

1. Project Information

Program	Microbial/CSP 2012
PMO Project	0
Seq Proj ID	1027094
Sequencing Project Name	NAG2 archaeon OSP418 JGI 000156CP-E12
JGI Project ID	0

2. Read Statistics

Illumina Std PE Statistics

File name	7666.6.80850.GTTTCG.fastq
Library	TGSZ
Number of reads	29,120,364
Sequencing depth [†]	874X
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	29,120,364	100
Contam removed	126	0.0
Artifact removed	1,102,668	3.8
Total removed	9,120,364	31.3
Total remaining	20,000,000	68.7

List of Contaminants Removed

Description	Num Reads	Pct Reads
gi 357579577 Canis_lupus_familiaris_chr3	54	0.00
human_chr2	54	0.00
gi 357579535 Canis_lupus_familiaris_chr20	18	0.00
gi 357579571 Canis_lupus_familiaris_chr5	14	0.00
human_chr4	10	0.00
human_chr1	8	0.00
human_chrX	6	0.00

human_chr10	6	0.00
human_chr18	2	0.00
human_chr11	2	0.00
human_chr19	2	0.00
human_chr7	2	0.00
human_chr3	2	0.00
human_chr6	2	0.00
human_chr15	2	0.00
gi 357579523 Canis_lupus_familiaris_chr27	2	0.00
human_chr17	2	0.00
human_chr5	2	0.00
human_chr16	2	0.00
human_chr12	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

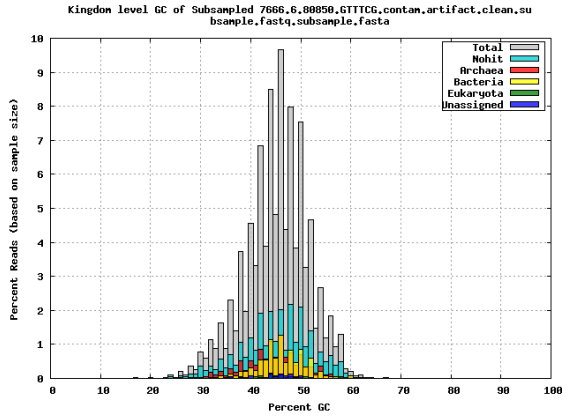
Description	Num Reads	Pct Reads
Input	29,120,364	100
Contam identified	8	0.0

List of Contaminants Identified

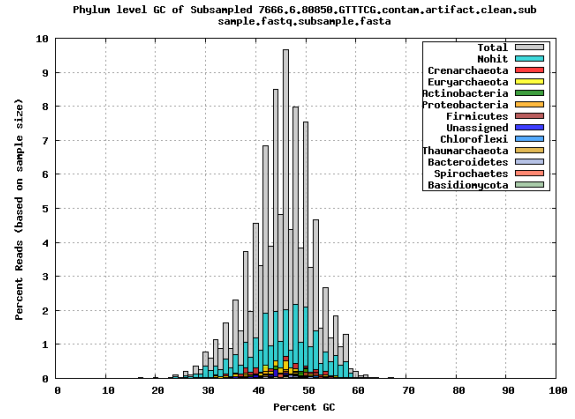
Description	Num Reads	Pct Reads
<i>Escherichia</i>	2	0.00
<i>Delftia</i>	2	0.00
<i>Klebsiella</i>	2	0.00
<i>Shigella</i>	2	0.00

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

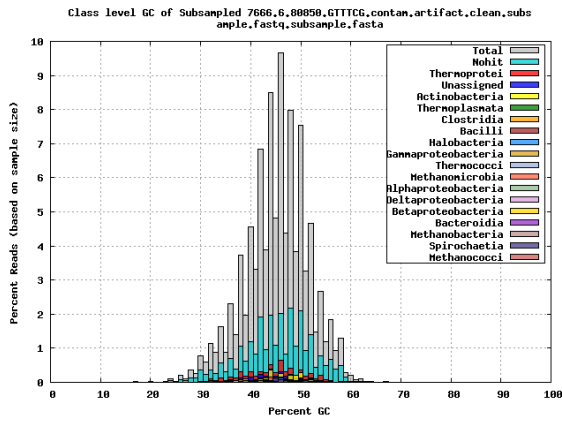
Kingdom



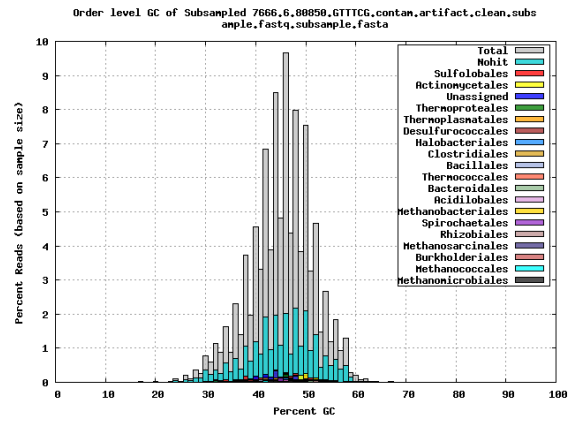
Phylum



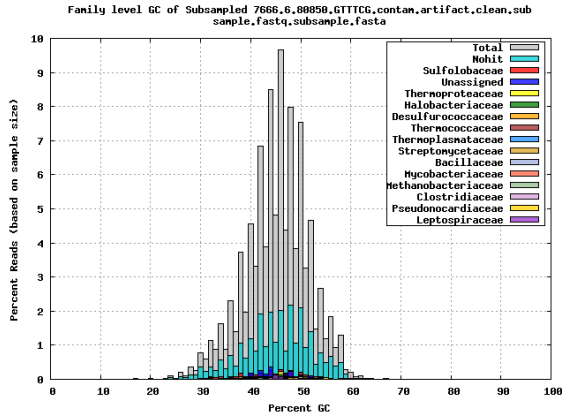
Class



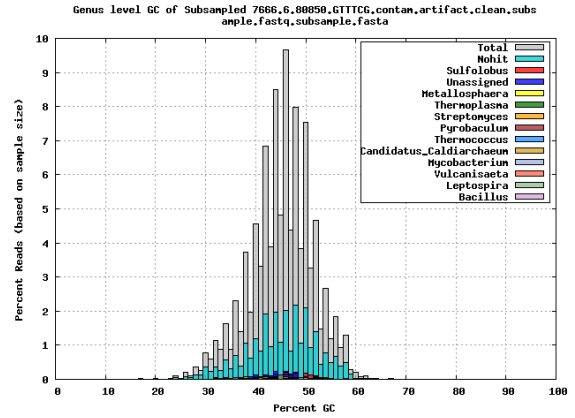
Order

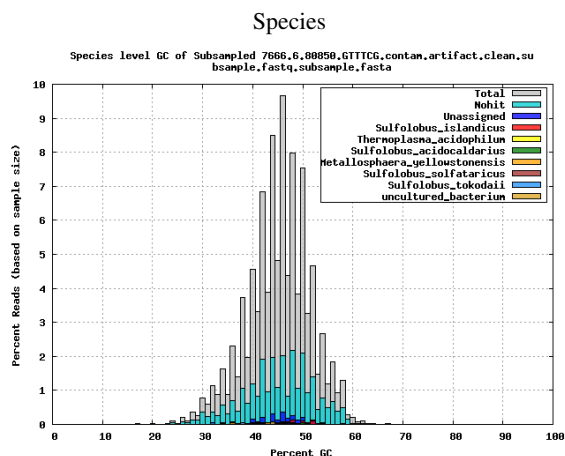


Family



Genus



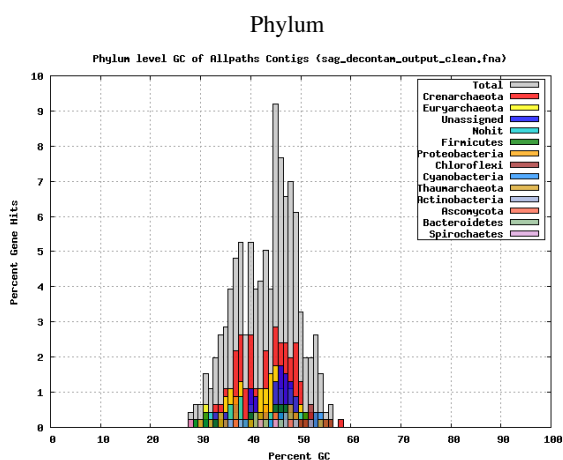
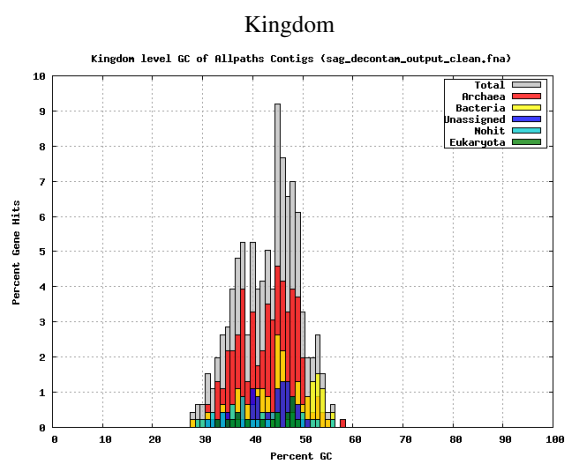


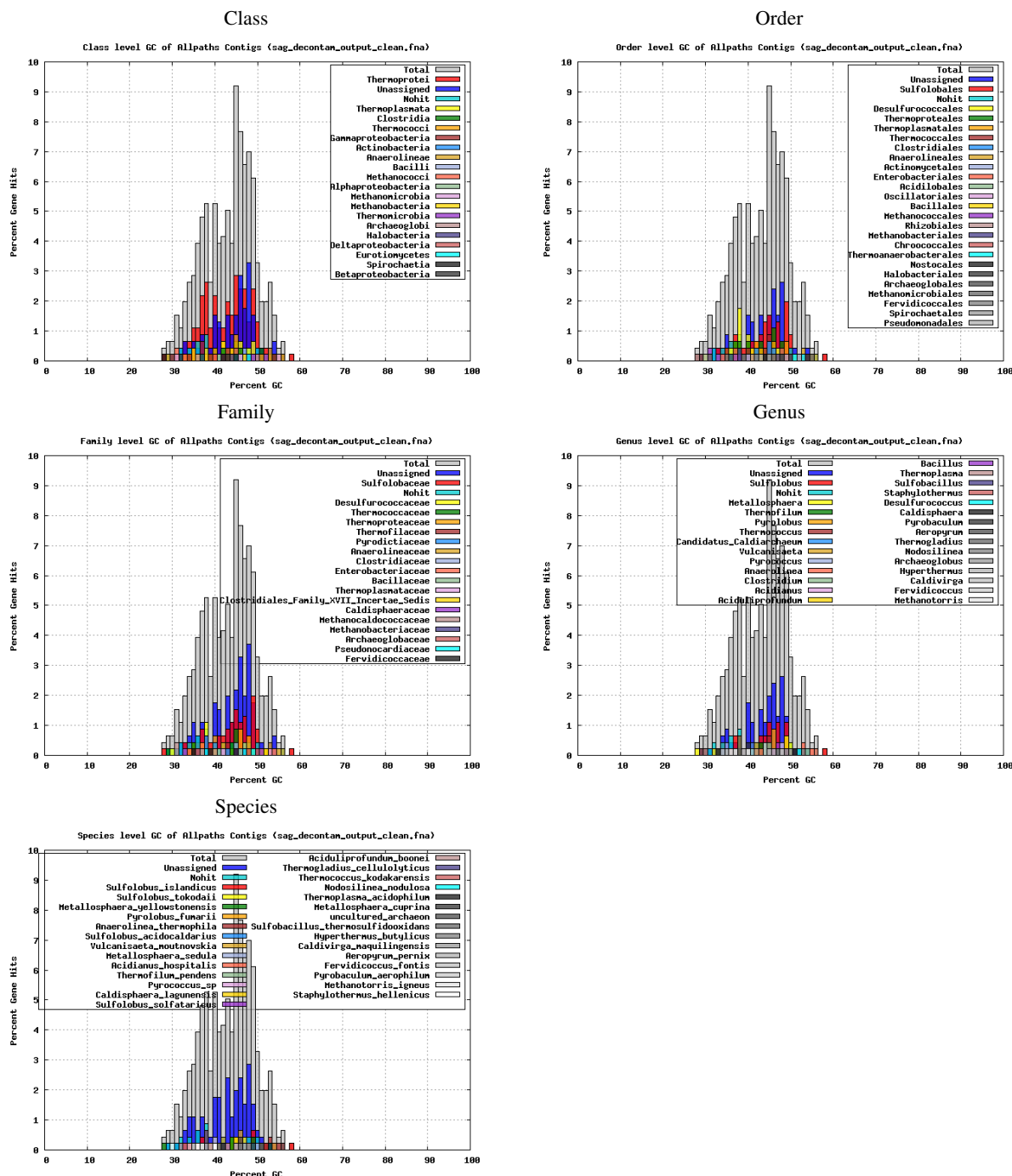
4. Assembly Statistics

Assembly method	SPAdes with auto decontamination
Scaffold total	52
Contig total	52
Scaffold sequence length	415.4 kb
Contig sequence length	415.4 kb (0.0% gap)
Scaffold N/L50	19/8.3 kb
Contig N/L50	19/8.3 kb
Largest Contig	20.8 kb
Number of scaffolds >50 kb	0
Pct of genome in scaffolds >50 kb	0.0
Pct of reads assembled (raw)	69.8
Pct of reads assembled (decontam)	44.5

5. Assembly QC Results

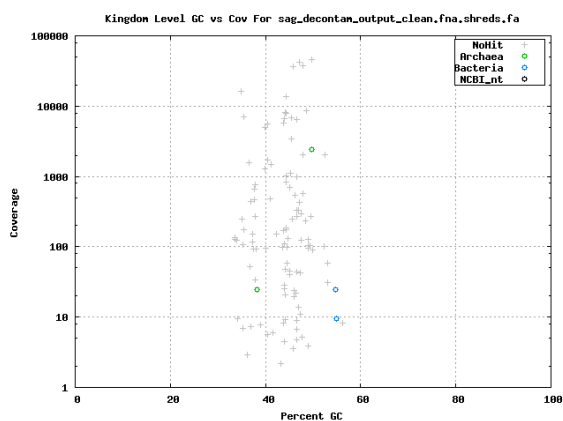
GC histogram of the predicted genes on each contig, overlaid with GC of hits based on BLASTP, shown for different taxonomic levels.



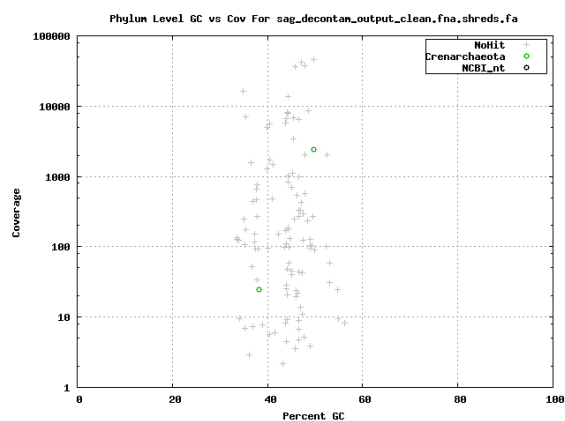


GC vs coverage based on GC of NCBI nt and Greengenes 16S rRNA gene hits to the assembly using megablast, shown for different taxonomic levels.

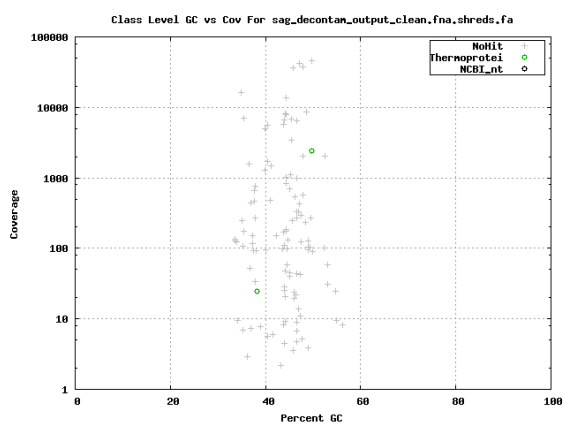
Kingdom



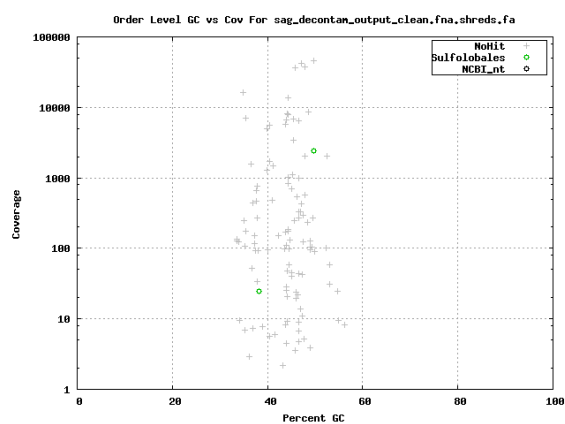
Phylum



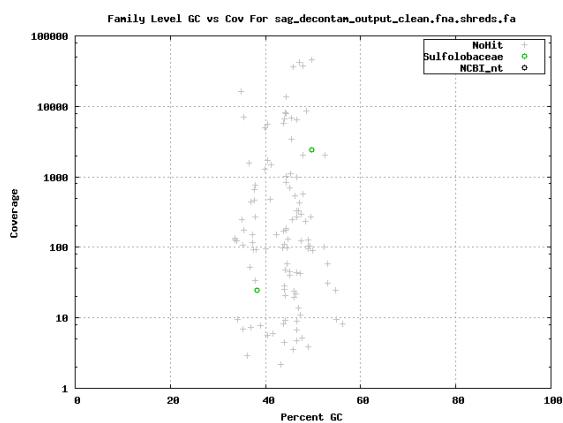
Class



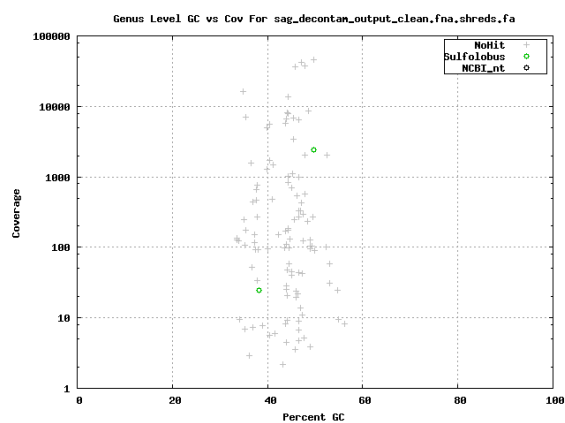
Order

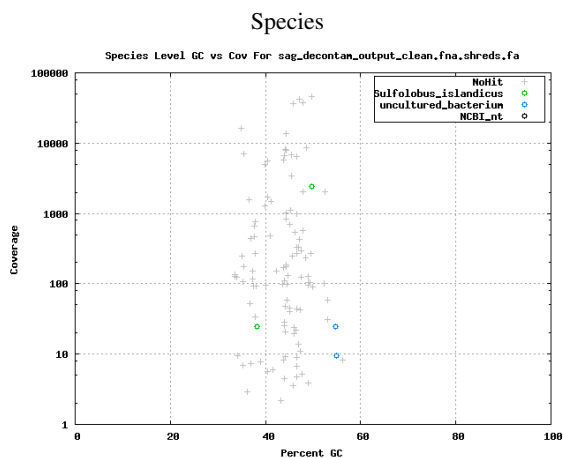


Family

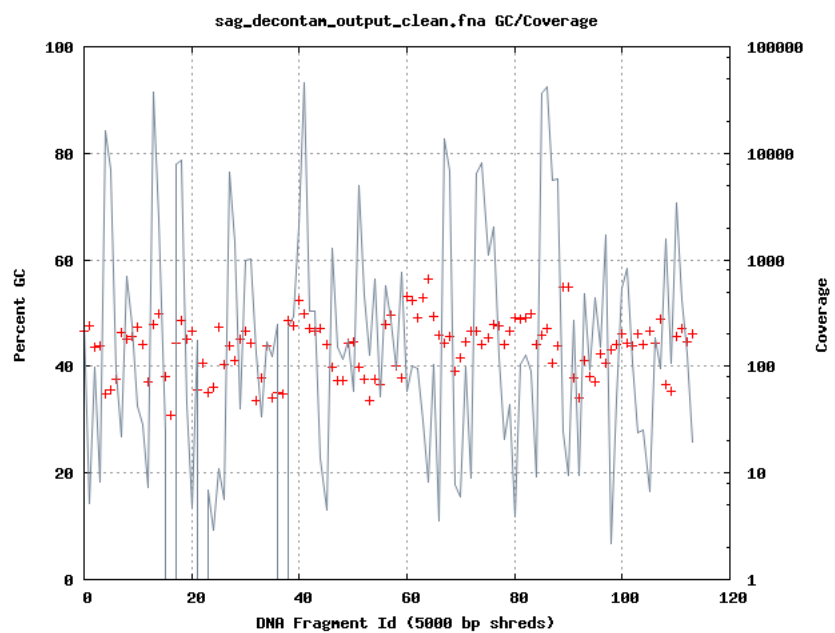


Genus

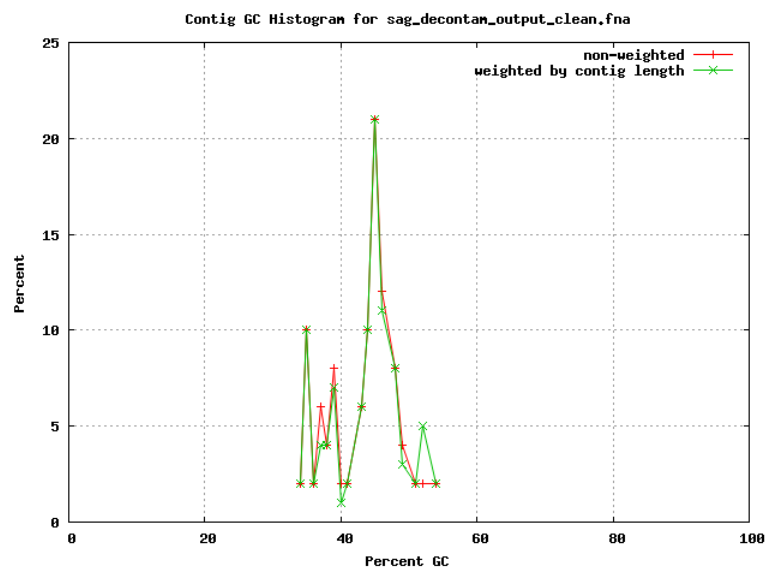




Coverage vs GC. Contigs were shredded into non-overlapping 5kbp and the GC of each shred was plotted as a point, colored by scaffold id. Coverage was calculated by mapping the fragment library to the final assembly and plotted as connected points.



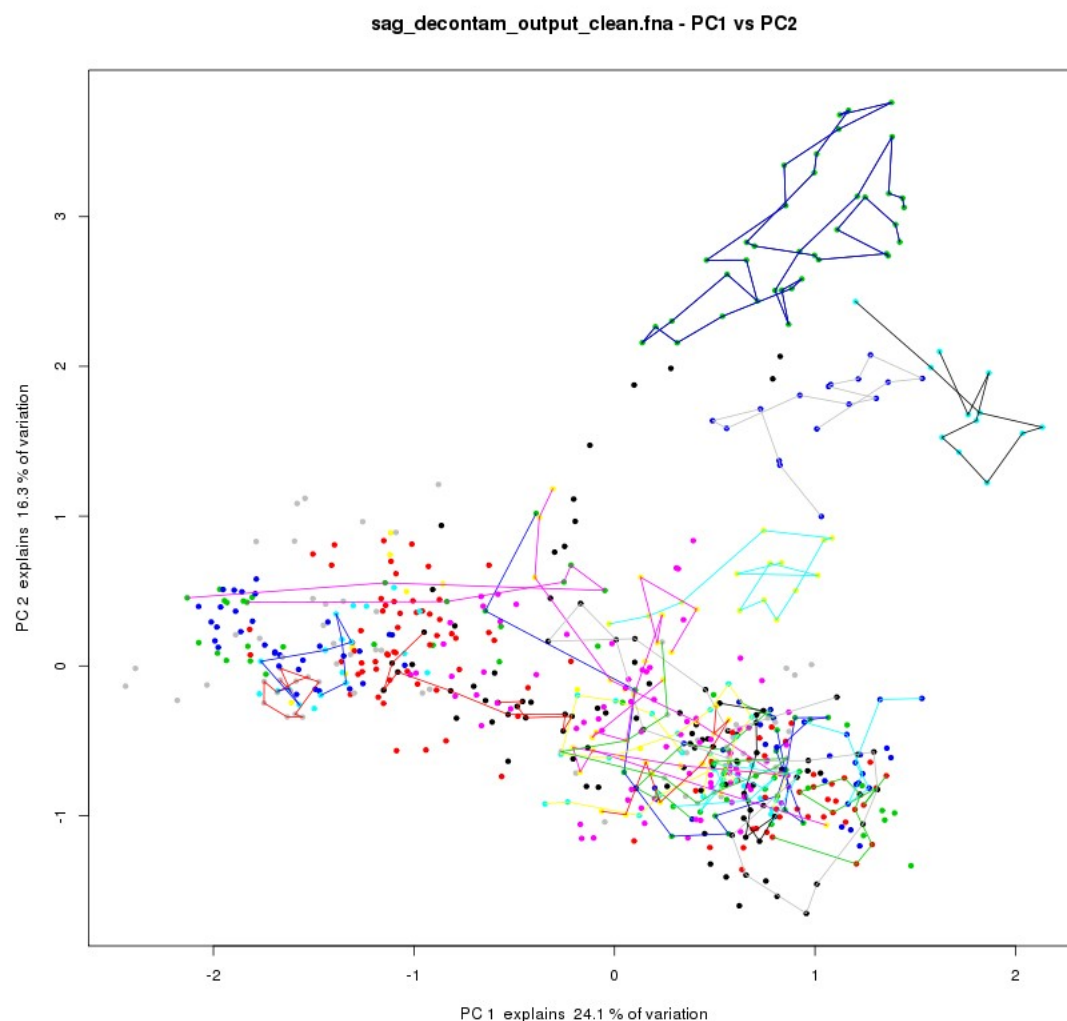
GC histogram of the contigs, including contig length weighted distribution.



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

Pct GC Bin	Contig Name
30	NODE.9.length.10023.cov.127.44.ID.17
35	NODE.4.length.13048.cov.408.372.ID.7, NODE.5.length.12144.cov.226.19.ID.9, NODE.11.length.9698.cov.283.817.ID.21, NODE.12.length.9608.cov.57.2369.ID.23, NODE.18.length.8371.cov.219.522.ID.35, NODE.24.length.7630.cov.8913.47.ID.51, NODE.29.length.6658.cov.797.46.ID.59, NODE.33.length.6276.cov.2694.45.ID.67, NODE.35.length.6233.cov.3.95403.ID.71, NODE.36.length.6194.cov.153.481.ID.73, NODE.40.length.5596.cov.284.758.ID.81, NODE.46.length.5408.cov.59.9813.ID.93, NODE.48.length.5353.cov.5.00113.ID.97, NODE.50.length.5139.cov.16.2565.ID.109 NODE.51.length.5101.cov.115.398.ID.111
40	NODE.2.length.17108.cov.302.418.ID.3, NODE.3.length.14383.cov.1886.54.ID.5, NODE.15.length.8833.cov.5762.45.ID.29, NODE.23.length.7689.cov.7482.31.ID.49, NODE.37.length.6002.cov.15.2181.ID.75, NODE.38.length.5979.cov.105.632.ID.77, NODE.43.length.5494.cov.5531.12.ID.87, NODE.44.length.5474.cov.1033.5.ID.89, NODE.49.length.5194.cov.67.3006.ID.103, NODE.54.length.4822.cov.110.886.ID.107
45	NODE.6.length.10941.cov.972.523.ID.11, NODE.7.length.10510.cov.20.5761.ID.13, NODE.8.length.10184.cov.24.2527.ID.15, NODE.10.length.9762.cov.26382.4.ID.19, NODE.14.length.9331.cov.37.885.ID.27, NODE.16.length.8697.cov.15247.1.ID.31, NODE.17.length.8478.cov.679.32.ID.33, NODE.19.length.8290.cov.179.88.ID.37, NODE.20.length.8184.cov.83.6514.ID.39, NODE.22.length.7828.cov.4875.78.ID.47, NODE.25.length.7385.cov.233.077.ID.53, NODE.26.length.7305.cov.320.618.ID.55, NODE.27.length.7188.cov.1667.47.ID.41, NODE.28.length.6775.cov.397.898.ID.57, NODE.30.length.6596.cov.69.9654.ID.61, NODE.31.length.6552.cov.113.034.ID.63, NODE.32.length.6501.cov.6.88148.ID.65, NODE.34.length.6244.cov.25.0439.ID.69, NODE.39.length.5846.cov.8.29097.ID.79, NODE.41.length.5520.cov.54.8675.ID.83, NODE.42.length.5506.cov.63.3614.ID.85, NODE.52.length.5082.cov.177.673.ID.113 NODE.53.length.4866.cov.16.0143.ID.119
50	NODE.1.length.20816.cov.45.9797.ID.1, NODE.13.length.9573.cov.14943.5.ID.25 NODE.21.length.7958.cov.12.4955.ID.45

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	15.19 %
archaea	26.75 %

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 29,120,364 reads totaling 4,368.1 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 52 contigs in 52 scaffolds, totalling 415.4 Kb 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.