, and average coverage values for the genomes, NRs and CRs were determined. One thousand sets of random coordinate regions of the same number and lengths as in each set were analyzed to estimate p-values. Results are reported as average coverage depth of NRs and CRs and the average difference from the genome depth-of-coverage.

Gene function analysis

UCC complete genome sequences were annotated by the IMG pipeline (57), which included COG assignment based on the December 2014 release of the 2003-2014 COGs (58). COGs assigned to more than one functional category were counted for each assigned category. Genes not assigned to a COG category were classified as ‘unassigned’. Ribosomal RNA (rRNA) gene features were identified by the IMG pipeline (59); transfer RNAs (tRNA) were identified with tRNAscan-SE (60); other non-coding RNAs (ncRNA) were identified using the Rfam database v11.0 (61) and infernal v1.1 software (62). For each gene set, the category counts were normalized to the total feature counts. Principle component analysis was performed and biplot of gene categories was generated using R package bpca v.1.2-2 (<http://cran.r-project.org/web/packages/bpca/>).

Statistical analysis.

Statistical tests were performed using modules within the Python package SciPy (63). The normality of the calculated variance distributions for each set of genomes was determined using the Shapiro-Wilk test (64). Genome sets with a normal distribution were compared to each other with the T-test for two independent variables (65). Genome sets without a normal distribution were compared to each other with the Mann-Whitney U test (66). p-values were adjusted for

data

multiple comparisons with the Benjamini-Hochberg procedure (67) correction with a false discovery rate of 25% (**Supplemental Table 3**).

RESULTS AND DISCUSSION

The power of metagenomics is that it allows exploration of diverse communities from which we cannot culture the component populations either because the proper growth conditions are unknown or difficult to replicate in a laboratory environment, or simply because there are too many organisms present to have the resources or time to pursue the effort. Because of this, there are very few examples of sequenced organisms isolated from the same sample from which metagenomic sequencing and binning has been done to generate MAGs. As such, a ‘gold standard’ for evaluation of MAG content has been difficult to come by. We have taken advantage of two enrichment cultures from which MAGs and isolate genomes have been derived to generate just such a ‘gold standard’ comparison framework. We have previously generated two uncyanobacterial consortial cultures (UCC) – enrichment cultures each containing a distinct cyanobacterial population and different, yet overlapping, communities of associated heterotrophs, each numbering <20 species – and performed metagenomic sequencing, assembly and binning.(47,48). Illumina 150 bp paired-end reads were generated from each community, and IDBA_ud was used to assemble the read sets separately and in co-assembly. The abundances of the organisms differed between the two communities, allowing us to bin the sequences by comparing sequence coverage values of contigs between the two UCCs in a predominantly manual process (inspired by the work of Dick, et al (68)). The resulting MAGs were manually curated to eliminate contaminating contigs and identify mis-binned contigs, correctly placing them when possible. In parallel, ten organisms were isolated from the UCCs and completely

sequenced. Comparison of the MAGs to the isolate genomes showed recovery of >90% of sequence for genomes with at least 10x coverage, with one exception, *Halomonas* sp. HL-93, which had 85% recovery from 11x coverage (**Table 1**). Co-linear sequence alignments indicated there were no assembly errors in the binned contigs (47, and data not shown). Based on the isolate-MAG comparisons, NRs were identified. *Porphyrobacter* HL-46 had the lowest metagenome coverage (3.6x). Its MAG comprised hundreds of short contigs and was determined to be ~40% complete. Thus, the NRs for HL-46 are assumed to be primarily caused by the random sampling of the shotgun sequencing methodology and not by any inherent content biases, allowing the HL-46 analyses to serve as a control.

To determine if NRs were not binned due to lack of assembly, we mapped the contigs from the assembly to the CR and NR regions of the genomes and looked at the contig coverage of the regions. As expected, the CRs showed an average contig coverage of 1.04 ± 0.14 , and most regions had only a single contig map to them (Fig S1). Many of the cases of multiple contigs mapping to a CR were due to short (<200 bp) contigs of repeat sequence which might be an artifact of the assembler (IDBA_ud). NRs show a strong positive correlation between region length and number of contigs mapping, with an average coverage of 0.94 ± 0.71 (Fig S2). This suggests poorer assembly of the NRs and higher repeat content, but also indicates that most NR sequence is present in the contig set, and thus the binning process is the main determinant of NRs.

Nucleotide composition of NRs frequently differs from the genome average

Bacteria and Archaea have evolved to have a fairly consistent %G+C across their genome (69), so much so that it has been proposed as a metric of classification at higher taxonomic levels (70).

It is not uncommon, however, to observe regions within a genome that differ significantly from the genome average (71). This variation can be the result of selective pressure for structural properties in non-coding genes, for instance ribosomal RNAs and other functional RNAs have been shown to vary in nucleotide composition in correlation with optimal growth temperature (72). In other cases, divergent %G+C indicates a region which has been acquired recently (in evolutionary time) from a non-related source (*i.e.*, horizontal gene transfer) (73). To investigate whether variant G+C confounds genome reconstruction, we compared the %G+C of NRs to that of CRs and the complete genome.

The genomes in this study had a range of %G+C values, from 42% (*A. marincola* HL-49) to 68% (*Erythrobacteraceae* bacterium HL-111), with most skewing toward the higher values (**Table 2**). We determined the %G+C for each CR and NR ≥ 200 bp in length and compared them to the %G+C for the complete genome. For genomes with more than one genomic element, each molecule was considered separately since extrachromosomal elements may have distinct nucleotide composition. For seven of the genomes, the %G+C for the NRs differed significantly ($p \leq 0.005$) from the genome average, while the CRs generally reflected the genome average (**Table 2**). The %G+C averages for NRs from HL-48 and HL-111 were significantly lower (45.76% and 64.26%, respectively) than the genomes' averages (58.98% and 68.12% respectively). Other genomes (HL-53, HL-55, HL-109) had some NRs with %G+C higher than the genome average and some NRs with lower values (**Figure 1**), despite having different average %G+C values (47.5%, 56.0% and 64.1% respectively). Extrachromosomal elements analyzed did not display a significant difference in the %G+C of their NRs from the molecule average. As expected, the values for the NRs and CRs of HL-46 showed no significant difference from the genome average (**Table 2**), however, HL-46's CRs and NRs did not display identical

%G+C profiles (**Figure 1**). There was a slight bias toward higher %G+C for the NRs and lower %G+C in the CRs, which could reflect a bias in the assembly algorithm.

Tetranucleotide frequency (TNF) has been shown to be capable of distinguishing higher taxonomic classifications, up to the species level (54,68). This resolving power has been leveraged in binning protocols (15,74–76). To investigate whether genomic regions with divergent TNF are poorly recovered in genome reconstruction, we compared the TNFs of CRs and NRs to that of the cognate complete genome using chi-squared analysis. In most cases, the chi-squared statistic was an order of magnitude higher for NRs versus CRs, and the differences were significant for all chromosomal sequences except for HL-46, HL-109, HL-93 and the small chromosome of HL-91 (**Table 3**).

One factor that could affect nucleotide composition effects on binning is the length of the region with divergent composition versus the length of the contig. If the variant region comprises most of the length of the contig being evaluated, the difference from the genome average will be pronounced, whereas if the divergent region is only a small percentage of the contig length, the signal will be muted. An examination of CR/NR length versus compositional variance (**Fig. S3**) revealed a strong, significant negative correlation between contig length and TNF chi square for CRs ($R^2=0.64$, $p\text{-value}<2.2\times10^{-16}$) and a weaker relationship for NRs ($R^2=0.14$, $p\text{-value}=4.9\times10^{-12}$). Taken together, the %G+C and TNF results show that genomic regions with divergent nucleotide composition are more likely to be missed during binning, and this effect is stronger for short contigs. The most effective way to overcome this problem is to enhance assembly such that regions with unusual content are included in significantly longer contigs, or, through clone linkage, identify strong, unique connections to binned contigs.

Repeated sequences segregate aberrantly

Sequence coverage profiles are frequently effective in discriminating contigs from different organisms (15). Samples taken under different conditions or at different times capture community states which have similar organismal composition but differing relative abundances. This difference translates to distinct coverage profiles for assembled contigs, and thus contigs with similar coverage profiles are assumed to originate from the same organism. In this data set, for example, we compared two cultures with near-identical heterotroph species composition, but different cyanobacteria acting as a conduit for energy and carbon (46,47). Other studies have compared samples taken at different times (75). Coverage analysis is more difficult for repeated regions of a genome, which will yield higher coverage values than the genome average and thus are more likely to be either not binned or binned improperly. Differential coverage analysis can mitigate this problem by identifying correlated changes in abundance of contigs with different coverage. Unlike nucleotide composition variance, however, unusual high-coverage signal due to repeat sequence is less likely to be diluted by incorporation into a larger contig because assemblers (especially standard de Bruijn graph assemblers using short-read data) tend to terminate contigs when repeats are encountered and/or assemble repeats into separate contigs (77).

To examine the impact of repeated sequences on genome reconstruction, we determined the repetitiveness of sequence information across CRs and NRs, determined from a self-versus-self similarity search, and compared those values to the genome average. Correspondence of repeated regions and NRs was strong (**Figures 2 and 3, Figure S4**). In HL-111, all NRs save one were present in at least two copies (**Figure 2**). For all reconstructions, save HL-46, the CRs had repeat content equal to or lower than the genome average.

Another phenomenon that can affect contig coverage in metagenomic assembly is multiple organisms sharing identical regions of DNA. Some regions are highly conserved between related species, an example being the ribosomal RNA operon, which is known to confound assemblers and segregation strategies (78). Alternatively, mobile elements such as plasmids or transposons can have a broad host range and invade and inhabit closely or even distantly related organisms (79). Such regions, even if not repeated within a genome, will exhibit anomalous coverage and thus could be either excluded or mis-binned. We examined the metagenomic read coverage depth to determine if NRs had anomalous profiles relative to the whole genome and the CRs. For most reconstructions, the NRs' coverage differed from the genome average and that of the CRs (**Table 4, Fig 2, Fig S4**). Only HL-46 and one of the HL-109 molecules did not have significant differences. Most NRs displayed higher or equivalent coverage values, however, several NRs in HL-48 and the two small plasmids associated with HL-91 showed lower metagenomic coverage values (**Figure S4**). A likely explanation for this is the presence in the consortia of sub-populations of these organisms that lack the plasmids.

Functional assessment of NR genes

To determine the extent to which regions missing from reconstructions might affect downstream metabolic or functional analyses and predictions for organisms and communities, we examined the gene content of the NRs and the functional roles of those genes. COG categorization was used as a basis for comparison because of its ability to identify, in particular, genes associated with mobile elements such as plasmids, phage and insertion sequences. In addition, we evaluated the distribution of non-coding RNA genes since some are known to be repeated within genomes

(multiple rRNA operons, for example), and others (tRNAs) are commonly associated with mobile elements (80).

For all the reconstructions, the gene content of the NRs differed from that of the CRs and complete genomes. Functional analysis of gene sequences shows that this difference was largely driven by genes encoding mobile element functions (COG category X) and RNA genes (**Figure 4**). The mobile element genes in the NR regions were predominantly transposases with some contribution from bacteriophage and plasmid genes (HL-91; HL-93). Most of the identified rRNA genes fell within NRs, with only HL-48 and HL-53 each having one rRNA contained in a CR. In addition, the NRs, including the two entire plasmids from HL-91 which were not binned, contained a higher percentage of genes that were not assigned to a COG category.

Evaluation of a complex metagenomic data set and common automated binning tools

To verify that our conclusions of genome reconstruction bias in the highly curated UCC data set were extendable to more complex data sets and for alternate, widely-used binning tools, we applied similar analyses to MAGs generated from the *Tara* Oceans metagenomic data using distinct genome reconstruction protocols. For this comparison, 4,557 MAGs generated from the *Tara* Oceans microbial metagenomic data reconstructed using three complementary methods were collected and analyzed. Three different automated binning methodologies were employed to generate the MAG data set: MetaBat (v0.26.3) (22,49), BinSanity (v1.0) (21,51), and CONCOCT (with manual refinement in anvi'o) (31,50). All three automated binning algorithms utilized read coverage and TNF to identify congruent contigs, with the intended role of the algorithms to reconstruct high confidence environmental genomes while avoiding over-binning (*i.e.*, removing elements that deviate from the mean values of the binned contigs). The MAGs

had a mean estimated completeness and contamination of 76.6% and 2.2%, respectively, as determined by CheckM v.1.1.1 (32). In comparison, 1,736 ‘representative’ and ‘reference’ complete genomes were collected from NCBI RefSeq.

Our results above predicted that the MAGs would have lower %G+C variance and TNF variance than the isolate complete genome data set. For the observed %G+C, MAGs tended to have lower variance ($p < 0.001$) than isolate genomes (**Figure 5A**). The exception was the Parks *et al.* MAGs, which had a much larger variance, even compared to the RefSeq genome set (mean vs mean, $p < 0.001$). This may be the result of the additional step applied to the MAGs by Parks *et al.*, whereby related MAGs with <3% mean %G+C difference were merged into a single representative MAG (22). For the Tully *et al.* and Delmont *et al.* MAGs, the lower variance observed compared to the RefSeq genomes is likely due to removal of contigs with deviant %G+C values during binning (21,50). The MAGs also had lower variance with regards to TNF compared to the RefSeq genomes ($p < 0.001$) (**Figure 5B**), again, likely due to genomic elements that deviated from the average value of the binned contigs having been removed during the binning steps. These observations support our conclusions regarding genome regions having divergent nucleotide composition being underrepresented in MAGs.

The Tara and NCBI Refseq data sets were then evaluated for repeat sequence content. Each MAG and isolate genome was compared to itself using NUCmer to identify the fraction of the genome composed of repeat regions (regions with $\geq 97\%$ sequence identity). MAGs universally had a smaller fraction of genomic information in repeat regions compared to isolate genomes ($p < 0.01$; **Figure 6**). The lack of repeat regions in MAGs is likely the result of repeated regions having inflated or depressed read coverage values relative to the mean of the genome, depending on the number of copies of the repeat region present in the genome and how stable

this number is across the population. Compared to the other *Tara* MAGs, the Tully *et al.* MAGs had a larger fraction of redundant genomic elements. It is unclear what aspect of the assembly and binning methodology has influenced these results. On average, the lengths of the repeat regions from the Tully *et al.* MAGs are longer than the repeat regions in the RefSeq genomes (mean: 1,052bp vs 868bp, respectively).

What's missing from reconstructed genomes?

Analysis of regions that were not recovered from genome reconstruction (NRs) showed both nucleotide compositional variance and intragenome repetitiveness. The %G+C and tetranucleotide frequencies of NRs tended to differ from that of complete genomes (**Tables 2 and 3, Figure 1**), and the sequence coverage differed. This met expectations since, in general, binning tools are designed around the assumption that sequences with similar properties belong together, thus any genome region that varies significantly from the genome average is likely to be incorrectly binned if it comprises the majority of a contig under consideration. Regions with atypical nucleotide content have been observed to contain genes upon which selective pressures are acting on nucleic acid structure, such as ribosomal RNAs and tRNAs (72,81,82), and exogenously introduced segments such as mobile elements (83,84). It is significant that many of the NRs displayed lower %G+C than the genome average, since it has been observed that laterally acquired regions tend to have lower %G+C than their hosts (83), as phage and insertion sequences tend to have A+T-enriched genomes (85). Notably, many genome regions with variant nucleotide composition were incorporated into longer contigs by the assembler, masking the variance and allowing correct binning. Conversely, the assembler collapsed repeated region sequences into single contigs, and thus they were not binned due to the inflated sequence

coverage values. Often, repeated sequences displayed divergent nucleotide composition, but the reciprocal was less frequent, indicating that repetitiveness is the stronger driver of binning failure. These results demonstrate that assembly efficiency is an important determining factor for correct binning, or conversely, any factor that results in shorter assemblies will result in poorer recovery of anomalous regions. Thus, it is advisable to include replication and positive controls in metagenomic sequencing protocols, particularly for highly diverse communities such as soils and riverbed sediments, to allow evaluation of assembly efficiency and accuracy.

Repeat regions identified in this study appeared to largely consist of insertion elements based on functional analysis and their relatively short size (1-2 kb). Failure of these regions to be correctly binned is unlikely to meaningfully affect functional predictions for a reconstructed genome. Their presence in a genome is more likely to affect metabolic reconstruction analysis by reducing assembly efficiency, resulting in more, shorter contigs and increasing the chance that these shorter contigs are not binned or incorrectly binned. Technological advances increasing read length beyond 2 kb will increase contig lengths, binning accuracy, and the likelihood of yielding closed genomes from environmental samples (8,86,87).

NRs were generally observed to be short, with a median length of less than 5 kb (**Table 1**) and containing only a handful of genes. Thus, even a MAG with many gaps (indicating a large number of NRs) may be missing only a small percentage of its genome. The conserved single-copy gene (CSCG) estimations for completeness appear for all intents and purposes to be a reasonable indication of how much information is absent (47). One caveat to this conclusion, however, is that extrachromosomal elements, plasmids and phages (integrated or otherwise) typically do not carry CSCG markers, and thus are essentially invisible in such analyses. The longer NRs observed in our analysis appear to comprise integrated plasmids or phage, and thus

any gap in a reconstruction could represent up to 50 kb (or more) of genetic material. Importantly, these represent introduced genetic material, which, while likely conveying a beneficial trait, are unlikely to carry functions that are integral to host metabolic function.

CONCLUSIONS

This analysis indicates that reconstructed genomes estimated to be near-complete can be assumed to contain nearly all genes important to metabolic reconstruction. The majority of identifiable genes present on NRs appear to be either highly conserved, non-coding genes that can be assumed to be present (such as the rRNA genes and tRNA genes) or are associated with mobile genetic elements. While many of these genes may be not be directly related to cellular metabolism (transposases, toxin/antitoxin systems, phage and plasmid functions), it should be noted that entire extrachromosomal elements may be missed by the binning process due to either alternate nucleotide composition, a higher number of copies per cell than the genome, or occupancy in only a subset of the population (such as the two molecules in HL-109). These elements frequently carry genes that alter the physiology or resistance of the host organism. For example, HL-109 and HL-111 have NRs that includes genes involved in glycan biosynthesis, suggesting alterations to the cell wall, while HL-91 has picked up a multidrug efflux transporter. As such, reconstructed genomes can be considered reliable foundations for metabolic reconstruction but should not be assumed to be comprehensive for the function of the organism.

ACKNOWLEDGEMENTS

The authors would like to thank Jim Fredrickson and Lori Nelson for critical evaluation of the manuscript during preparation. W.C.N. and J.M.M. were funded through the U.S. Department of

Energy (DOE) Genome Sciences Program (GSP), Office of Biological and Environmental Research (OBER), and this work was a contribution of the Pacific Northwest National Laboratory (PNNL) Foundational Scientific Focus Area. B.J.T. was funded through the Center for Dark Energy Biosphere Investigations (OCE- 0939654). This is C-DEBI Contribution XXX.

REFERENCES

1. Staley JT, Konopka A. Measurement of in Situ Activities of Nonphotosynthetic Microorganisms in Aquatic and Terrestrial Habitats. *Annu Rev Microbiol.* 1985;39(1):321–46.
2. Venter JC, Remington K, Heidelberg JF, Halpern AL, Rusch D, Eisen JA, et al. Environmental genome shotgun sequencing of the Sargasso Sea. *Science.* 2004 Apr 2;304(5667):66–74.
3. DeLong EF, Preston CM, Mincer T, Rich V, Hallam SJ, Frigaard N-U, et al. Community Genomics Among Stratified Microbial Assemblages in the Ocean’s Interior. *Science.* 2006 Jan 27;311(5760):496–503.
4. Costello EK, Lauber CL, Hamady M, Fierer N, Gordon JI, Knight R. Bacterial Community Variation in Human Body Habitats Across Space and Time. *Science.* 2009 Dec 18;326(5960):1694–7.
5. Caporaso JG, Lauber CL, Walters WA, Berg-Lyons D, Huntley J, Fierer N, et al. Ultra-high-throughput microbial community analysis on the Illumina HiSeq and MiSeq platforms. *ISME J.* 2012 Aug;6(8):1621–4.
6. Zhou J, He Z, Yang Y, Deng Y, Tringe SG, Alvarez-Cohen L. High-throughput metagenomic technologies for complex microbial community analysis: open and closed formats. *MBio.* 2015 Jan 27;6(1).
7. Rinke C, Schwientek P, Sczyrba A, Ivanova NN, Anderson IJ, Cheng J-F, et al. Insights into the phylogeny and coding potential of microbial dark matter. *Nature.* 2013 Jul;499(7459):431–7.
8. White RA, Bottos EM, Roy Chowdhury T, Zucker JD, Brislawn CJ, Nicora CD, et al. Molecule Long-Read Sequencing Facilitates Assembly and Genomic Binning from Complex Soil Metagenomes. *mSystems.* 2016 Jun;1(3).
9. Iverson V, Morris RM, Frazar CD, Berthiaume CT, Morales RL, Armbrust EV. Untangling Genomes from Metagenomes: Revealing an Uncultured Class of Marine Euryarchaeota. *Science.* 2012 Feb 3;335(6068):587–90.
10. Sharon I, Morowitz MJ, Thomas BC, Costello EK, Relman DA, Banfield JF. Time series community genomics analysis reveals rapid shifts in bacterial species, strains, and phage during infant gut colonization. *Genome Res.* 2013 Jan 1;23(1):111–20.

11. Delmont TO, Eren AM, Vineis JH, Post AF. Genome reconstructions indicate the partitioning of ecological functions inside a phytoplankton bloom in the Amundsen Sea, Antarctica. *Front Microbiol* [Internet]. 2015 [cited 2020 Jul 5];6. Available from: <https://www.frontiersin.org/articles/10.3389/fmicb.2015.01090/full>
12. Lesniewski RA, Jain S, Anantharaman K, Schloss PD, Dick GJ. The metatranscriptome of a deep-sea hydrothermal plume is dominated by water column methanotrophs and lithotrophs. *ISME J*. 2012 Dec;6(12):2257–68.
13. Ram RJ, VerBerkmoes NC, Thelen MP, Tyson GW, Baker BJ, Blake RC, et al. Community Proteomics of a Natural Microbial Biofilm. *Science*. 2005 Jun 24;308(5730):1915–20.
14. Sangwan N, Xia F, Gilbert JA. Recovering complete and draft population genomes from metagenome datasets. *Microbiome*. 2016 Mar 8;4:8.
15. Tyson GW, Chapman J, Hugenholtz P, Allen EE, Ram RJ, Richardson PM, et al. Community structure and metabolism through reconstruction of microbial genomes from the environment. *Nature*. 2004 Mar 4;428(6978):37–43.
16. Brown CT, Hug LA, Thomas BC, Sharon I, Castelle CJ, Singh A, et al. Unusual biology across a group comprising more than 15% of domain Bacteria. *Nature*. 2015 Jul 9;523(7559):208–U173.
17. Anantharaman K, Breier JA, Dick GJ. Metagenomic resolution of microbial functions in deep-sea hydrothermal plumes across the Eastern Lau Spreading Center. *ISME J*. 2016 Jan;10(1):225–39.
18. Baker BJ, Lazar CS, Teske AP, Dick GJ. Genomic resolution of linkages in carbon, nitrogen, and sulfur cycling among widespread estuary sediment bacteria. *Microbiome*. 2015;3:14.
19. Li M, Baker BJ, Anantharaman K, Jain S, Breier JA, Dick GJ. Genomic and transcriptomic evidence for scavenging of diverse organic compounds by widespread deep-sea archaea. *Nat Commun*. 2015 Nov 17;6:8933.
20. Nobu MK, Narihiro T, Rinke C, Kamagata Y, Tringe SG, Woyke T, et al. Microbial dark matter ecogenomics reveals complex synergistic networks in a methanogenic bioreactor. *ISME J*. 2015 Aug;9(8):1710–22.
21. Tully BJ, Graham ED, Heidelberg JF. The reconstruction of 2,631 draft metagenome-assembled genomes from the global oceans. *Sci Data*. 2018 Jan 16;5:170203.
22. Parks DH, Rinke C, Chuvochina M, Chaumeil PA, Woodcroft B, Evans PN, et al. Recovery of nearly 8,000 metagenome-assembled genomes substantially expands the tree of life. *Nat Microbiol*. 2018 Feb;3(2):253–253.
23. Pasolli E, Asnicar F, Manara S, Zolfo M, Karcher N, Armanini F, et al. Extensive Unexplored Human Microbiome Diversity Revealed by Over 150,000 Genomes from Metagenomes Spanning Age, Geography, and Lifestyle. *Cell*. 2019 Jan;176(3):649–662.e20.
24. Almeida A, Mitchell AL, Boland M, Forster SC, Gloor GB, Tarkowska A, et al. A new genomic blueprint of the human gut microbiota. *Nature*. 2019 Apr;568(7753):499–504.

25. Mobberley JM, Lindemann SR, Bernstein HC, Moran JJ, Renslow RS, Babauta J, Hu D, Beyenal H, Nelson WC. Organismal and spatial partitioning of energy and macronutrient transformations within a hypersaline mat. *FEMS Microbiol Ecol* [Internet]. 2017 Apr 1 [cited 2020 Sep 2];93(4). Available from: <https://academic.oup.com/femsec/article/doi/10.1093/femsec/fix028/3071443>
26. Stewart RD, Auffret MD, Warr A, Walker AW, Roehe R, Watson M. Compendium of 4,941 rumen metagenome-assembled genomes for rumen microbiome biology and enzyme discovery. *Nat Biotechnol*. 2019;37(8):953–61.
27. Pedron R, Esposito A, Bianconi I, Pasolli E, Tett A, Asnicar F, Cristofolini M, Segata N, Jousson O. Genomic and metagenomic insights into the microbial community of a thermal spring. *Microbiome*. 2019 Dec;7(1):8.
28. Wong HL, White RA, Visscher PT, Charlesworth JC, Vázquez-Campos X, Burns BP. Disentangling the drivers of functional complexity at the metagenomic level in Shark Bay microbial mat microbiomes. *ISME J*. 2018 Nov;12(11):2619–39.
29. Daly RA, Borton MA, Wilkins MJ, Hoyt DW, Kountz DJ, Wolfe RA, Welch SA, Marcus DN, Trexler RV, MacRae JD, Krzycki JA, Cole DR, Mouser PJ, Wrighton KC. Microbial metabolisms in a 2.5-km-deep ecosystem created by hydraulic fracturing in shales. *Nat Microbiol*. 2016 Oct;1(10):16146.
30. Danczak RE, Johnston MD, Kenah C, Slattery M, Wrighton KC, Wilkins MJ. Members of the Candidate Phyla Radiation are functionally differentiated by carbon- and nitrogen-cycling capabilities. *Microbiome*. 2017 Dec;5(1):112.
31. Eren AM, Esen OC, Quince C, Vineis JH, Morrison HG, Sogin ML, et al. Anvi'o: an advanced analysis and visualization platform for 'omics data. *PeerJ*. 2015;3:e1319.
32. Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW. CheckM: assessing the quality of microbial genomes recovered from isolates, single cells, and metagenomes. *Genome Res*. 2015 Jul;25(7):1043–55.
33. Waterhouse RM, Seppey M, Simão FA, Manni M, Ioannidis P, Klioutchnikov G, et al. BUSCO Applications from Quality Assessments to Gene Prediction and Phylogenomics. *Mol Biol Evol*. 2018 Mar 1;35(3):543–8.
34. Chen L-X, Anantharaman K, Shaiber A, Eren AM, Banfield JF. Accurate and Complete Genomes from Metagenomes. *bioRxiv*. 2019 Oct 29;808410.
35. Hugoson E, Lam WT, Guy L. miComplete: weighted quality evaluation of assembled microbial genomes. Hancock J, editor. *Bioinformatics*. 2019 Aug 22;btz664.
36. Hardwick SA, Chen WY, Wong T, Kanakamedala BS, Deveson IW, Ongley SE, et al. Synthetic microbe communities provide internal reference standards for metagenome sequencing and analysis. *Nat Commun*. 2018 Dec;9(1):3096.
37. Angly FE, Willner D, Rohwer F, Hugenholtz P, Tyson GW. Grinder: a versatile amplicon and shotgun sequence simulator. *Nucleic Acids Res*. 2012 Jul 1;40(12):e94–e94.

38. Richter DC, Ott F, Auch AF, Schmid R, Huson DH. MetaSim—A Sequencing Simulator for Genomics and Metagenomics. Field D, editor. PLoS ONE. 2008 Oct 8;3(10):e3373.
39. McElroy KE, Luciani F, Thomas T. GemSIM: general, error-model based simulator of next-generation sequencing data. BMC Genomics. 2012;13(1):74.
40. Johnson S, Trost B, Long JR, Pittet V, Kusalik A. A better sequence-read simulator program for metagenomics. BMC Bioinformatics. 2014 Sep;15(S9):S14.
41. Jia B, Xuan L, Cai K, Hu Z, Ma L, Wei C. NeSSM: A Next-Generation Sequencing Simulator for Metagenomics. Janssen PJ, editor. PLoS ONE. 2013 Oct 4;8(10):e75448.
42. Fritz A, Hofmann P, Majda S, Dahms E, Dröge J, Fiedler J, et al. CAMISIM: simulating metagenomes and microbial communities. Microbiome. 2019 Dec;7(1):17.
43. Sczyrba A, Hofmann P, Belmann P, Koslicki D, Janssen S, Dröge J, et al. Critical Assessment of Metagenome Interpretation—a benchmark of metagenomics software. Nat Methods. 2017 Nov;14(11):1063–71.
44. Darling ACE. Mauve: Multiple Alignment of Conserved Genomic Sequence With Rearrangements. Genome Res. 2004 Jun 14;14(7):1394–403.
45. Rocha EPC. Neutral Theory, Microbial Practice: Challenges in Bacterial Population Genetics. Mol Biol Evol. 2018 Jun 1;35(6):1338–47.
46. Cole JK, Hutchison JR, Renslow RS, Kim YM, Chrisler WB, Engelmann HE, et al. Phototrophic biofilm assembly in microbial-mat-derived unicyanobacterial consortia: model systems for the study of autotroph-heterotroph interactions. Front Microbiol. 2014;5:109.
47. Nelson WC, Maezato Y, Wu YW, Romine MF, Lindemann SR. Identification and Resolution of Microdiversity through Metagenomic Sequencing of Parallel Consortia. Appl Env Microbiol. 2015;82(1):255–67.
48. Romine MF, Rodionov DA, Maezato Y, Osterman AL, Nelson WC. Underlying mechanisms for syntrophic metabolism of essential enzyme cofactors in microbial communities. ISME J. 2017;
49. Kang DD, Froula J, Egan R, Wang Z. MetaBAT, an efficient tool for accurately reconstructing single genomes from complex microbial communities. PeerJ. 2015;3:e1165.
50. Delmont TO, Quince C, Shaiber A, Esen ÖC, Lee ST, Rappé MS, et al. Nitrogen-fixing populations of Planctomycetes and Proteobacteria are abundant in surface ocean metagenomes. Nat Microbiol. 2018 Jul;3(7):804–13.
51. Graham ED, Heidelberg JF, Tully BJ. BinSanity: unsupervised clustering of environmental microbial assemblies using coverage and affinity propagation. PeerJ. 2017 Mar 8;5:e3035.
52. O’Leary NA, Wright MW, Brister JR, Ciufo S, Haddad D, McVeigh R, et al. Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation. Nucleic Acids Res. 2016 Jan 4;44(D1):D733–745.

- 569 53. Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, et al. Versatile and open
570 software for comparing large genomes. *Genome Biol.* 2004;5(2).
- 571 54. Teeling H, Meyerdieks A, Bauer M, Amann R, Glockner FO. Application of tetranucleotide
572 frequencies for the assignment of genomic fragments. *Env Microbiol.* 2004 Sep;6(9):938–47.
- 573 55. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Methods.* 2012
574 Apr;9(4):357-U54.
- 575 56. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence Alignment/Map
576 format and SAMtools. *Bioinformatics.* 2009 Aug 15;25(16):2078–9.
- 577 57. Huntemann M, Ivanova NN, Mavromatis K, Tripp HJ, Paez-Espino D, Palaniappan K, et al. The
578 standard operating procedure of the DOE-JGI Microbial Genome Annotation Pipeline (MGAP v.4).
579 *Stand Genomic Sci.* 2015 Oct 26;10.
- 580 58. Galperin MY, Makarova KS, Wolf YI, Koonin EV. Expanded microbial genome coverage and
581 improved protein family annotation in the COG database. *Nucleic Acids Res.* 2015 Jan
582 28;43(D1):D261–9.
- 583 59. Markowitz VM, Chen IM, Palaniappan K, Chu K, Szeto E, Pillay M, et al. IMG 4 version of the
584 integrated microbial genomes comparative analysis system. *Nucleic Acids Res.* 2014
585 Jan;42(Database issue):D560-7.
- 586 60. Lowe TM, Eddy SR. tRNAscan-SE: A program for improved detection of transfer RNA genes in
587 genomic sequence. *Nucleic Acids Res.* 1997 Mar 1;25(5):955–64.
- 588 61. Burge SW, Daub J, Eberhardt R, Tate J, Barquist L, Nawrocki EP, et al. Rfam 11.0: 10 years of
589 RNA families. *Nucleic Acids Res.* 2013 Jan;41(D1):D226–32.
- 590 62. Nawrocki EP, Eddy SR. Infernal 1.1: 100-fold faster RNA homology searches. *Bioinformatics.*
591 2013 Nov 15;29(22):2933–5.
- 592 63. Virtanen P, Gommers R, Oliphant TE, Haberland M, Reddy T, Cournapeau D, et al. SciPy 1.0--
593 Fundamental Algorithms for Scientific Computing in Python. *ArXiv190710121 Phys [Internet].*
594 2019 Jul 23 [cited 2020 Jan 2]; Available from: <http://arxiv.org/abs/1907.10121>
- 595 64. Shapiro SS, Wilk MB. An Analysis of Variance Test for Normality (Complete Samples).
596 *Biometrika.* 1965 Dec;52(3/4):591.
- 597 65. Welch BL. THE GENERALIZATION OF 'STUDENT'S' PROBLEM WHEN SEVERAL
598 DIFFERENT POPULATION VARLANCES ARE INVOLVED. *Biometrika.* 1947 Jan 1;34(1-
599 2):28–35.
- 600 66. Mann HB, Whitney DR. On a Test of Whether one of Two Random Variables is Stochastically
601 Larger than the Other. *Ann Math Stat.* 1947 Mar;18(1):50–60.
- 602 67. Benjamini Y, Hochberg Y. Controlling the False Discovery Rate: A Practical and Powerful
603 Approach to Multiple Testing. *J R Stat Soc Ser B Methodol.* 1995;57(1):289–300.

- 604 68. Dick GJ, Andersson AF, Baker BJ, Simmons SL, Thomas BC, Yelton AP, et al. Community-wide
605 analysis of microbial genome sequence signatures. *Genome Biol.* 2009;10(8):R85.
- 606 69. Karlin S, Campbell AM, Mrazek J. Comparative DNA analysis across diverse genomes. *Annu Rev*
607 *Genet.* 1998;32:185–225.
- 608 70. Wayne LG, Brenner DJ, Colwell RR, Grimont PAD, Kandler O, Krichevsky MI, et al. Report of the
609 Ad Hoc Committee on Reconciliation of Approaches to Bacterial Systematics. *Int J Syst Evol*
610 *Microbiol.* 1987;37(4):463–4.
- 611 71. Bohlin J, Snipen L, Hardy SP, Kristoffersen AB, Lagesen K, Donsvik T, et al. Analysis of intra-
612 genomic GC content homogeneity within prokaryotes. *BMC Genomics.* 2010 Aug 6;11:464.
- 613 72. Galtier N, Lobry JR. Relationships between genomic G+C content, RNA secondary structures, and
614 optimal growth temperature in prokaryotes. *J Mol Evol.* 1997 Jun;44(6):632–6.
- 615 73. Wixon J. Featured organism: reductive evolution in bacteria: *Buchnera* sp., *Rickettsia prowazekii*
616 and *Mycobacterium leprae*. *Comp Funct Genomics.* 2001;2(1):44–8.
- 617 74. Wu YW, Simmons BA, Singer SW. MaxBin 2.0: an automated binning algorithm to recover
618 genomes from multiple metagenomic datasets. *Bioinformatics.* 2016 Feb 15;32(4):605–7.
- 619 75. Albertsen M, Hugenholtz P, Skarshewski A, Nielsen KL, Tyson GW, Nielsen PH. Genome
620 sequences of rare, uncultured bacteria obtained by differential coverage binning of multiple
621 metagenomes. *Nat Biotechnol.* 2013 Jun;31(6):533–8.
- 622 76. Imelfort M, Parks D, Woodcroft BJ, Dennis P, Hugenholtz P, Tyson GW. GroopM: an automated
623 tool for the recovery of population genomes from related metagenomes. *PeerJ.* 2014;2:e603.
- 624 77. Pop M. Genome assembly reborn: recent computational challenges. *Brief Bioinform.* 2009
625 Jul;10(4):354–66.
- 626 78. Ghurye JS, Cepeda-Espinoza V, Pop M. Metagenomic Assembly: Overview, Challenges and
627 Applications. *Yale J Biol Med.* 2016;89(3):353.
- 628 79. Frost LS, Leplae R, Summers AO, Toussaint A. Mobile genetic elements: the agents of open source
629 evolution. *Nat Rev Microbiol.* 2005 Sep;3(9):722–32.
- 630 80. Hacker J, Kaper JB. Pathogenicity islands and the evolution of microbes. *Annu Rev Microbiol.*
631 2000;54:641–79.
- 632 81. Hurst LD, Merchant AR. High guanine–cytosine content is not an adaptation to high temperature: a
633 comparative analysis amongst prokaryotes. *Proc R Soc Lond B Biol Sci.* 2001 Mar
634 7;268(1466):493–7.
- 635 82. Schattner P. Searching for RNA genes using base-composition statistics. *Nucleic Acids Res.* 2002
636 May 1;30(9):2076–82.
- 637 83. Daubin V, Lerat E, Perriere G. The source of laterally transferred genes in bacterial genomes.
638 *Genome Biol.* 2003;4(9):R57.

84. Garcia-Vallve S, Romeu A, Palau J. Horizontal gene transfer in bacterial and archaeal complete genomes. *Genome Res.* 2000 Nov;10(11):1719–25.
85. Rocha EPC, Danchin A. Base composition bias might result from competition for metabolic resources. *Trends Genet.* 2002 Jun;18(6):291–4.
86. Frank JA, Pan Y, Tooming-Klunderud A, Eijsink VGH, McHardy AC, Nederbragt AJ, et al. Improved metagenome assemblies and taxonomic binning using long-read circular consensus sequence data. *Sci Rep.* 2016 Jul;6(1):25373.
87. Bertrand D, Shaw J, Kalathiyappan M, Ng AHQ, Kumar MS, Li C, et al. Hybrid metagenomic assembly enables high-resolution analysis of resistance determinants and mobile elements in human microbiomes. *Nat Biotechnol.* 2019 Aug;37(8):937–44.

Table 1 (on next page)

Reconstructed genome coverage and completeness

Table 1 Reconstructed genome coverage and completeness

Genome	Genome NCBI accessions	MAG NCBI accessions	MG Cov ^a	%CR ^b	NR ^c	mean NR length (bp)	NR length range
HL-46	EI34DRAFT_7210	GCA_001314525.1	3.9x	40%	284	4742	1007..42318
	EI34DRAFT_6181 ^d		3.9x	25%	7	18136	1108..49149
HL-48	CY41DRAFT	GCA_001314875.1	69x	95%	29	1892	330..53737
HL-49	K302DRAFT	GCA_001314815.1	9.7x	91%	89	3234	209..25366
HL-53	Ga0003345	GCA_001314555.1	113x	98%	15	1564	952..6133
HL-55	K417DRAFT	GCA_001314845.1	11x	95%	34	3574	417..45387
HL-58	CD01DRAFT	GCA_001314605.1	128x	99%	13	1124	959..12996
HL-91	Ga0058931_14	GCA_001314645.1	226x	97%	20	3129	135..11341
	Ga0058931_11 ^d		227x	97%	6	2188	914..4391
	Ga0058931_13 ^d		158x	0%	1	113349	113349
	Ga0058931_12 ^d		160x	0%	1	97917	97917
HL-93	Ga0071314	GCA_001314745.1	11x	85%	98	3605	232..78515
HL-109	Ga0071312_11	GCA_001314785.1	612x	87%	20	1835	204..63971
	Ga0071312_12		669x	92%	28	1285	506..52589
	Ga0071312_13 ^d		615x	95%	3	6053	1908..10088
HL-111	Ga0071316	GCA_001314765.1	18x	95%	39	1589	501..20407

^a Metagenomic read coverage

^b Percentage of the genome represented in the MAG

^c Number of not-binned regions

^d Predicted to be an extrachromosomal element

Table 2(on next page)

%G+C analysis

Table 2. Comparison of %G+C for genomes, CRs and NRs

		Genome	CRs			NRs		
molecule		mean	mean	distance	p-value	mean	distance	p-value
HL-46	EI34DRAFT_7210	64.42	63.96±1.94	1.55±1.25	0.997	65.12±2.13	1.61±1.56	0.263
HL-46	EI34DRAFT_6181	59.94	60.78±2.27	1.97±1.41	0.856	60.97±1.78	1.92±0.75	0.605
HL-48	CY41DRAFT	58.98	59.00±1.52	1.01±1.13	0.996	45.76±19.69	13.22±19.69	<0.001^a
HL-49	K302DRAFT	42.22	42.24±1.71	1.15±1.27	0.434	42.73±3.37	2.44±2.38	0.001
HL-53	Ga0003345	47.50	46.95±1.61	0.96±1.40	0.031	48.83±3.55	3.70±0.82	<0.001
HL-55	K417DRAFT	56.26	55.87±1.97	1.42±1.41	0.025	55.44±3.30	3.00±1.59	0.001
HL-58	CD01DRAFT	57.56	56.83±2.61	1.69±2.12	0.047	56.11±3.69	3.93±0.51	0.016
HL-91	Ga0058931_11	61.75	62.05±0.25	0.31±0.23	0.954	60.39±3.17	2.79±2.02	0.053
HL-91	Ga0058931_12	60.37	nd ^b	nd	nd	nd	nd	nd
HL-91	Ga0058931_13	61.77	nd	nd	nd	nd	nd	nd
HL-91	Ga0058931_14	61.84	60.99±1.90	1.33±1.60	0.030	59.11±2.96	3.52±1.96	0.005
HL-93	Ga0071314_11	55.88	56.75±2.20	1.75±1.59	1.000	56.08±4.42	3.6±2.57	<0.001
HL-109	Ga0071312_11	64.09	64.55±1.46	1.12±1.05	0.715	60.96±3.02	3.28±2.85	0.073
HL-109	Ga0071312_12	64.07	63.89±1.41	0.92±1.09	0.169	63.11±2.21	1.94±1.43	0.593
HL-109	Ga0071312_13	65.34	65.47±0.07	0.13±0.07	0.778	61.68±2.24	3.66±2.24	0.009
HL-111	Ga0071316_11	68.12	68.20±1.44	0.99±1.05	0.465	64.26±1.39	3.86±1.39	<0.001

^a Bold type indicates significant results ($P \leq 0.005$).

^b Not determined because the entire molecule was missing from the reconstructed genome.

Table 3(on next page)

Tetranucleotide frequency analysis

1 Table 3. Tetranucleotide frequency χ^2 analysis.

		CR			NR		
molecule		mean	sd	p-value	mean	sd	p-value
HL-46	EI34DRAFT_6181	0.2323	0.1883	0.154	0.1518	0.1429	0.983
HL-46	EI34DRAFT_7210	0.2042	0.0696	0.896	0.1701	0.1332	0.975
HL-48	CY41DRAFT	0.0276	0.0577	0.387	0.4425	0.2689	<0.001 ^a
HL-49	K302DRAFT	0.0522	0.0431	0.757	0.2340	0.2164	<0.001
HL-53	Ga0003345	0.0261	0.0451	0.001	0.3851	0.1525	<0.001
HL-55	K417DRAFT	0.0458	0.0726	0.086	0.2774	0.2168	0.004
HL-58	CD01DRAFT	0.0761	0.1451	0.008	0.2974	0.969	0.004
HL-91	Ga0058931_11	0.0266	0.0213	0.313	0.3043	0.1416	0.011
HL-91	Ga0058931_12	nd ^b	nd	nd	nd	nd	nd
HL-91	Ga0058931_13	nd	nd	nd	nd	nd	nd
HL-91	Ga0058931_14	0.0557	0.0647	0.004	0.3614	0.2052	<0.001
HL-93	Ga0071314_11	0.0925	0.0738	0.993	0.2254	0.1595	0.062
HL-109	Ga0071312_11	0.0262	0.0401	0.396	0.3148	0.1842	0.087
HL-109	Ga0071312_12	0.0216	0.0281	0.076	0.2907	0.1913	0.231
HL-109	Ga0071312_13	0.0048	0.0019	0.538	0.3651	0.2299	0.016
HL-111	Ga0071316_11	0.0396	0.0561	0.322	0.4504	0.1640	<0.001

^a Bold text indicates significant result

^b Not determined because the entire molecule was missing from the reconstructed genome.

Table 4(on next page)

Genomic redundancy

1 **Table 4. Metagenomic redundancy.**

		Genome	CR			NR		
molecule		mean	mean	distance	p-value	mean	distance	p-value
HL-46	EI34DRAFT_6181	2.76	2.78	0.26±0.15	0.992	2.42	0.43±0.31	0.264
HL-46	EI34DRAFT_7210	5.98	4.43	2.95±2.34	0.978	4.99	4.01±13.35	0.649
HL-48	CY41DRAFT	72.40	69.29	3.65±2.45	1.000	140.93	100.97±153.33	<0.001^a
HL-49	K302DRAFT	8.97	8.67	0.51±0.58	0.999	11.38	4.16±17.52	0.002
HL-53	Ga0003345	441.81	446.29	24.51±18.21	0.073	517.07	115.56±59.71	<0.001
HL-55	K417DRAFT	16.76	15.35	7.06±10.45	0.679	117.37	110.35±333.81	<0.001
HL-58	CD01DRAFT	128.28	127.85	9.10±15.71	1.000	180.14	60.44±27.54	<0.001
HL-91	Ga0058931_11	231.39	228.46	3.64±2.25	0.786	311.6	91.27±97.44	0.001
HL-91	Ga0058931_12	163.24	nd ^b	nd	nd	nd	nd	nd
HL-91	Ga0058931_13	168.27	nd	nd	nd	nd	nd	nd
HL-91	Ga0058931_14	227.56	231.77	8.18±6.59	0.220	273.03	97.82±117.47	<0.001
HL-93	Ga0071314_11	50.87	50.03	4.04±2.92	1.000	65.73	16.16±35.87	<0.001
HL-109	Ga0071312_11	3103.11	3098.73	97.47±72.15	0.748	3072.59	323.24±240.86	0.005
HL-109	Ga0071312_12	2821.18	2822.26	113.08±78.18	0.124	2778.03	352.81±436.28	0.003
HL-109	Ga0071312_13	2853.84	2901.40	47.56±9.73	0.179	2097.01	756.83±256.91	0.018
HL-111	Ga0071316_11	90.14	88.03	3.98±4.31	0.993	98.25	38.42±104.87	0.027

2 a Bold text indicates significant result

3 b Not determined because the entire molecule was missing from the reconstructed genome.

4

5

Figure 1

Distributions of %G+C for MDR and CDR genomic regions.

G+C composition was determined for individual regions identified as CDRs or MDRs. Bar height represents the percentage of regions in the category. Black bars, CDRs; white bars, MDRs.

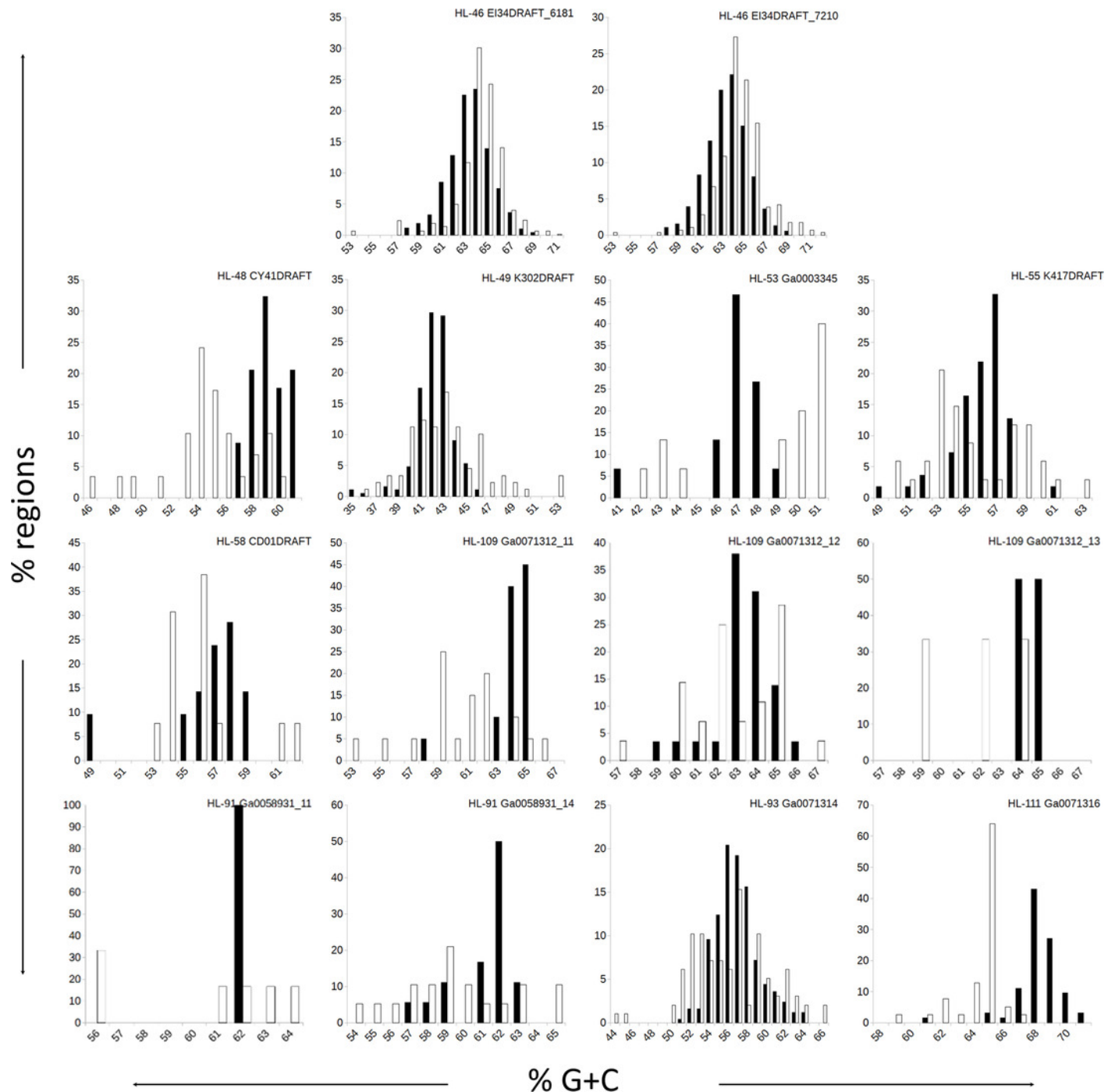
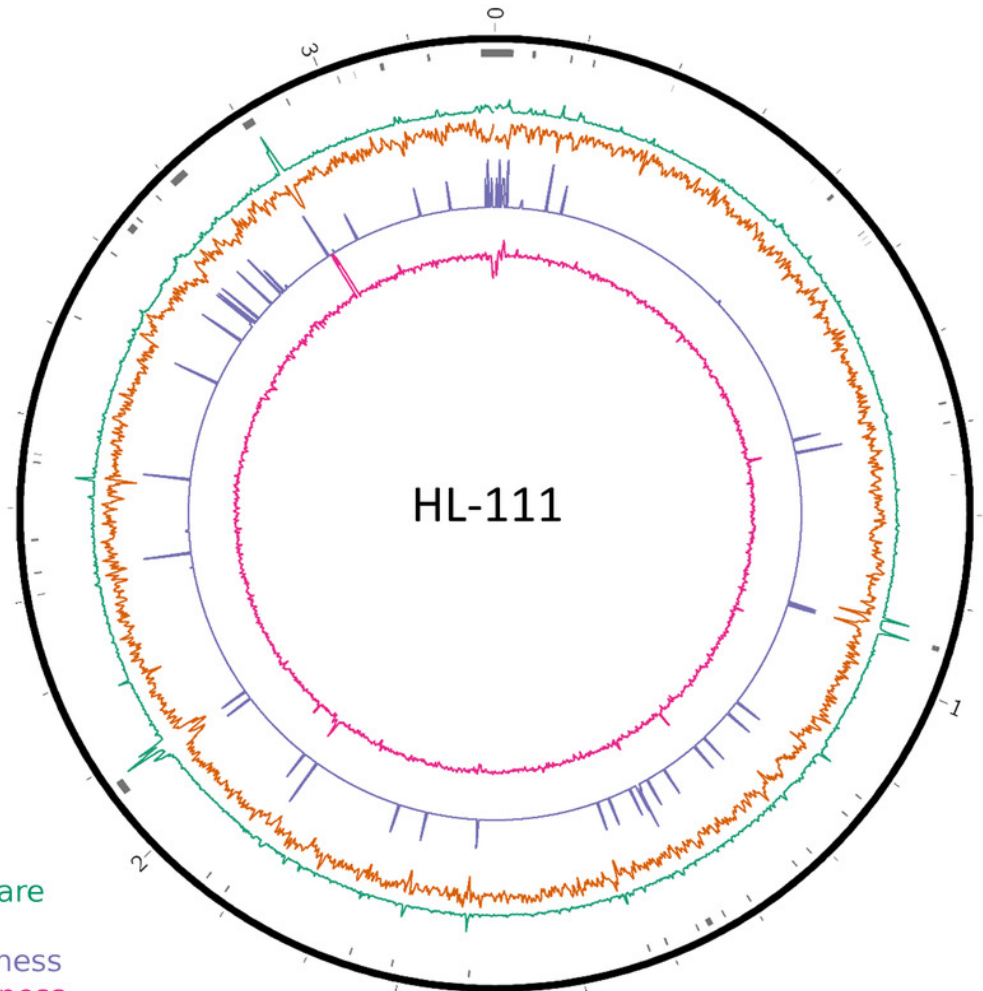


Figure 2

Analysis of HL-111 genome.

Ring 1 (outermost, black) – genome sequence; ring 2 (grey bars) – missed detection regions (MDRs); ring 3 (teal) – tetranucleotide frequency (TNF) distance χ^2 values; ring 4 (orange) – %G+C; ring 5 (blue) – intragenome redundancy; ring 6 (magenta) – metagenome redundancy. Values were calculated across 2000 nt windows with a step size of 1000 nt. For TNF, χ^2 was calculated for the windows using the whole molecule frequencies as the expected. Data for other genomes analyzed is presented in Figure S1. Circular plots were generated using Circos v0.69.3 (Krzywinski, Schein et al. 2009) .



Outermost to innermost

1. Scale in Mb
2. Genome
3. Not-binned regions
4. Tetranucleotide chi square
5. %G+C
6. Intragenome repetitiveness
7. Metagenome repetitiveness

Ga0071316 11

Figure 3

Repeat content of genomes versus MAGs

Box plot representation of the total fraction of each genome/MAG in a repeat region as determined by NUCmer ($\geq 97\%$ identity; center line, median; box limits, upper and lower quartiles; whiskers, $1.5 \times$ interquartile range; diamonds, outliers). UCC MAG and genome comparison were significantly different ($p = 0.01$; Mann-Whitney U).

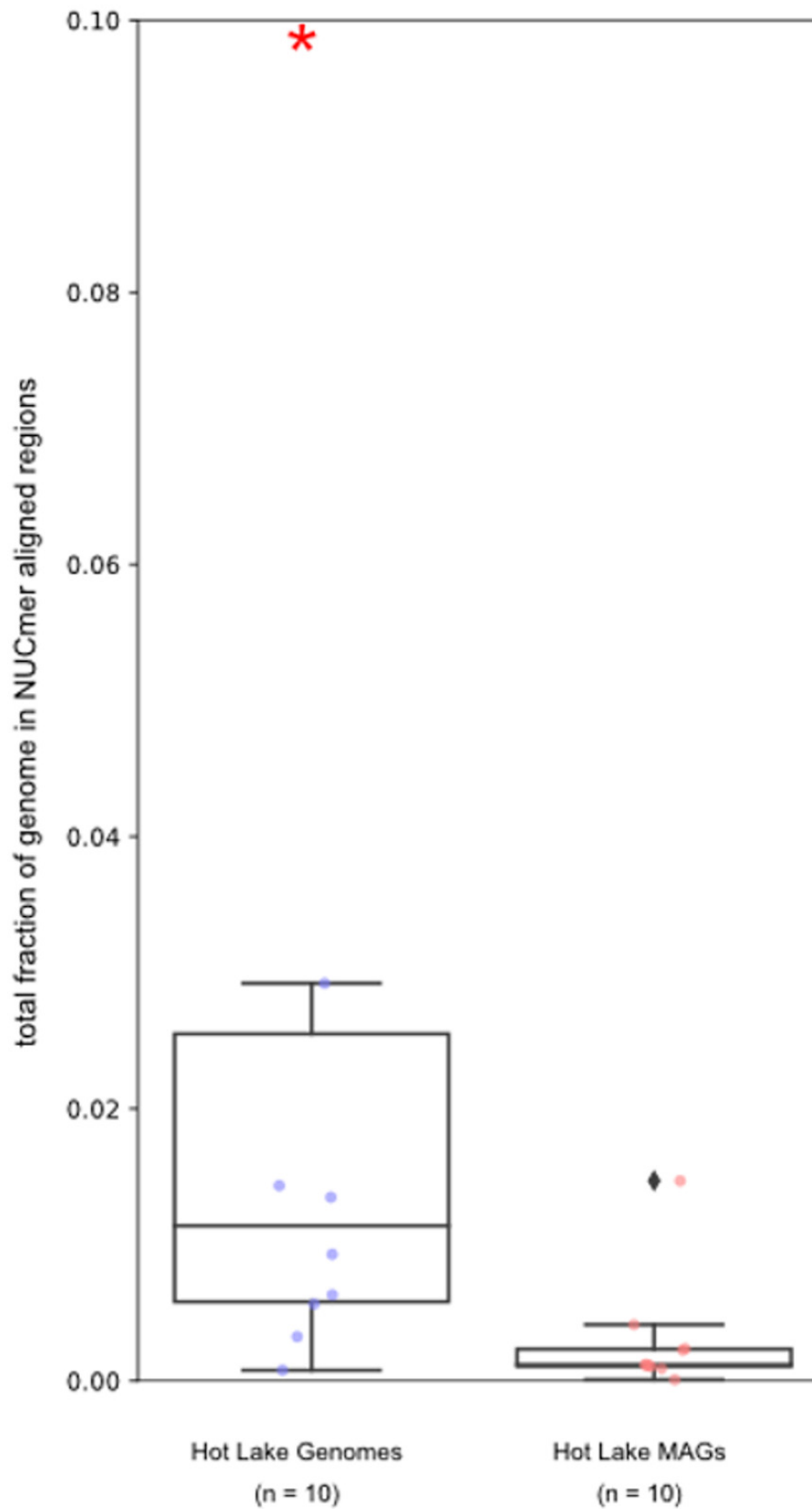


Figure 4

Functional categorization of genes present on MDRs.

The gene features of each genome region were assigned to functional COG categories or as non-coding genes (rRNA; tRNA; ncRNA). Organisms' gene sets were compared using Principal Component Analysis. Organisms are represented by colors (HL-46, yellow; HL-48, purple; HL-49, blue; HL-53, light blue; HL-55, gray; HL-58, orange; HL-91, black; HL-93, pink; HL-109, red; HL-111, green). The genome region categories are represented by shapes (whole isolate genomes, circles; CDRs, squares; MDRs, triangles; extrachromosomal elements, diamonds). COG categories: *A* - RNA processing and modification; *B* - Chromatic structure and dynamics; *C* - Energy production and conversion; *D* - Cell cycle control, cell division, chromosome partitioning; *E* - Amino acid transport and metabolism; *F* - Nucleotide transport and metabolism; *G* - Carbohydrate transport and metabolism; *H* - Coenzyme transport and metabolism; *I* - Lipid transport and metabolism; *J* - Translation, ribosomal structure and biogenesis; *K* - Transcription; *L* - DNA replication, recombination and repair; *M* - Cell wall/membrane/envelope biogenesis; *N* - Cell motility; *O* - Post-translational modification, protein turnover, chaperones; *P* - Inorganic ion transport and metabolism; *Q* - Secondary metabolites biosynthesis, transport and catabolism; *R* - General function prediction; *S* - Function unknown; *T* - Signal transduction mechanisms; *U* - Intracellular trafficking, secretion and vesicular transport; *V* - Defense mechanisms; *W* - Extracellular structures; *X* - Mobilome, transposons, phages; *Y* - Nuclear structure; *Z* - Cytoskeleton.

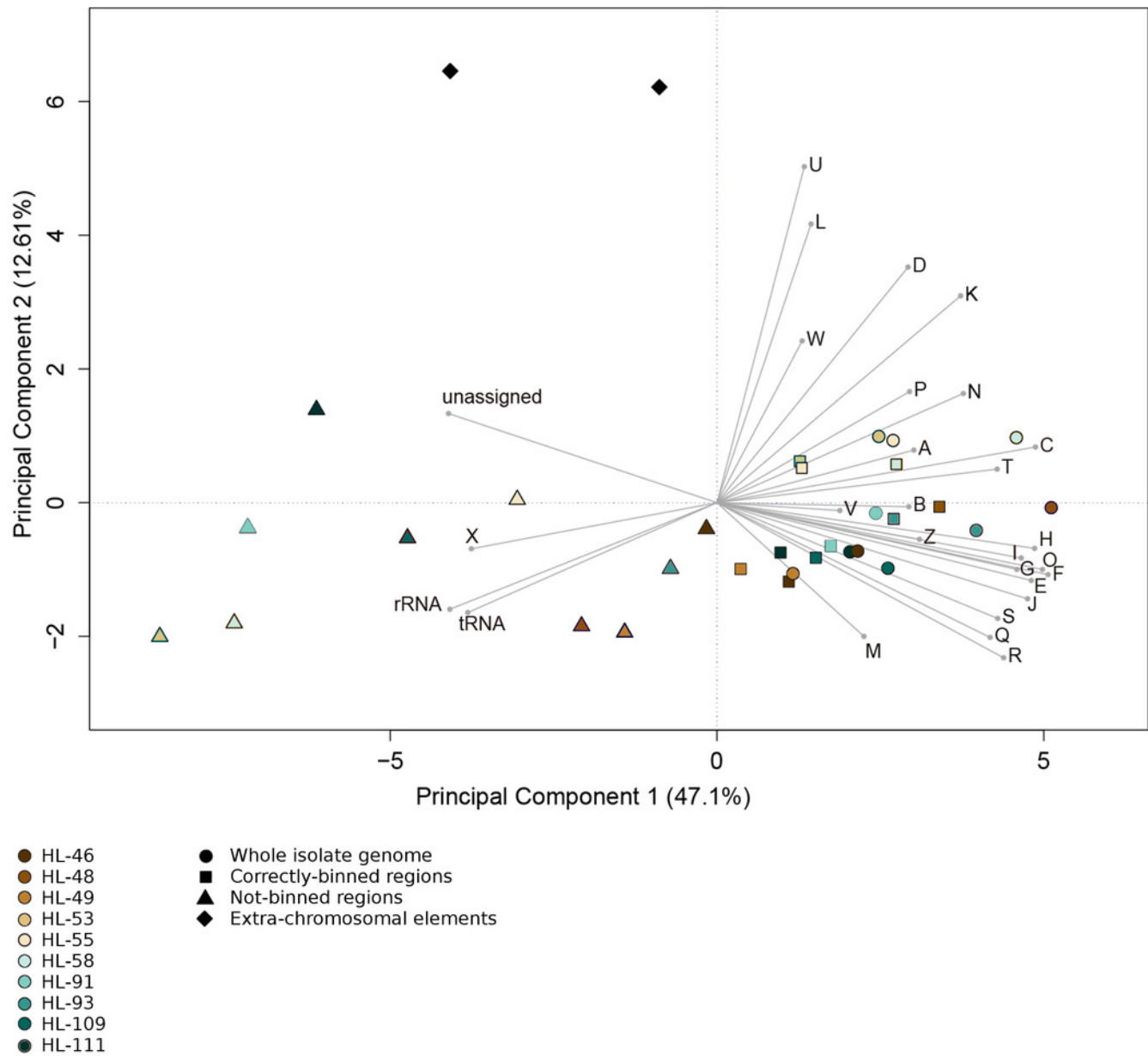


Figure 5

Tara Ocean MAG nucleotide composition analysis

(A) %G+C variance analysis. Box plot representation of the %G+C variance for each 2,000bp segment of genome/MAG (sliding window step: 500bp; center line, median; box limits, upper and lower quartiles; whiskers, 1.5×interquartile range; diamonds, outliers). Comparisons between *Tara* Oceans MAG datasets and RefSeq genomes were significantly different ($p < 0.001$; Mann-Whitney U with Benjamini-Hochberg False Discovery Rate Correction (BH FDR)).

(B) Tetranucleotide analysis. Box plot representation of the variance in Pearson correlation values of the tetranucleotide Z-scores for a pair-wise comparison of each 10kb segment of genome/MAG (sliding window step: 5kb; center line, median; box limits, upper and lower quartiles; whiskers 1.5x interquartile range; diamonds, outliers). Comparisons between *Tara* Oceans MAG datasets and RefSeq genomes were significantly different ($p < 0.001$; Mann-Whitney U with BH FDR Correction). Red asterisks denote the existence of outliers outside of the displayed range.

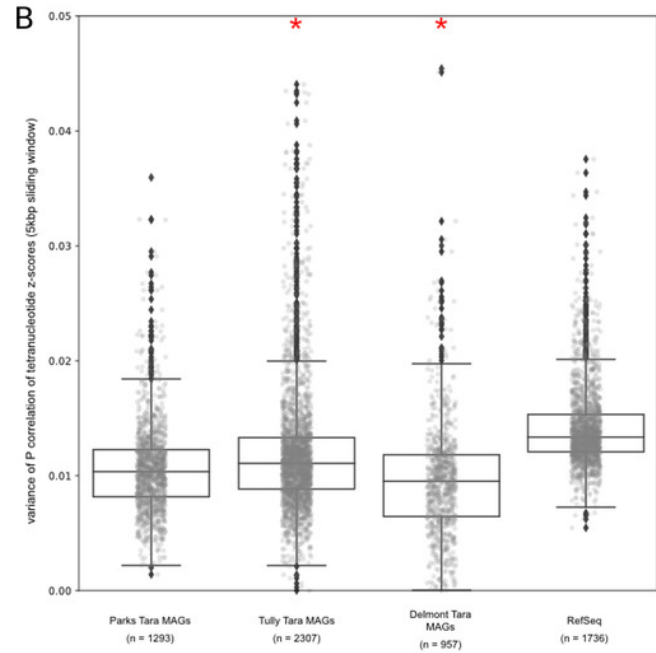
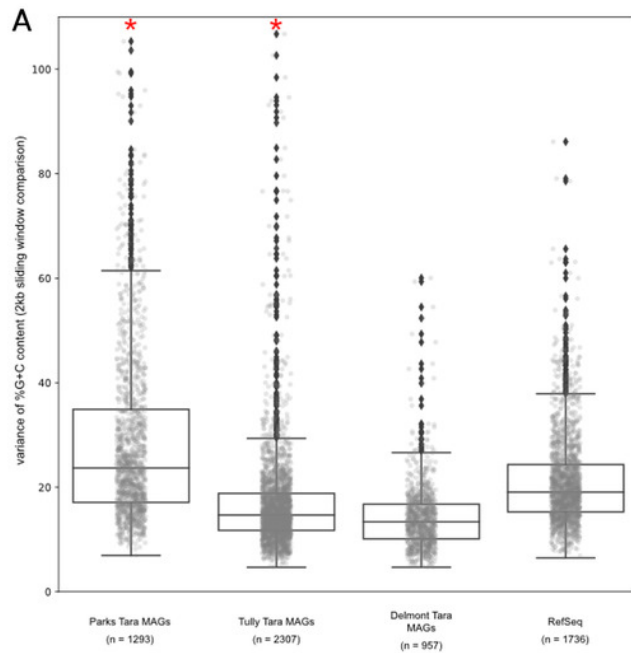


Figure 6

Tara Ocean MAG repeat content

Box plot representation of the total fraction of each genome/MAG in a repeat region as determined by NUCmer ($\geq 97\%$ identity; center line, median; box limits, upper and lower quartiles; whiskers, $1.5 \times$ interquartile range; diamonds, outliers). Comparisons between *Tara* Oceans MAG datasets and RefSeq genomes were significantly different ($p < 0.001$; Mann-Whitney U with BH FDR Correction). Red asterisks denote the existence of outliers outside of the displayed dataset.

