

1. Project Information

Program	Microbial/CSP 2012
PMO Project	0
Seq Proj ID	1027073
Sequencing Project Name	Geoarchaeota archaeon JGI 000156CP-M9
JGI Project ID	0

2. Read Statistics

Illumina Std PE Statistics

File name	7667.5.80864.AGTTCC.fastq
Library	TGSP
Number of reads	25,732,728
Sequencing depth [†]	772X
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	25,732,728	100
Contam removed	76	0.0
Artifact removed	770,354	3.0
Total removed	5,732,728	22.3
Total remaining	20,000,000	77.7

List of Contaminants Removed

Description	Num Reads	Pct Reads
gi 357579577 Canis_lupus_familiaris_chr3	52	0.00
human_chr2	46	0.00
gi 357579535 Canis_lupus_familiaris_chr20	20	0.00
human_chr16	2	0.00
human_chr3	2	0.00
gi 357579571 Canis_lupus_familiaris_chr5	2	0.00

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

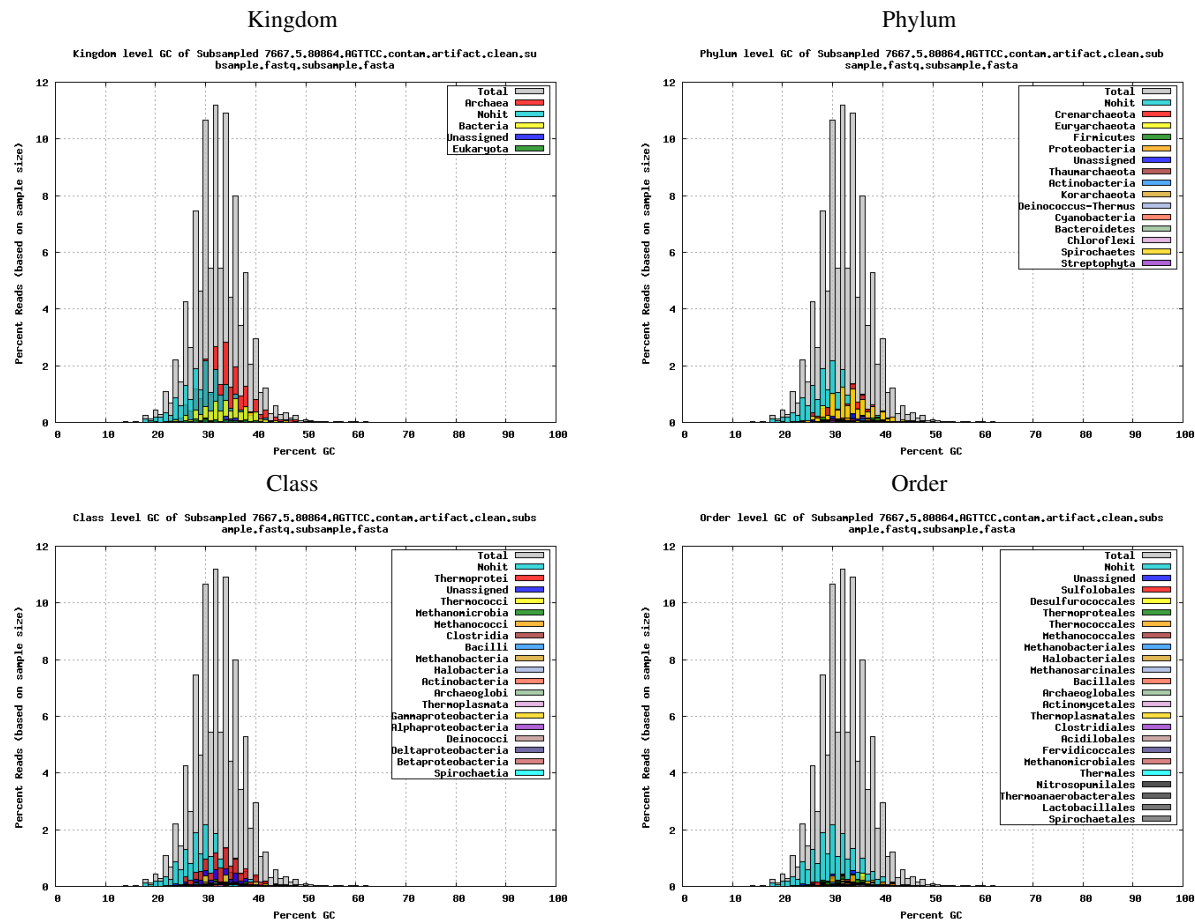
Illumina Std PE Contamination Identification Statistics

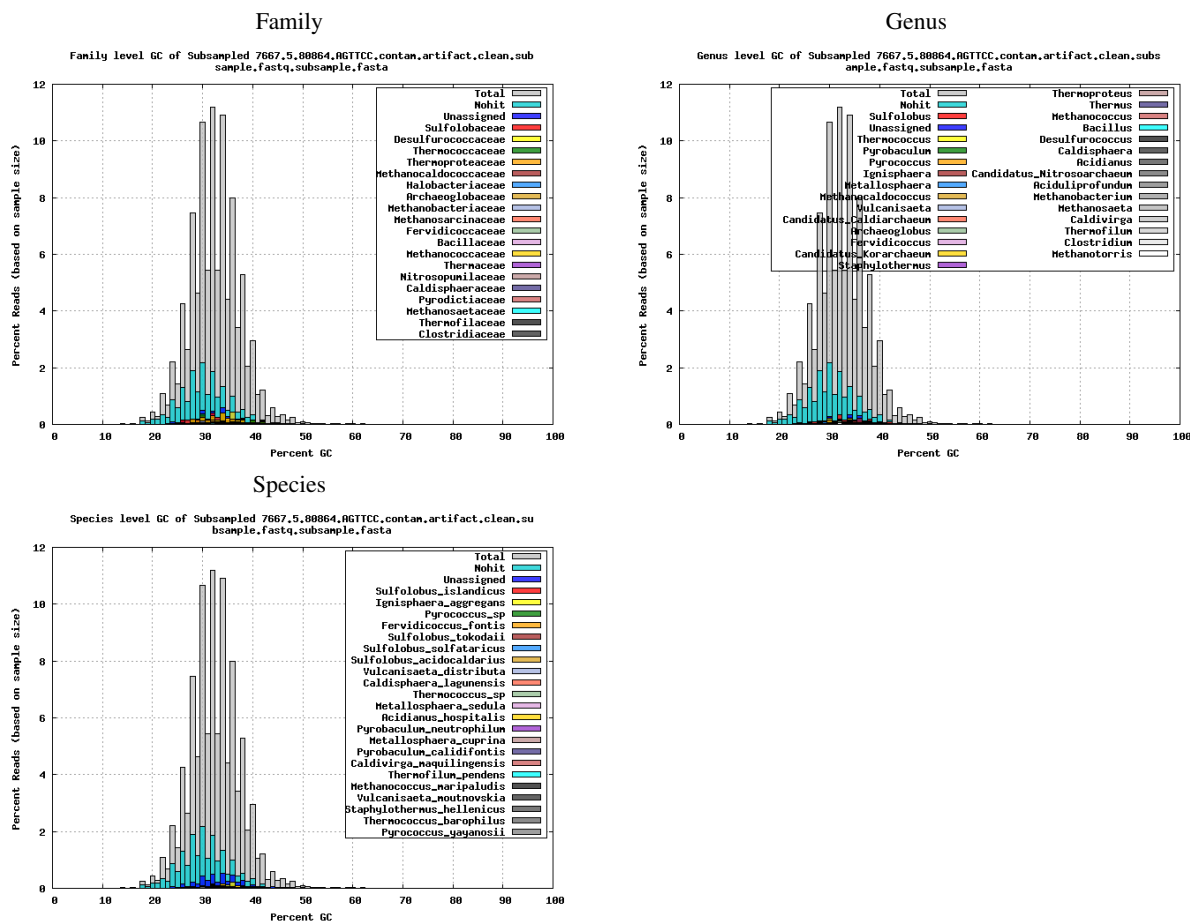
Description	Num Reads	Pct Reads
Input	25,732,728	100
Contam identified	6	0.0

List of Contaminants Identified

Description	Num Reads	Pct Reads
<i>Delftia</i>	2	0.00
<i>Pseudomonas</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.





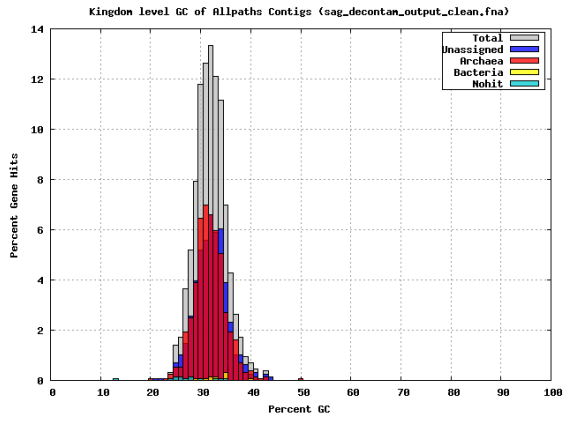
4. Assembly Statistics

Assembly method	SPAdes with auto decontamination
Scaffold total	57
Contig total	57
Scaffold sequence length	1.1 Mb
Contig sequence length	1.1 Mb (0.0% gap)
Scaffold N/L50	8/36.1 kb
Contig N/L50	8/36.1 kb
Largest Contig	145.4 kb
Number of scaffolds >50 kb	4
Pct of genome in scaffolds >50 kb	38.2
Pct of reads assembled (raw)	93.4
Pct of reads assembled (decontam)	92.7

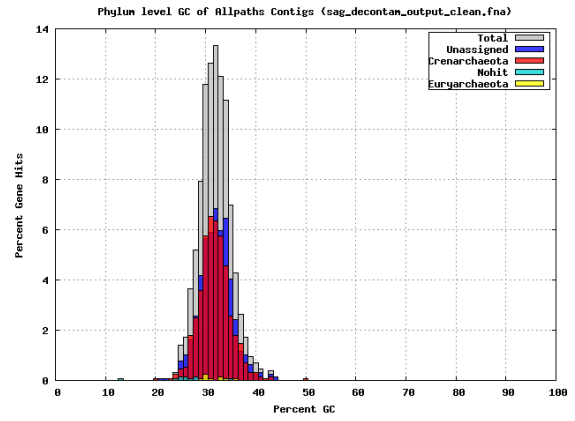
5. Assembly QC Results

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

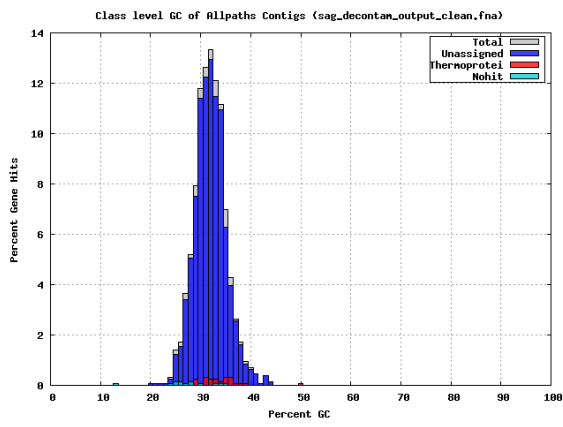
Kingdom



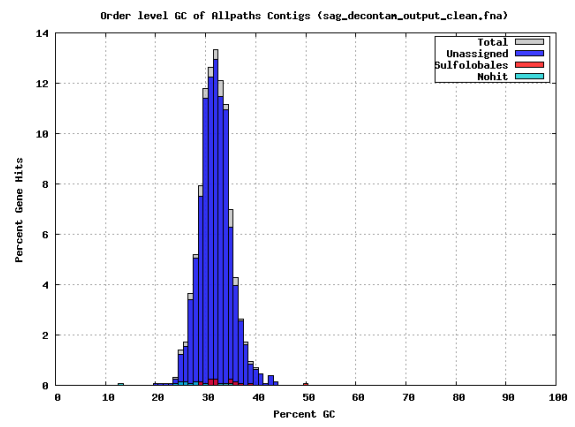
Phylum



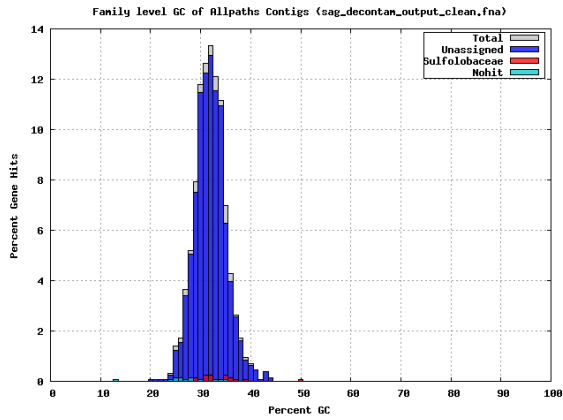
Class



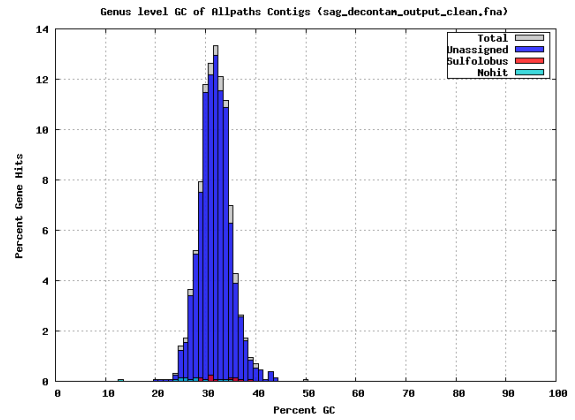
Order



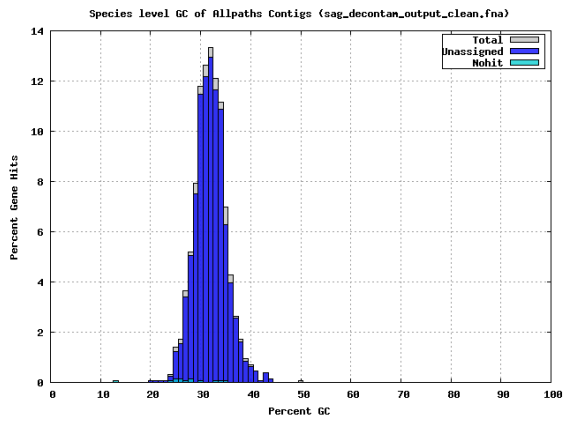
Family



Genus

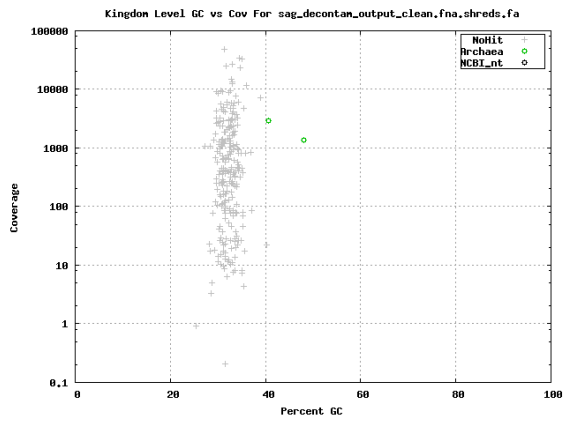


Species

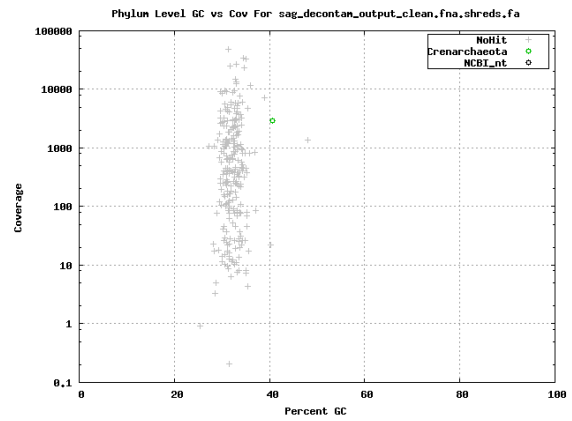


GC vs coverage based on GC of NCBI nt and Greengenes 16S rRNA gene hits to the assembly using megablast, shown for different taxonomic levels.

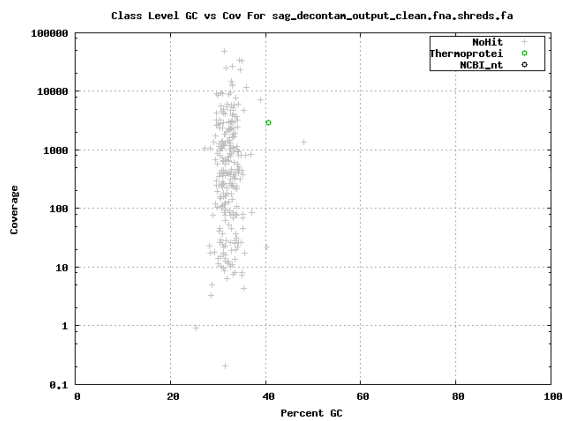
Kingdom



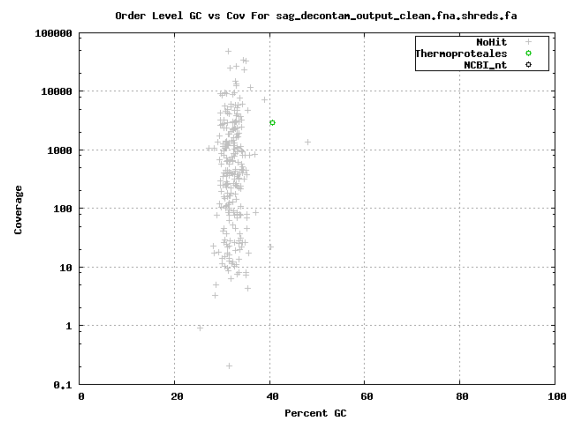
Phylum

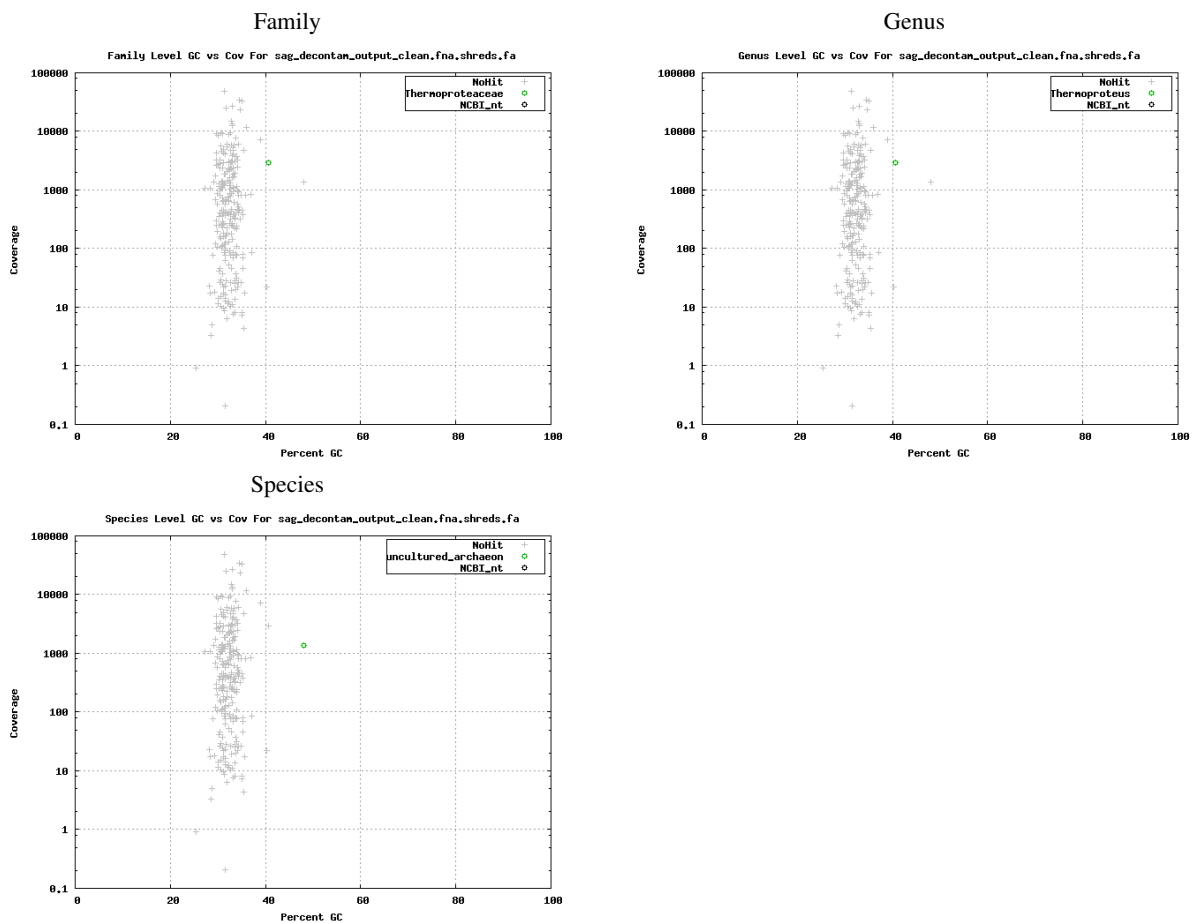


Class

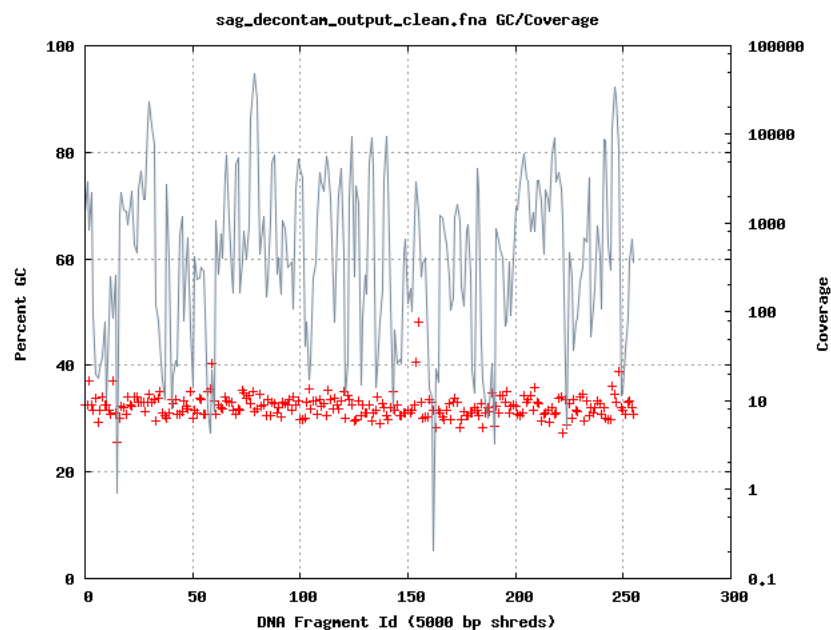


Order

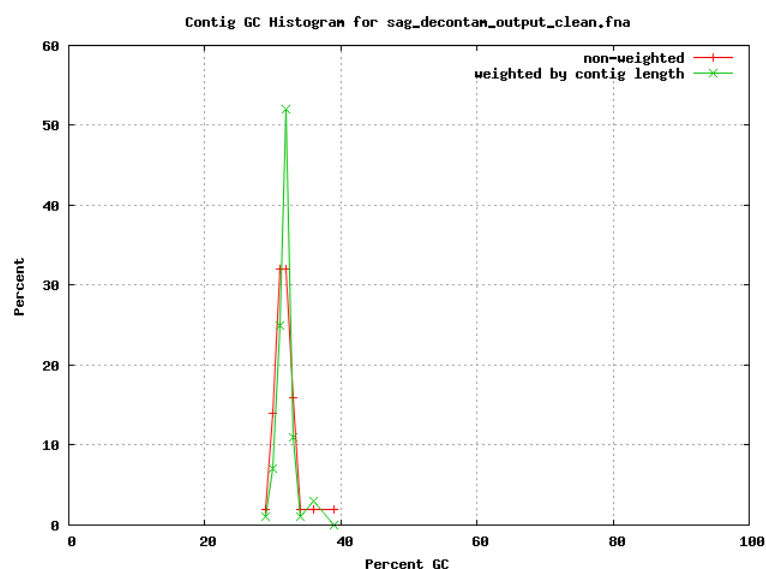




Coverage vs GC. Contigs were shredded into non-overlapping 5kbp and the GC of each shred was plotted as a point, colored by scaffold id. Coverage was calculated by mapping the fragment library to the final assembly and plotted as connected points.



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

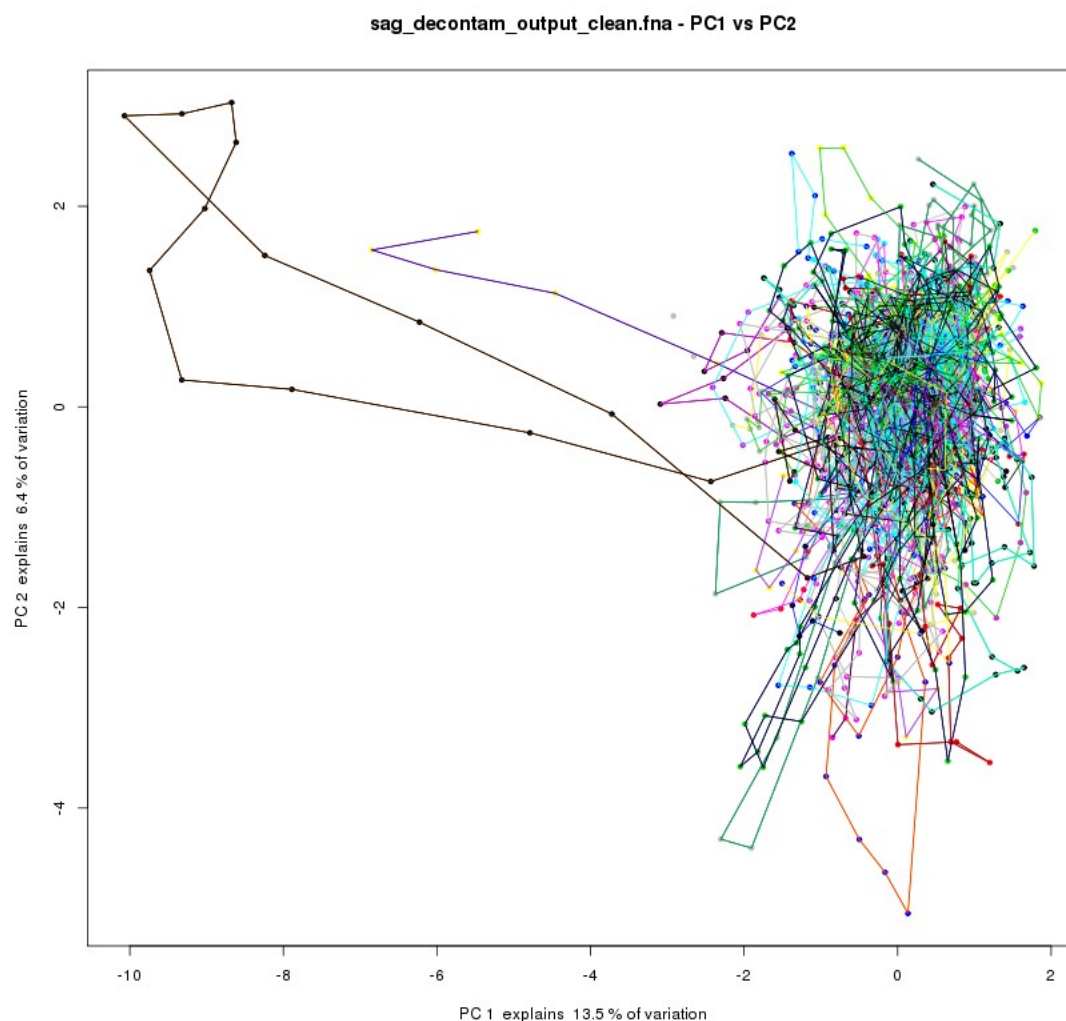


List of contigs and average percent GC, grouped in bins of 5:

Pct GC Bin	Contig Name
25	NODE.40.length.7868.cov.13.9345.ID.79
30	NODE.1.length.145385.cov.1529.69.ID.1, NODE.2.length.139752.cov.3662.9.ID.3, NODE.3.length.84587.cov.421.012.ID.5, NODE.4.length.65284.cov.1575.28.ID.7, NODE.5.length.40584.cov.4981.28.ID.9, NODE.6.length.40512.cov.138.222.ID.11, NODE.7.length.38412.cov.379.268.ID.13, NODE.8.length.36080.cov.1062.54.ID.15, NODE.10.length.22578.cov.2280.4.ID.19, NODE.11.length.22027.cov.201.773.ID.21, NODE.12.length.20289.cov.2654.34.ID.23, NODE.13.length.20251.cov.12053.5.ID.25, NODE.14.length.18691.cov.293.633.ID.27, NODE.15.length.18195.cov.460.549.ID.29,

	NODE_16.length_16004.cov_169.318.ID_31, NODE_17.length_14798.cov_113.197.ID_33, NODE_18.length_14542.cov_757.081.ID_35, NODE_19.length_14361.cov_2357.58.ID_37, NODE_20.length_14302.cov_273.771.ID_39, NODE_21.length_13899.cov_2446.5.ID_41, NODE_22.length_13510.cov_500.04.ID_43, NODE_23.length_13462.cov_4425.89.ID_45, NODE_24.length_13107.cov_1043.68.ID_47, NODE_25.length_12940.cov_313.313.ID_49, NODE_26.length_12049.cov_4460.63.ID_51, NODE_27.length_11576.cov_16.471.ID_53, NODE_28.length_11572.cov_891.775.ID_55, NODE_29.length_11324.cov_106.012.ID_57, NODE_30.length_11148.cov_19.88.ID_59, NODE_31.length_10304.cov_11.4155.ID_61, NODE_32.length_9525.cov_611.115.ID_63, NODE_33.length_9461.cov_249.489.ID_65, NODE_34.length_8905.cov_24.5447.ID_67, NODE_35.length_8707.cov_82.2549.ID_69, NODE_36.length_8662.cov_53.1031.ID_71, NODE_37.length_8425.cov_2321.82.ID_73, NODE_38.length_8431.cov_1338.77.ID_75, NODE_39.length_8205.cov_14.8694.ID_77, NODE_41.length_7793.cov_5.88964.ID_81, NODE_42.length_7726.cov_7.80589.ID_83, NODE_43.length_7676.cov_82.6294.ID_85, NODE_44.length_7517.cov_18.4648.ID_87, NODE_45.length_7100.cov_44.5568.ID_89, NODE_46.length_7092.cov_22.5957.ID_91, NODE_47.length_6722.cov_186.57.ID_93, NODE_48.length_6737.cov_26.7326.ID_95, NODE_49.length_6635.cov_11.116.ID_97, NODE_50.length_5899.cov_38.6894.ID_99, NODE_51.length_5855.cov_49.9991.ID_101, NODE_53.length_5640.cov_177.767.ID_105, NODE_54.length_5155.cov_7.25549.ID_107, NODE_55.length_5106.cov_168.797.ID_109, NODE_56.length_4828.cov_10.2407.ID_111, NODE_59.length_4692.cov_52.1037.ID_117
35	NODE_9.length_29996.cov_652.743.ID_17, NODE_52.length_5680.cov_13.6411.ID_103

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	52.76 %
archaea	95.34 %

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 25,732,728 reads totaling 3,859.9 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 57 contigs in 57 scaffolds, totalling 1.1 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.