

1. Project Information

Program	Microbial/CSP 2012
PMO Project	0
Seq Proj ID	1027067
Sequencing Project Name	Geoarchaeota archaeon JGI 000156CP-C13
JGI Project ID	0

2. Read Statistics

Illumina Std PE Statistics

File name	7667.5.80864.CTTGTA.fastq
Library	TGSN
Number of reads	25,141,154
Sequencing depth [†]	754X
Read type	2x150 bp

[†] A genome size of 5.0 Mbp was assumed in this calculation.

3. Read QC Results

The following are the results of reads screened against contaminants. Pairs of matching reads were removed from the dataset.

Illumina Std PE Read Filter Statistics

Description	Num Reads	Pct Reads
Input	25,141,154	100
Contam removed	156	0.0
Artifact removed	783,208	3.1
Total removed	5,141,154	20.4
Total remaining	20,000,000	79.6

List of Contaminants Removed

Description	Num Reads	Pct Reads
gi 357579577 Canis_lupus_familiaris_chr3	98	0.00
human_chr2	92	0.00
gi 357579535 Canis_lupus_familiaris_chr20	22	0.00
gi 357579571 Canis_lupus_familiaris_chr5	10	0.00
human_chr1	4	0.00
human_chr7	4	0.00
human_chr8	4	0.00

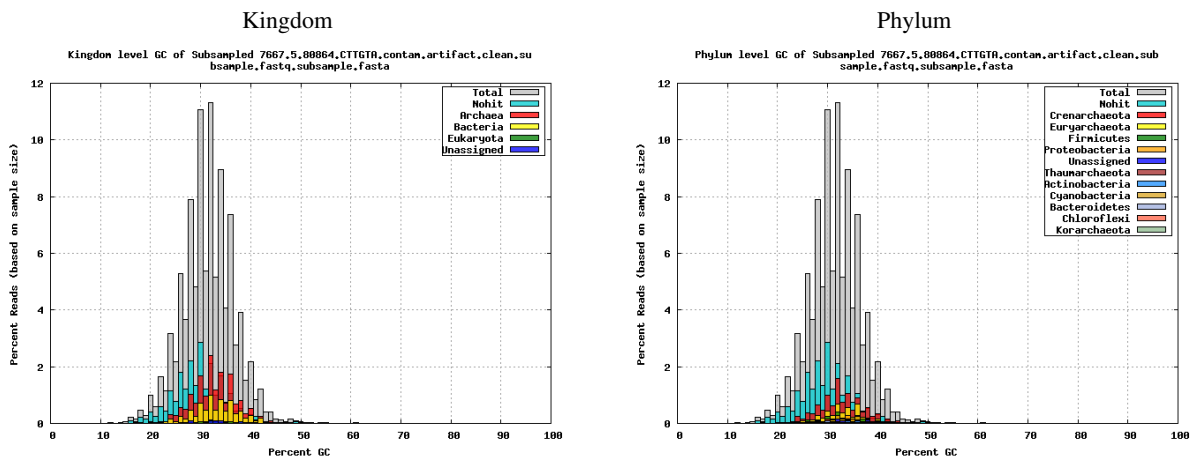
human_chr20	4	0.00
human_chr22	2	0.00
human_chr19	2	0.00
human_chr9	2	0.00
human_chr6	2	0.00
gi 357579575 Canis_lupus_familiaris_chr4	2	0.00
human_chr17	2	0.00
human_chr5	2	0.00
human_chrX	2	0.00
human_chr12	2	0.00
gi 362110686 Felis_catus_breed_Abyssinian_chrA1	2	0.00

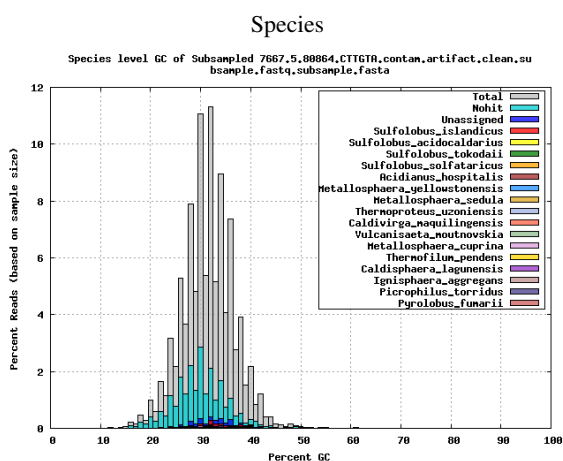
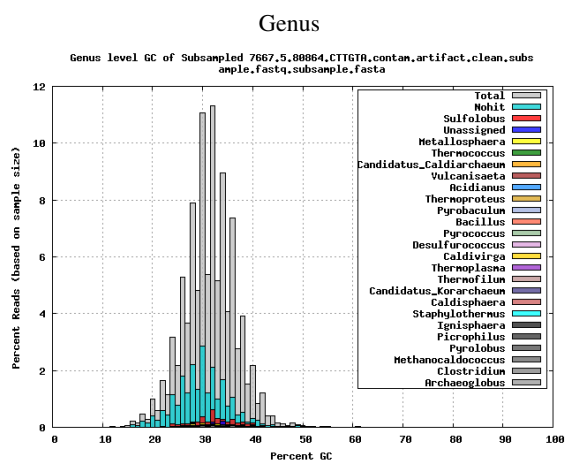
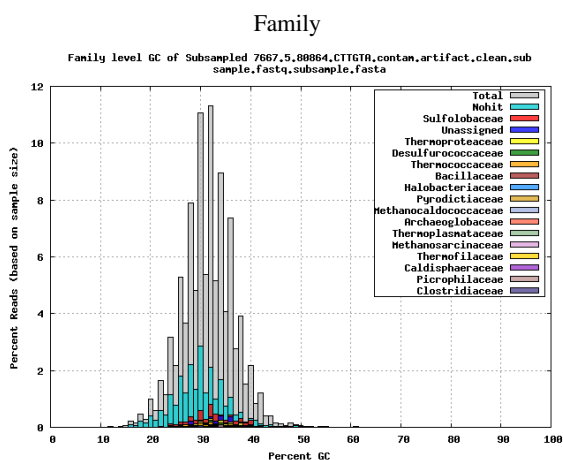
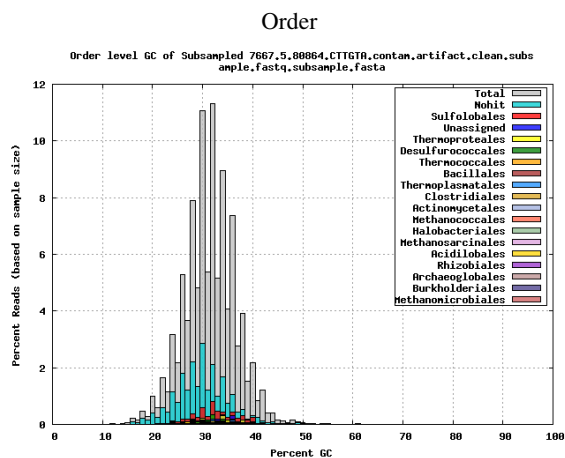
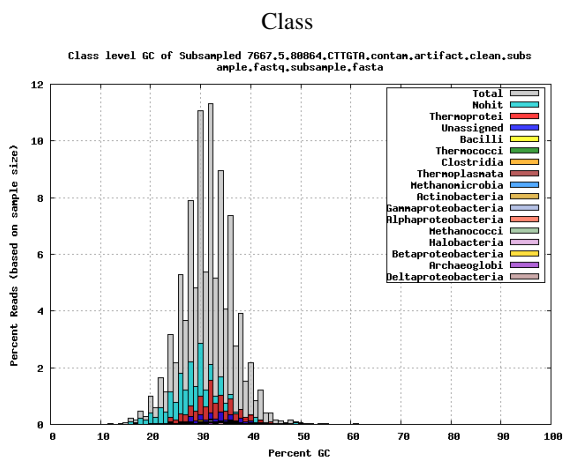
The following are the results of reads screened against potential reagent and process contaminants but were not removed from the dataset.

Illumina Std PE Contamination Identification Statistics

Description	Num Reads	Pct Reads
Input	25,141,154	100
Contam identified	0	0.0

GC histogram of the reads subsampled to 10k, overlaid with GC of hits based on BLASTX, shown for different taxonomic levels.



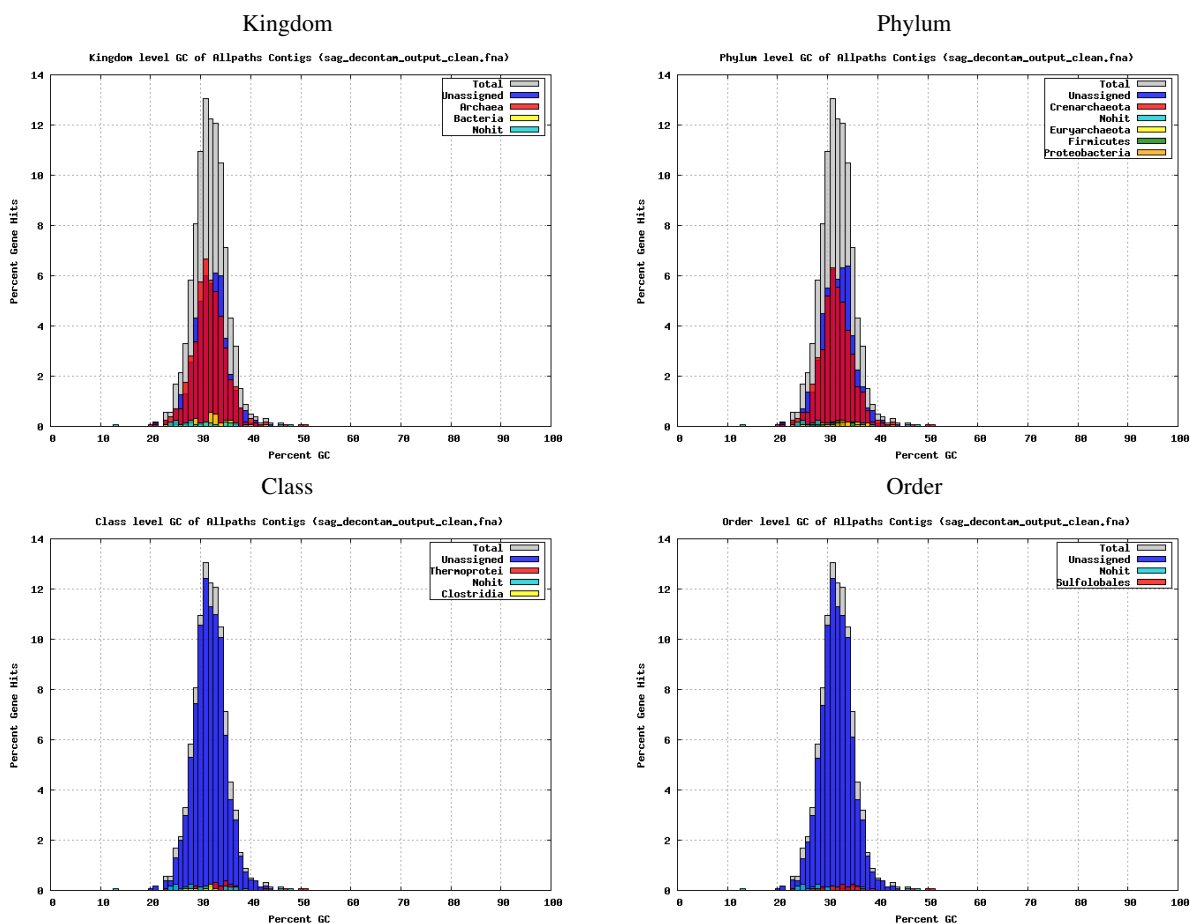


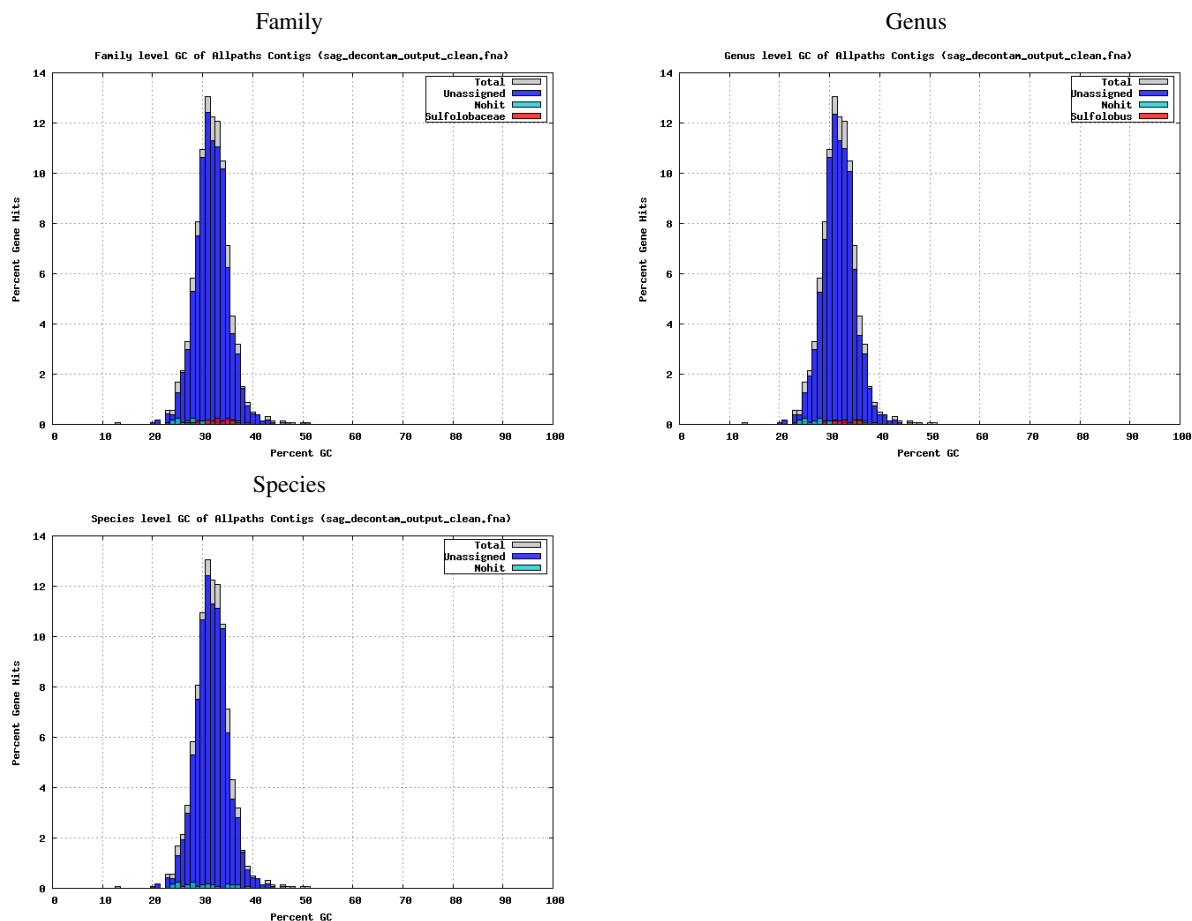
4. Assembly Statistics

Assembly method	SPAdes with auto decontamination
Scaffold total	66
Contig total	66
Scaffold sequence length	1.4 Mb
Contig sequence length	1.4 Mb (0.0% gap)
Scaffold N/L50	14/32.5 kb
Contig N/L50	14/32.5 kb
Largest Contig	97.4 kb
Number of scaffolds >50 kb	7
Pct of genome in scaffolds >50 kb	32.1
Pct of reads assembled (raw)	93.5
Pct of reads assembled (decontam)	93.0

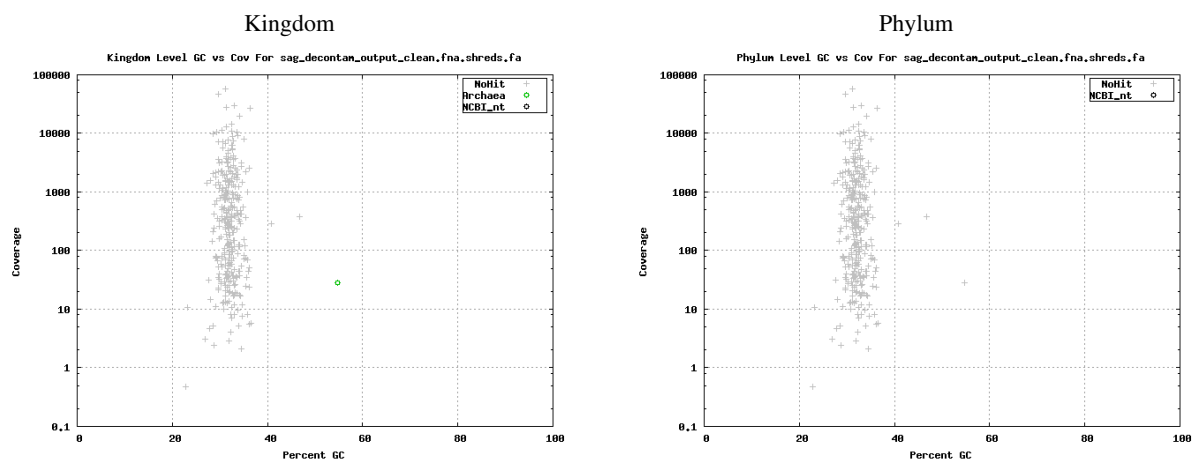
5. Assembly QC Results

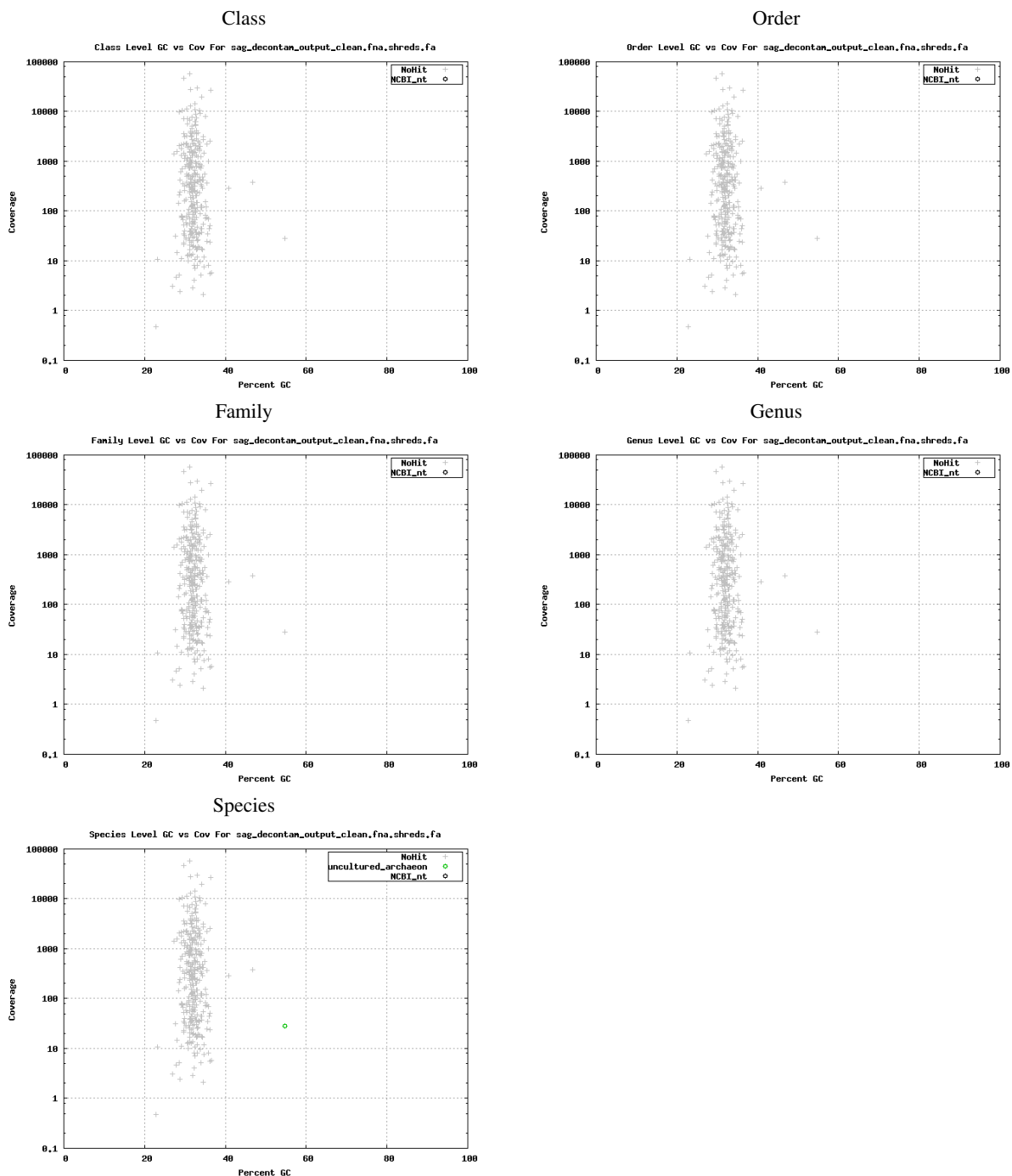
GC histogram of the predicted genes on each contig, overlaid with GC of hits based on BLASTP, shown for different taxonomic levels.



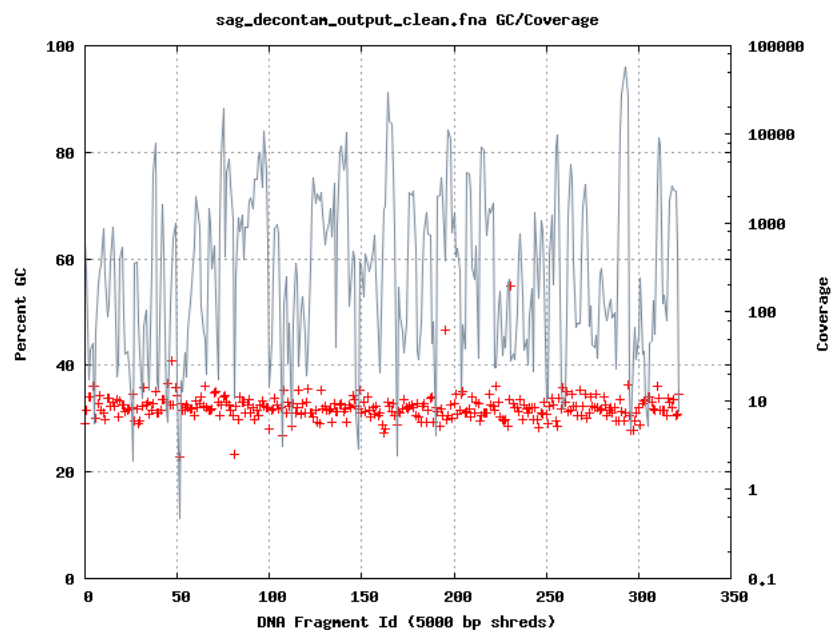


GC vs coverage based on GC of NCBI nt and Greengenes 16S rRNA gene hits to the assembly using megablast, shown for different taxonomic levels.

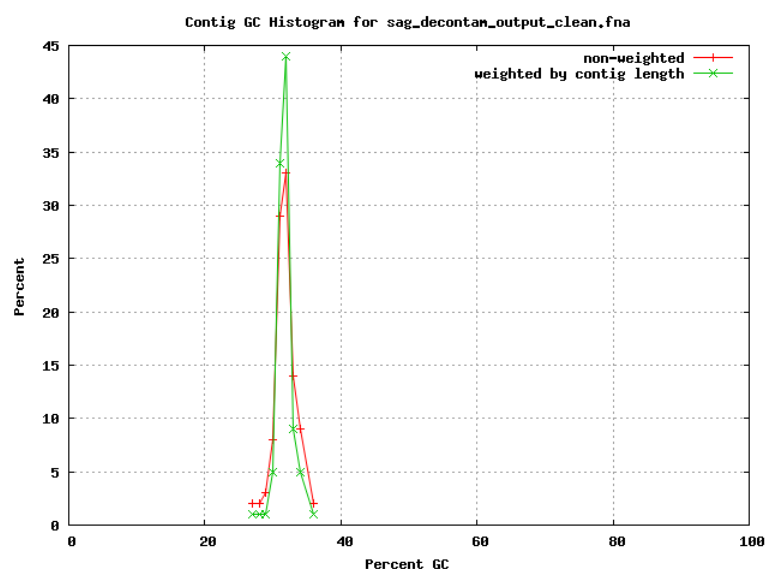




Coverage vs GC. Contigs were shredded into non-overlapping 5kbp and the GC of each shred was plotted as a point, colored by scaffold id. Coverage was calculated by mapping the fragment library to the final assembly and plotted as connected points.



GC histogram of the contigs, including contig length weighted distribution.

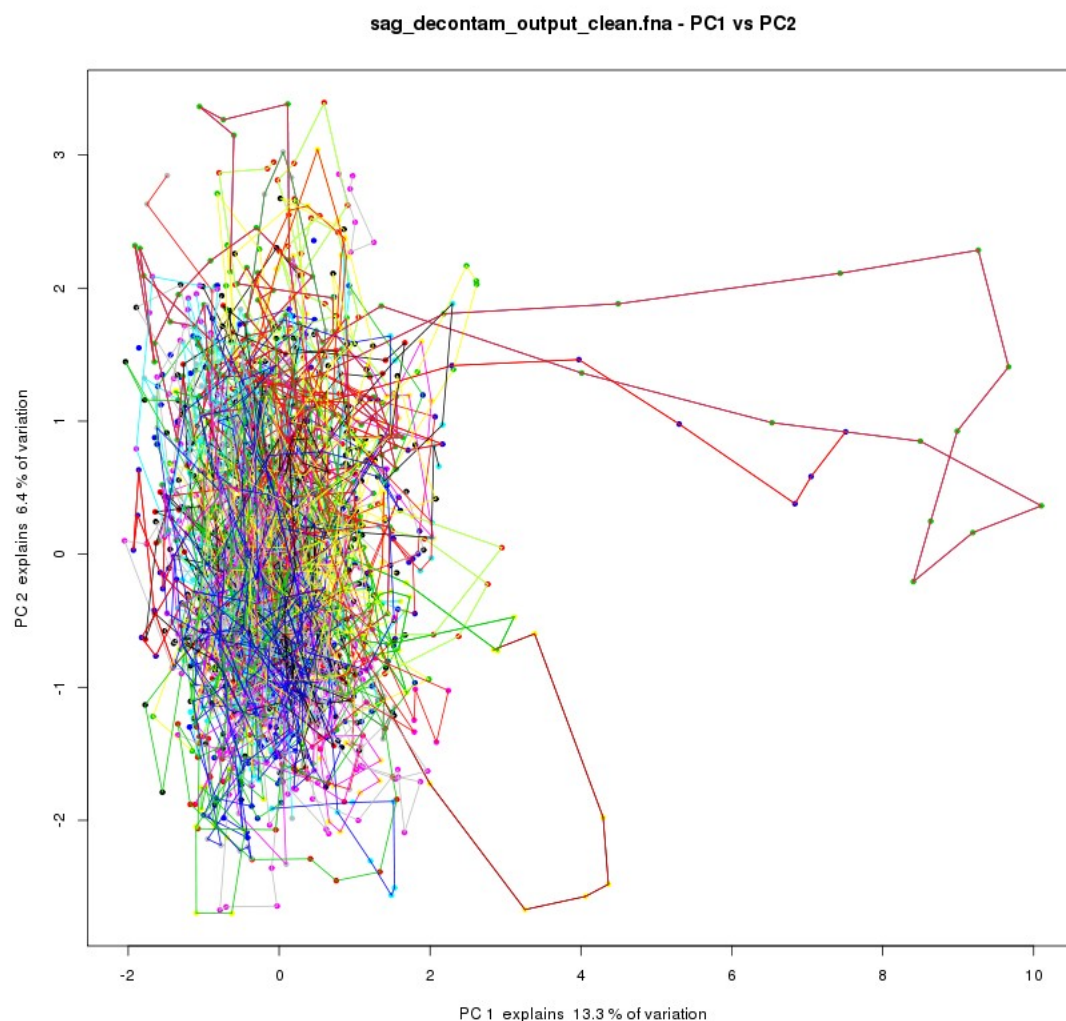


List of contigs and average percent GC, grouped in bins of 5:

Pct GC Bin	Contig Name
25	NODE.42.length.10847.cov.663.839.ID.83, NODE.44.length.10207.cov.5700.28.ID.87, NODE.51.length.7434.cov.29.2347.ID.101, NODE.54.length.7051.cov.347.768.ID.107
30	NODE.1.length.97427.cov.96.1335.ID.1, NODE.2.length.72741.cov.1413.66.ID.3, NODE.3.length.69530.cov.567.83.ID.5, NODE.4.length.59964.cov.221.421.ID.7, NODE.5.length.59368.cov.1026.88.ID.9, NODE.6.length.53955.cov.1703.51.ID.11, NODE.7.length.50437.cov.9750.34.ID.13, NODE.8.length.46193.cov.411.317.ID.15, NODE.9.length.42870.cov.2624.03.ID.17, NODE.10.length.42697.cov.754.857.ID.19,

	NODE.11.length.40139.cov.3322.35.ID.21, NODE.12.length.36440.cov.1903.36.ID.23, NODE.13.length.34008.cov.111.407.ID.25, NODE.14.length.32526.cov.287.537.ID.27, NODE.15.length.30720.cov.1111.3.ID.29, NODE.16.length.28763.cov.408.019.ID.31, NODE.17.length.28647.cov.319.764.ID.33, NODE.18.length.27169.cov.456.869.ID.35, NODE.19.length.25417.cov.8135.65.ID.37, NODE.20.length.24240.cov.1299.84.ID.39, NODE.21.length.23783.cov.230.19.ID.41, NODE.22.length.23754.cov.864.144.ID.43, NODE.23.length.23623.cov.230.091.ID.45, NODE.24.length.23360.cov.124.492.ID.47, NODE.25.length.21888.cov.339.169.ID.49, NODE.26.length.21857.cov.134.886.ID.51, NODE.27.length.17222.cov.16.3829.ID.53, NODE.28.length.16526.cov.20.6844.ID.55, NODE.29.length.16039.cov.75.1689.ID.57, NODE.30.length.15691.cov.188.588.ID.59, NODE.31.length.15511.cov.148.537.ID.61, NODE.32.length.15162.cov.376.94.ID.63, NODE.33.length.14850.cov.44.8602.ID.65, NODE.34.length.14626.cov.70.5512.ID.67, NODE.35.length.13874.cov.13.9135.ID.69, NODE.36.length.13624.cov.445.824.ID.71, NODE.37.length.13068.cov.4282.45.ID.73, NODE.38.length.13023.cov.84.9055.ID.75, NODE.39.length.11843.cov.109.406.ID.77, NODE.40.length.11797.cov.4187.84.ID.79, NODE.41.length.11157.cov.83.3294.ID.81, NODE.43.length.10684.cov.42.9131.ID.85, NODE.46.length.9658.cov.257.723.ID.91, NODE.47.length.9392.cov.3806.18.ID.93, NODE.48.length.8935.cov.67.6847.ID.95, NODE.49.length.8883.cov.17.611.ID.97, NODE.50.length.8790.cov.16.2425.ID.99, NODE.52.length.7193.cov.26.6751.ID.103, NODE.53.length.7186.cov.1098.27.ID.105, NODE.56.length.6790.cov.5.05494.ID.111, NODE.57.length.6619.cov.6.86213.ID.113, NODE.58.length.6253.cov.31.4076.ID.115, NODE.59.length.6218.cov.33.5543.ID.117, NODE.60.length.5881.cov.7.95434.ID.123, NODE.61.length.5849.cov.2598.64.ID.125, NODE.62.length.5758.cov.42.3191.ID.127, NODE.63.length.5734.cov.8.96584.ID.129, NODE.64.length.5683.cov.11.4128.ID.131, NODE.65.length.5432.cov.4.65501.ID.133, NODE.66.length.5257.cov.5.40811.ID.135, NODE.67.length.4888.cov.1132.61.ID.137
35	NODE.45.length.9943.cov.180.073.ID.89

Principal component analysis of tetramer frequencies of contigs. Detectable variations are highlighted in color.



Estimated genome recovery derived from analysis of universal single-copy genes detected in final assembly.

HMM	Pct Recovered
bacteria	48.76 %
archaea	86.42 %

6. Sequence Data Availability

The following sequence fasta files can be downloaded from our JGI portal website.

<http://www.jgi.doe.gov/genome-projects>

Filename	Description
sag_decontam_output_clean.fna	SPAdes with auto decontamination

7. Annotation Data Availability

The annotation of the assembled contigs can be found within IMG.

<http://img.jgi.doe.gov>

8. Methods

Single Cell Minimal Draft

Genome sequencing and assembly

The draft genome of was generated at the DOE Joint genome Institute (JGI) using the Illumina technology [1]. An Illumina std shotgun library was constructed and sequenced using the Illumina HiSeq 2000 platform which generated 25,141,154 reads totaling 3,771.2 Mb. All general aspects of library construction and sequencing performed at the JGI can be found at <http://www.jgi.doe.gov>. All raw Illumina sequence data was passed through DUK, a filtering program developed at JGI, which removes known Illumina sequencing and library preparation artifacts [2]. Following steps were then performed for assembly: (1) artifact filtered Illumina reads were assembled using SPAdes [3] (version 3.0.0), (3) Parameters for assembly steps were `-t 16 -m 120 -sc -careful -12`. The final draft assembly contained 66 contigs in 66 scaffolds, totalling 1.4 Mb in size. The final assembly was based on 3,000.0 Mb of Illumina data. Based on a presumed genome size of 5.0 Mb, the average input read coverage used for the assembly was 600.0X.

Genome annotation

Genes were identified using Prodigal [4], followed by a round of manual curation using GenePRIMP [5] for finished genomes and Draft genomes in fewer than 20 scaffolds. The predicted CDSs were translated and used to search the National Center for Biotechnology Information (NCBI) nonredundant database, UniProt, TIGRFam, Pfam, KEGG, COG, and InterPro databases. The tRNAScanSE tool [6] was used to find tRNA genes, whereas ribosomal RNA genes were found by searches against models of the ribosomal RNA genes built from SILVA [7]. Other non-coding RNAs such as the RNA components of the protein secretion complex and the RNase P were identified by searching the genome for the corresponding Rfam profiles using INFERNAL [8]. Additional gene prediction analysis and manual functional annotation was performed within the Integrated Microbial Genomes (IMG) platform [9] developed by the Joint Genome Institute, Walnut Creek, CA, USA [10].

1. Bennett S. Solexa Ltd. Pharmacogenomics. 2004;5(4):433–8.
2. Mingkun L, Copeland A, Han J. DUK, unpublished, 2011.
3. Bankevich A, et.al, SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing. J Comput Biol 2012; 19:455–77.
4. Hyatt D, Chen GL, Lacascio PF, Land ML, Larimer FW, Hauser LJ. Prodigal: prokaryotic gene recognition and translation initiation site identification. BMC Bioinformatics 2010; 11:119.
5. Pati A, Ivanova NN, Mikhailova N, Ovchinnikova G, Hooper SD, Lykidis A, Kyrpides NC. GenePRIMP: a gene prediction improvement pipeline for prokaryotic genomes. Nat Methods 2010; 7:455–457.
6. Lowe TM, Eddy SR. tRNAscan-SE: a program for improved detection of transfer RNA genes in genomic sequence. Nucleic Acids Res 1997; 25:955–964.
7. Pruesse E, Quast C, Knittel, Fuchs B, Ludwig W, Peplies J, Glckner FO. SILVA: a comprehensive online resource for quality checked and aligned ribosomal RNA sequence data compatible with ARB. Nuc Acids Res 2007; 35: 2188–7196.
8. INFERNAL. Inference of RNA alignments. <http://infernal.janelia.org>.
9. The Integrated Microbial Genomes (IMG) platform. <http://www.ncbi.nlm.nih.gov/pubmed/24165883>
10. Markowitz VM, Mavromatis K, Ivanova NN, Chen IMA, Chu K, Kyrpides NC. IMG ER: a system for microbial genome annotation expert review and curation. Bioinformatics 2009; 25:2271–2278.