

1. Project Information

Program	Microbial/CSP 2012
PMO Project	0
Seq Proj ID	1027085
Sequencing Project Name	NAG2 archaeon MK4 SK304 JGI 000156CP-M15
JGI Project ID	0

2. Read Statistics

Illumina Std PE Statistics

File name	7666.6.80850.GTCCGC.fastq
Library	TGSW
Number of reads	24,831,480
Sequencing depth [†]	745X
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	24,831,480	100
Contam removed	82	0.0
Artifact removed	127,202	0.5
Total removed	4,831,480	19.5
Total remaining	20,000,000	80.5

List of Contaminants Removed

Description	Num Reads	Pct Reads
human_chr2	42	0.00
gi 357579577 Canis_lupus_familiaris_chr3	38	0.00
gi 357579535 Canis_lupus_familiaris_chr20	18	0.00
gi 357579571 Canis_lupus_familiaris_chr5	10	0.00
human_chr1	2	0.00
human_chr5	2	0.00
gi 357579553 Canis_lupus_familiaris_chr10	2	0.00

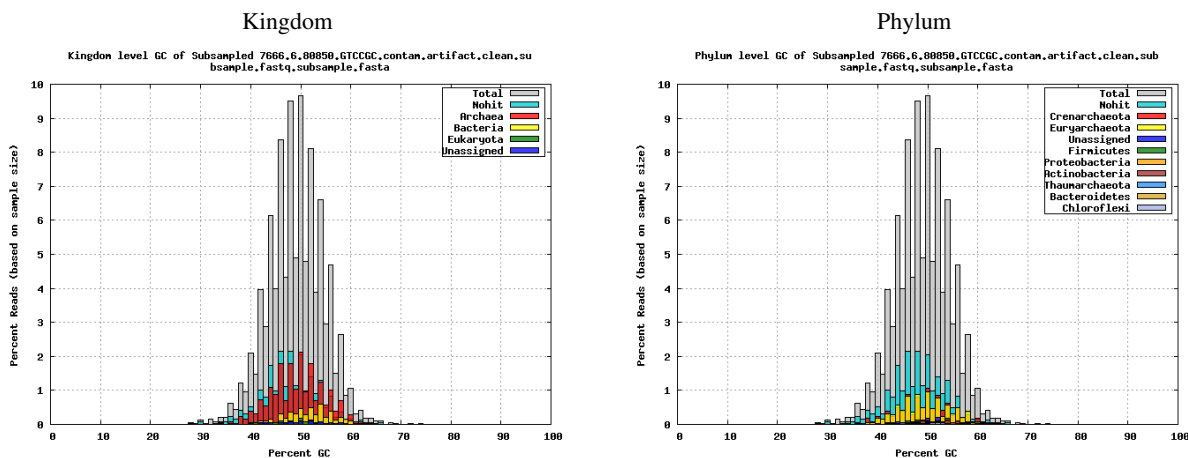
human_chr19	2	0.00
human_chr7	2	0.00
human_chr16	2	0.00
gi 357579531 Canis_lupus_familiaris_chr23	2	0.00
human_chr8	2	0.00
gi 357579550 Canis_lupus_familiaris_chr12	2	0.00
human_chr4	2	0.00
human_chr15	2	0.00

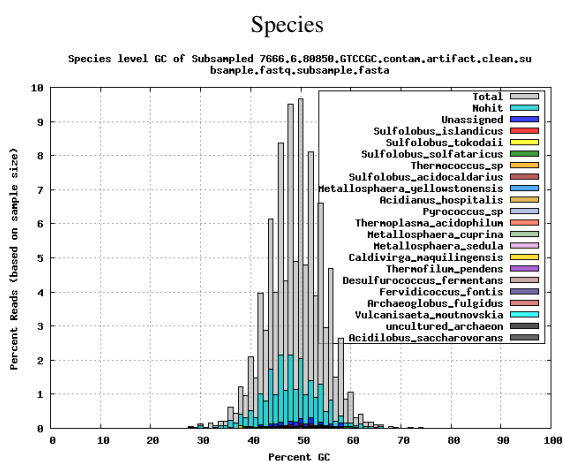
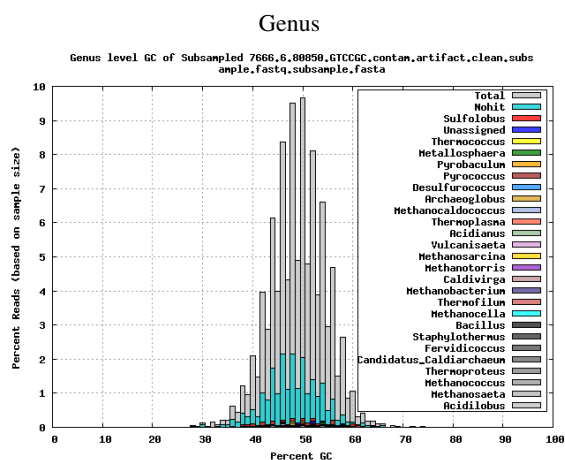
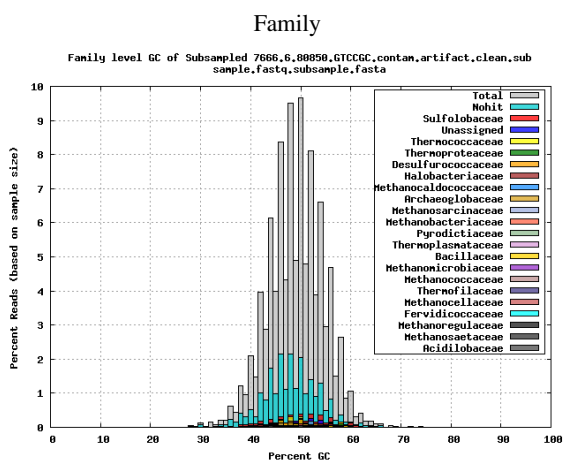
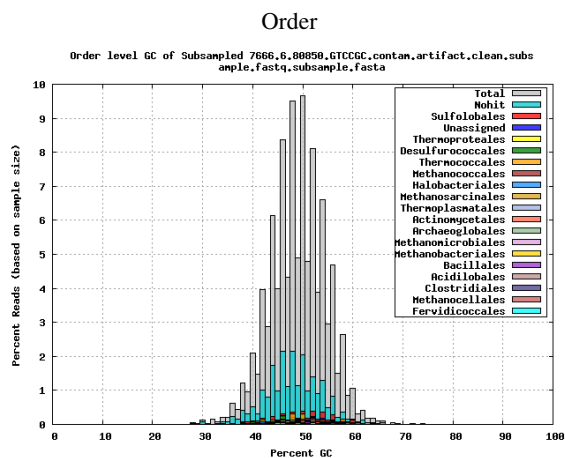
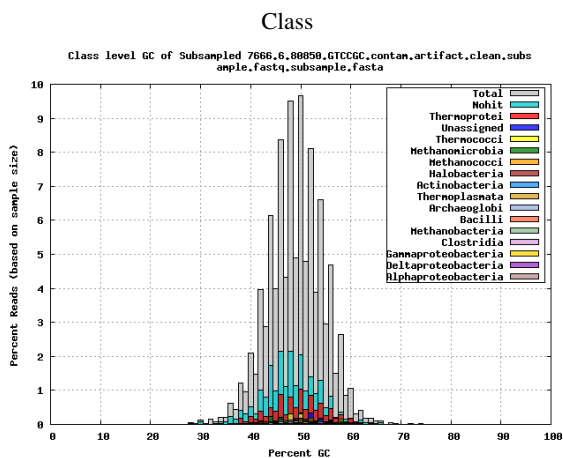
The following are the results of reads screened against potential reagent and process contaminants but were not removed from the dataset.

Illumina Std PE Contamination Identification Statistics

Description	Num Reads	Pct Reads
Input	24,831,480	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.



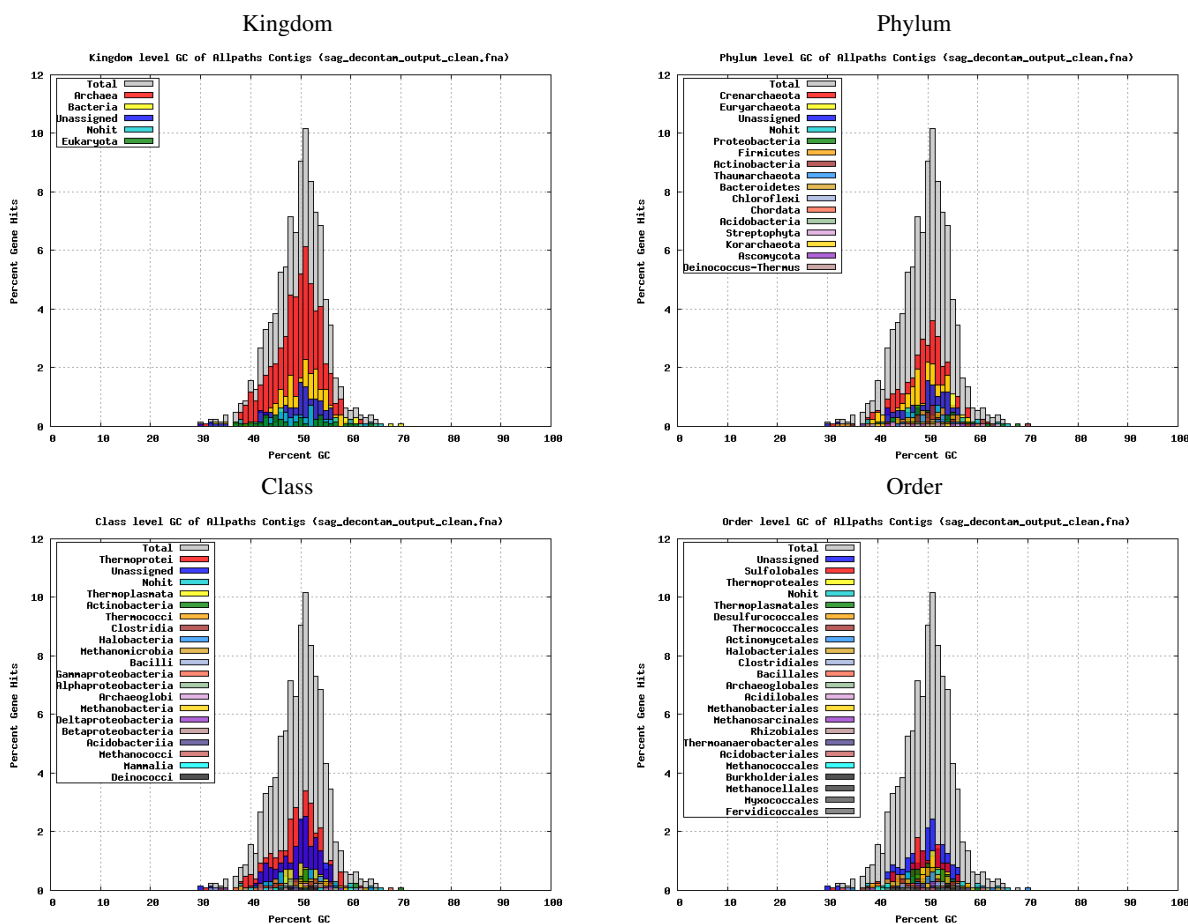


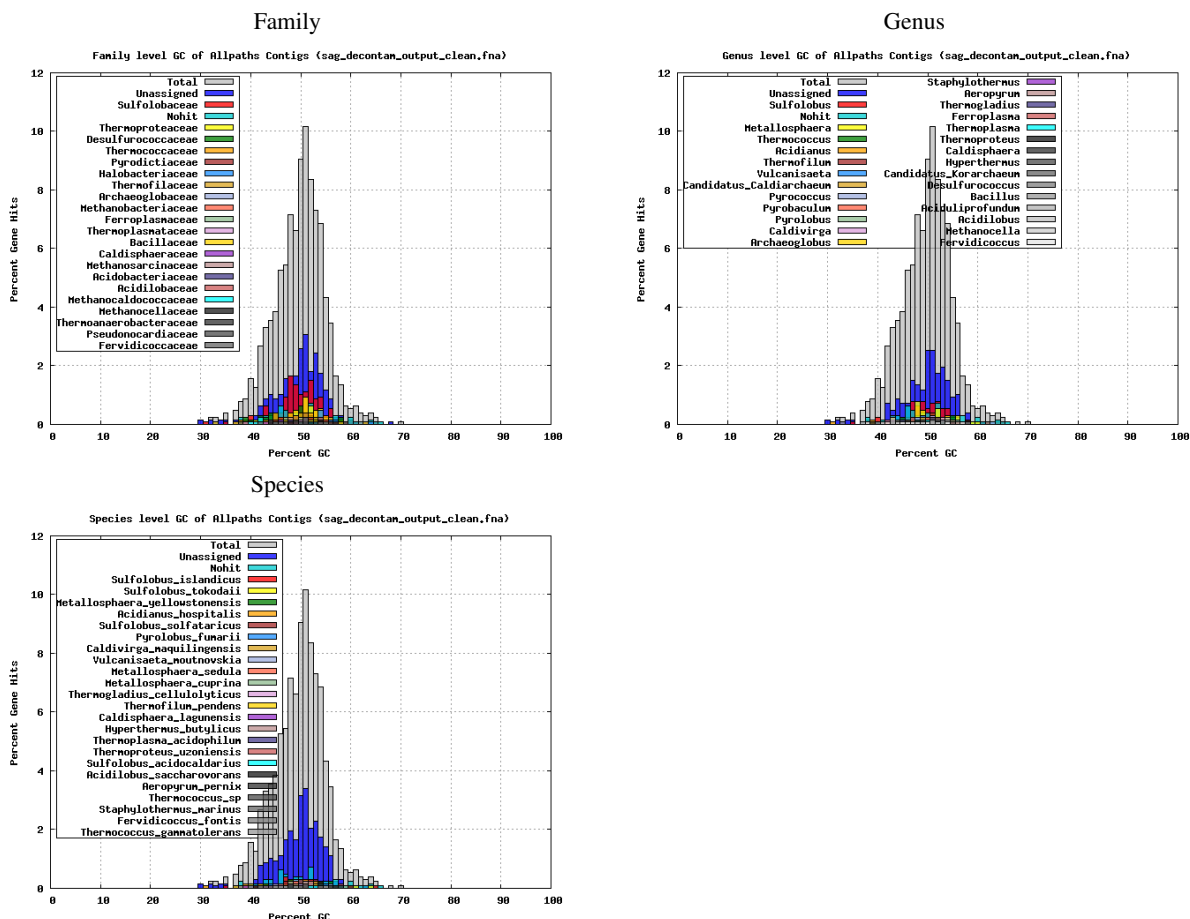
4. Assembly Statistics

Assembly method	SPAdes with auto decontamination
Scaffold total	111
Contig total	111
Scaffold sequence length	1.2 Mb
Contig sequence length	1.2 Mb (0.0% gap)
Scaffold N/L50	34/12.0 kb
Contig N/L50	34/12.0 kb
Largest Contig	40.3 kb
Number of scaffolds >50 kb	0
Pct of genome in scaffolds >50 kb	0.0
Pct of reads assembled (raw)	86.5
Pct of reads assembled (decontam)	84.2

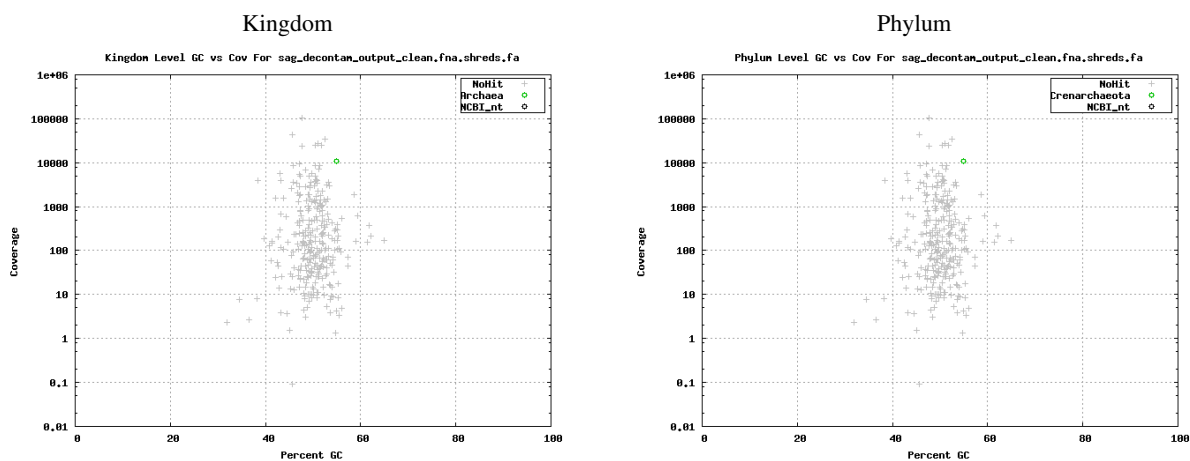
5. Assembly QC Results

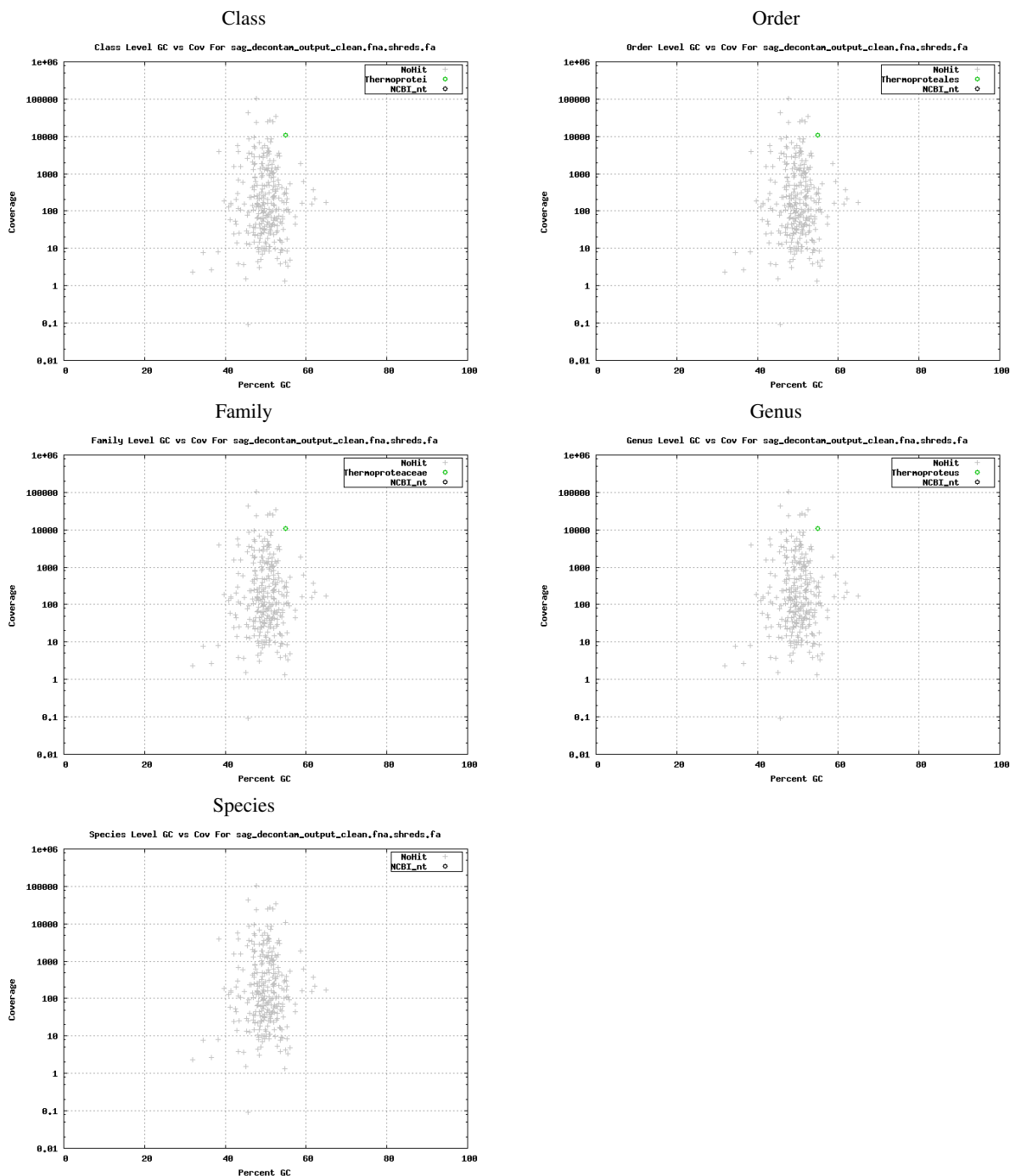
GC histogram of the predicted genes on each contig, overlaid with GC of hits based on BLASTP, shown for different taxonomic levels.



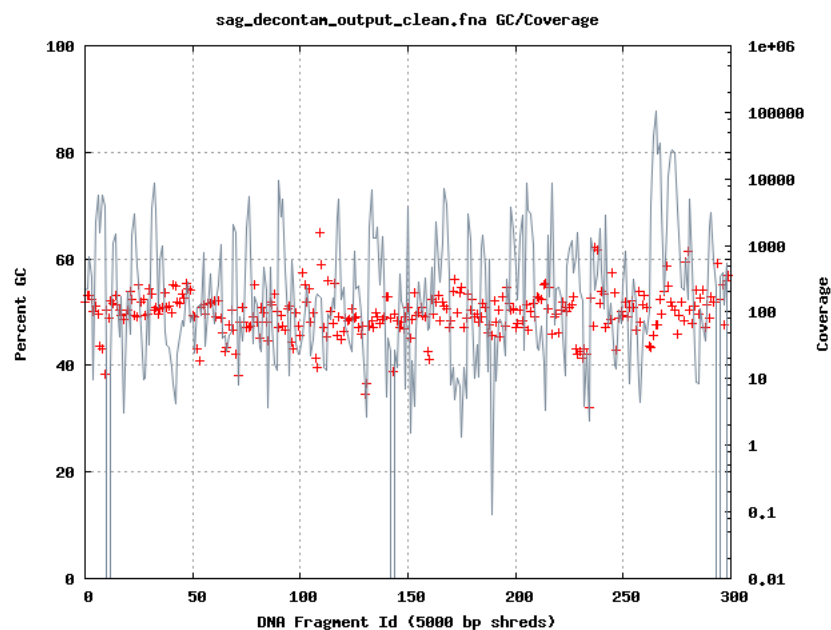


GC vs coverage based on GC of NCBI nt and Greengenes 16S rRNA gene hits to the assembly using megablast, shown for different taxonomic levels.

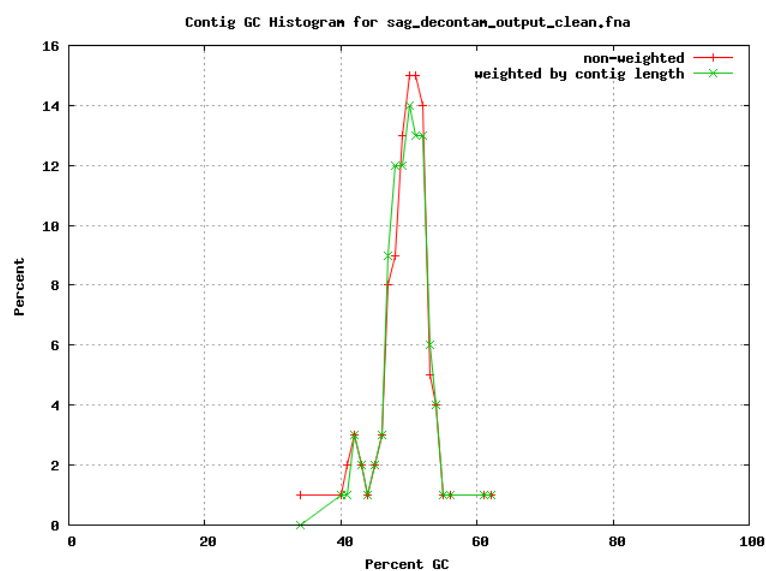




Coverage vs GC. Contigs were shredded into non-overlapping 5kbp and the GC of each shred was plotted as a point, colored by scaffold id. Coverage was calculated by mapping the fragment library to the final assembly and plotted as connected points.



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

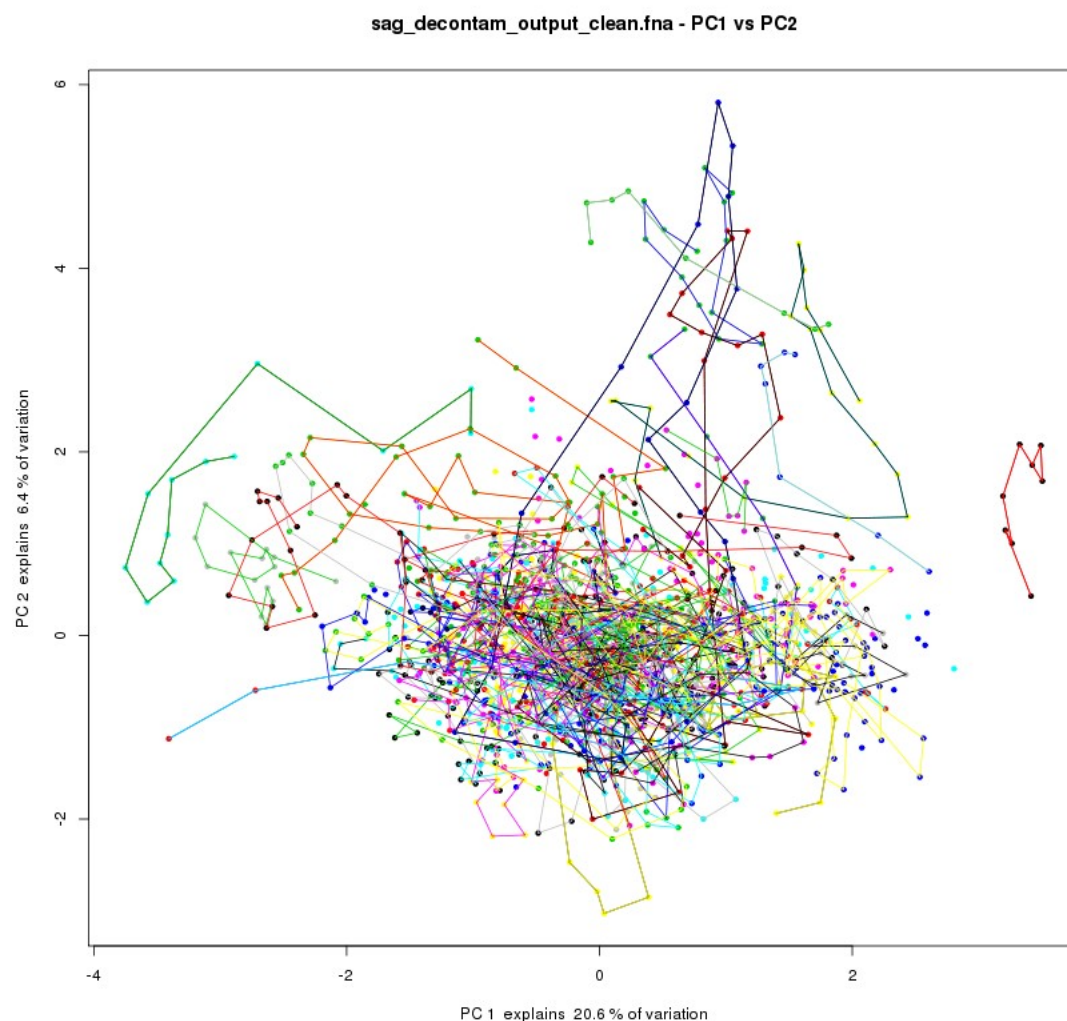


List of contigs and average percent GC, grouped in bins of 5:

Pct GC Bin	Contig Name
30	NODE_100.length_5576.cov_4.84659_ID.225
40	NODE_8.length_21297.cov_426.01_ID.15, NODE_35.length_11924.cov_27.093_ID.75, NODE_39.length_11168.cov_2535.54_ID.85, NODE_40.length_11128.cov_1125.81_ID.87, NODE_48.length_9892.cov_313.155_ID.103, NODE_49.length_9470.cov_38.544_ID.105, NODE_58.length_8342.cov_83.3798_ID.123, NODE_67.length_7837.cov_113.876_ID.141 NODE_103.length_5319.cov_15.7451_ID.231

45	<p> NODE_1.length.40345.cov.1260.75.ID.1, NODE_3.length.24839.cov.22707.ID.5, NODE_4.length.24094.cov.113.719.ID.7, NODE_5.length.23874.cov.2233.1.ID.9, NODE_7.length.21499.cov.222.921.ID.13, NODE_11.length.19823.cov.1351.86.ID.21, NODE_12.length.17866.cov.844.896.ID.23, NODE_13.length.17634.cov.3308.71.ID.25, NODE_16.length.16852.cov.629.835.ID.31, NODE_17.length.16835.cov.30.3779.ID.33, NODE_18.length.16767.cov.172.403.ID.35, NODE_19.length.16185.cov.1710.59.ID.37, NODE_21.length.14117.cov.143.833.ID.47, NODE_25.length.13369.cov.1320.14.ID.53, NODE_28.length.13156.cov.82.0706.ID.59, NODE_38.length.11183.cov.19.285.ID.79, NODE_43.length.10645.cov.69.1045.ID.93, NODE_45.length.10109.cov.233.512.ID.97, NODE_59.length.8223.cov.10.0109.ID.125, NODE_62.length.8153.cov.21.2736.ID.131, NODE_66.length.7884.cov.19.6832.ID.139, NODE_69.length.7435.cov.3809.97.ID.149, NODE_71.length.7235.cov.53.4305.ID.157, NODE_73.length.7097.cov.23.0263.ID.161, NODE_74.length.7090.cov.51.4263.ID.163, NODE_78.length.6959.cov.2502.06.ID.175, NODE_80.length.6775.cov.15.6839.ID.179, NODE_84.length.6453.cov.67.2737.ID.187, NODE_85.length.6441.cov.31.671.ID.151, NODE_89.length.6228.cov.198.843.ID.193, NODE_90.length.6221.cov.1508.81.ID.195, NODE_92.length.6025.cov.10.8469.ID.203, NODE_96.length.5822.cov.8.92509.ID.145, NODE_99.length.5725.cov.17.9557.ID.219, NODE_101.length.5557.cov.89.466.ID.227, NODE_106.length.5231.cov.382.4.ID.237, NODE_111.length.4824.cov.6.66765.ID.245, NODE_112.length.4686.cov.50.6435.ID.247 </p>
50	<p> NODE_2.length.29991.cov.9874.04.ID.3, NODE_6.length.23061.cov.100.435.ID.11, NODE_9.length.21257.cov.227.757.ID.17, NODE_10.length.20072.cov.239.241.ID.19, NODE_14.length.17503.cov.113.428.ID.27, NODE_15.length.17368.cov.92.0258.ID.29, NODE_22.length.13973.cov.668.494.ID.49, NODE_23.length.13882.cov.1179.57.ID.51, NODE_24.length.13754.cov.564.481.ID.43, NODE_26.length.13351.cov.278.023.ID.55, NODE_27.length.13289.cov.343.612.ID.57, NODE_29.length.12785.cov.3675.56.ID.61, NODE_30.length.12647.cov.142.489.ID.63, NODE_31.length.12543.cov.27.9301.ID.65, NODE_32.length.12504.cov.42.3613.ID.67, NODE_33.length.12184.cov.724.25.ID.69, NODE_34.length.12039.cov.78.9249.ID.73, NODE_36.length.11916.cov.50.5089.ID.77, NODE_37.length.11504.cov.38.0885.ID.83, NODE_41.length.10994.cov.1705.88.ID.89, NODE_42.length.10812.cov.109.512.ID.91, NODE_44.length.10256.cov.26.0733.ID.95, NODE_46.length.10103.cov.312.259.ID.99, NODE_47.length.10054.cov.22.2501.ID.101, NODE_50.length.9440.cov.539.223.ID.107, NODE_51.length.9391.cov.441.346.ID.109, NODE_52.length.9313.cov.1472.43.ID.111, NODE_53.length.9134.cov.13039.1.ID.113, NODE_54.length.9027.cov.191.995.ID.115, NODE_55.length.8544.cov.43.278.ID.117, NODE_57.length.8488.cov.1034.41.ID.121, NODE_61.length.8174.cov.221.596.ID.129, NODE_63.length.8120.cov.506.251.ID.133, NODE_64.length.8102.cov.11.2022.ID.135, NODE_65.length.8065.cov.320.403.ID.137, NODE_68.length.7448.cov.107.043.ID.147, NODE_70.length.7256.cov.211.916.ID.155, NODE_72.length.7176.cov.98.4957.ID.159, NODE_75.length.7050.cov.628.018.ID.169, NODE_76.length.7013.cov.174.748.ID.171, NODE_77.length.7008.cov.34.6935.ID.173, NODE_79.length.6859.cov.59.7551.ID.177, NODE_81.length.6609.cov.4343.17.ID.181, NODE_82.length.6482.cov.818.335.ID.183, NODE_83.length.6469.cov.43.2323.ID.185, NODE_87.length.6275.cov.19.4617.ID.165, NODE_88.length.6256.cov.2090.73.ID.191, NODE_91.length.6088.cov.8.9463.ID.201, NODE_94.length.6007.cov.24.2566.ID.207, NODE_95.length.5826.cov.5.72102.ID.209, NODE_97.length.5808.cov.18.8356.ID.211, NODE_98.length.5771.cov.15.9365.ID.217, NODE_102.length.5395.cov.11.5569.ID.229, NODE_104.length.5257.cov.35.1738.ID.233, NODE_105.length.5245.cov.4.90424.ID.235, NODE_107.length.5195.cov.2520.47.ID.239, NODE_108.length.5179.cov.30.5859.ID.197, NODE_109.length.5041.cov.12.2302.ID.241 NODE_110.length.4832.cov.9.69269.ID.243 </p>
55	<p> NODE_20.length.14458.cov.224.031.ID.41, NODE_93.length.6013.cov.5.15089.ID.205 </p>
60	<p> NODE_56.length.8519.cov.191.491.ID.119, NODE_60.length.8218.cov.109.393.ID.127 </p>

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	31.18 %
archaea	74.76 %

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 24,831,480 reads totaling 3,724.7 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 111 contigs in 111 scaffolds, totalling 1.2 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.