

1. Project Information

Program	Microbial/CSP 2012
PMO Project	0
Seq Proj ID	1027100
Sequencing Project Name	Leptolyngbya sp. HL7711_P1F1 JGI 000148CP-K16
JGI Project ID	0

2. Read Statistics

Illumina Std PE Statistics

File name	7666.6.80850.GAGTGG.fastq
Library	TGTB
Number of reads	24,239,914
Sequencing depth [†]	727X
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	24,239,914	100
Contam removed	74	0.0
Artifact removed	1,098,826	4.5
Total removed	4,239,914	17.5
Total remaining	20,000,000	82.5

List of Contaminants Removed

Description	Num Reads	Pct Reads
gi 357579577 Canis_lupus_familiaris_chr3	58	0.00
human_chr2	56	0.00
gi 357579535 Canis_lupus_familiaris_chr20	12	0.00
gi 357579571 Canis_lupus_familiaris_chr5	6	0.00
human_chr14	2	0.00
human_chr8	2	0.00

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

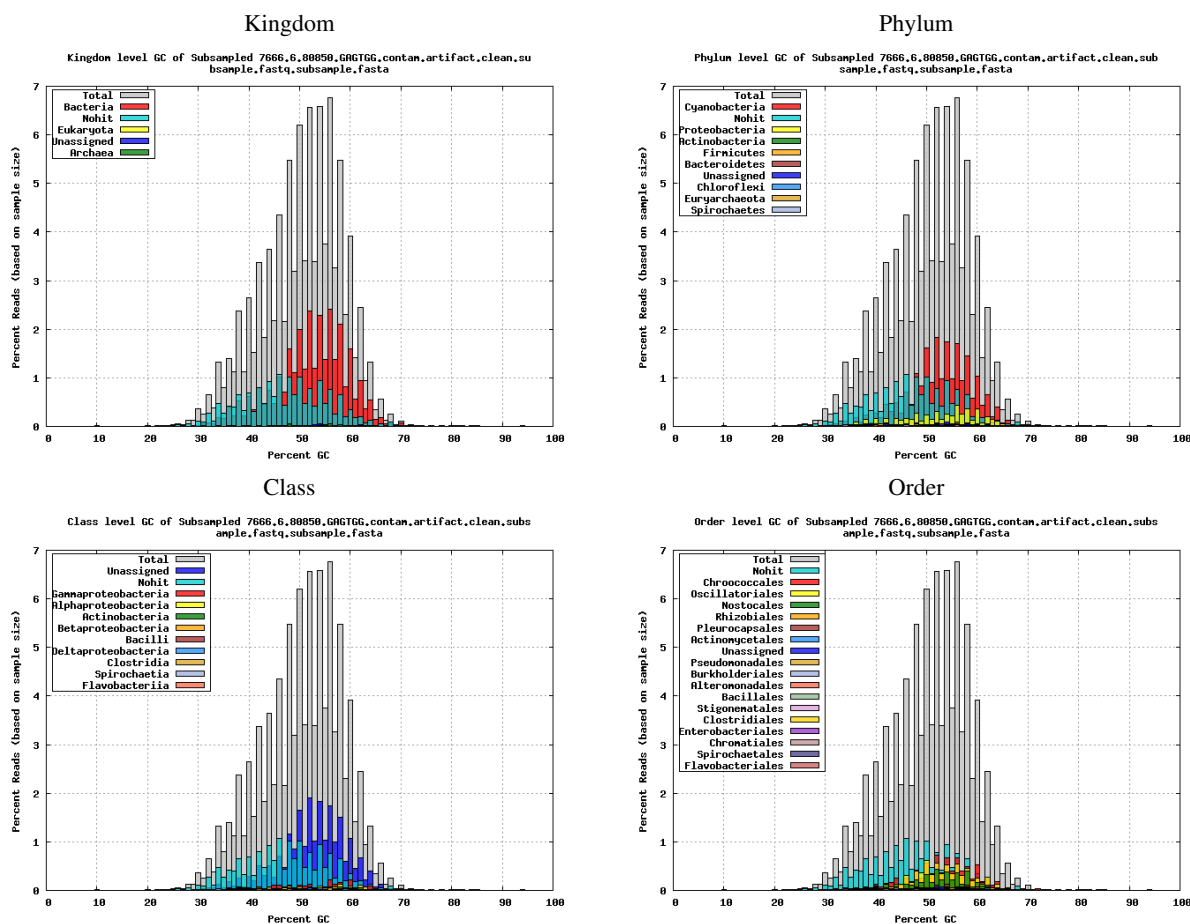
Illumina Std PE Contamination Identification Statistics

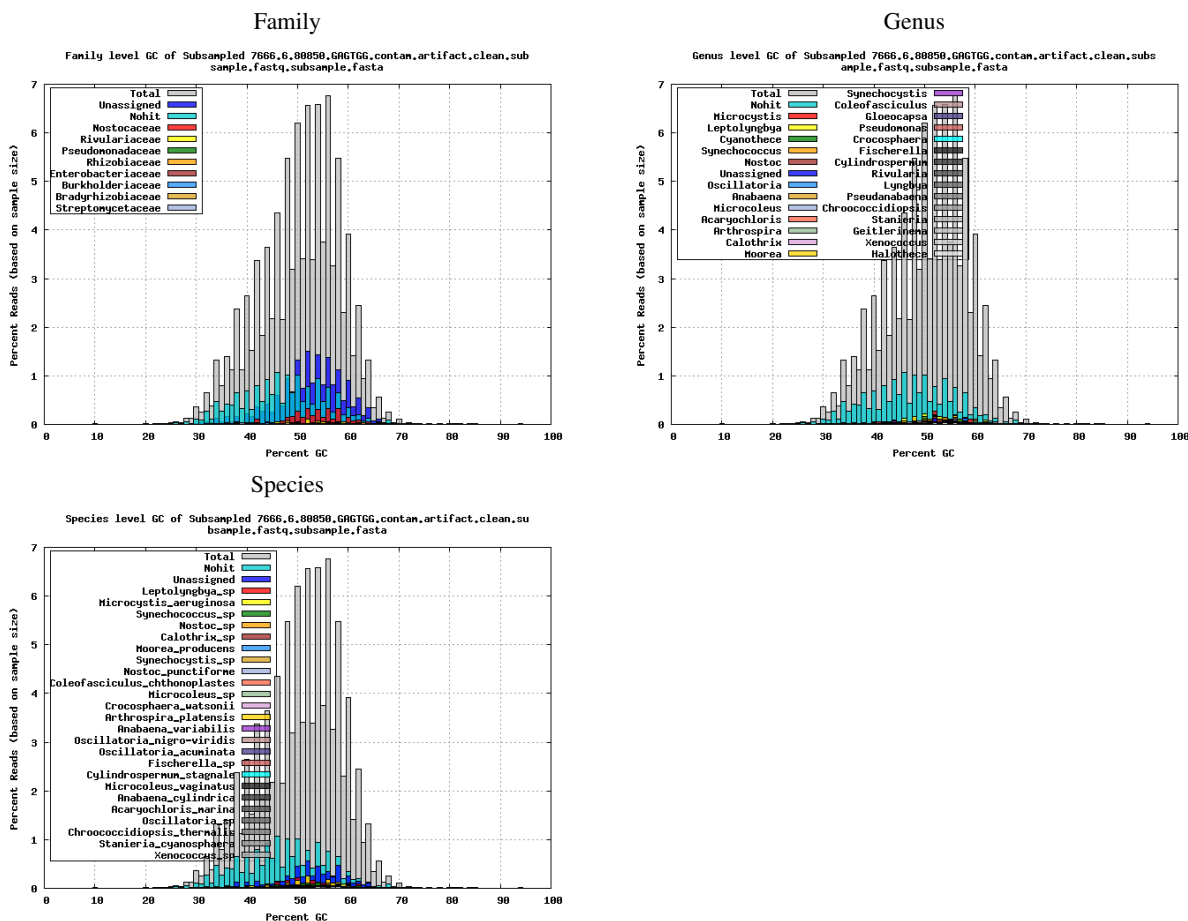
Description	Num Reads	Pct Reads
Input	24,239,914	100
Contam identified	14	0.0

List of Contaminants Identified

Description	Num Reads	Pct Reads
<i>Pseudomonas</i>	10	0.00
<i>Delftia</i>	2	0.00
<i>Ralstonia</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.





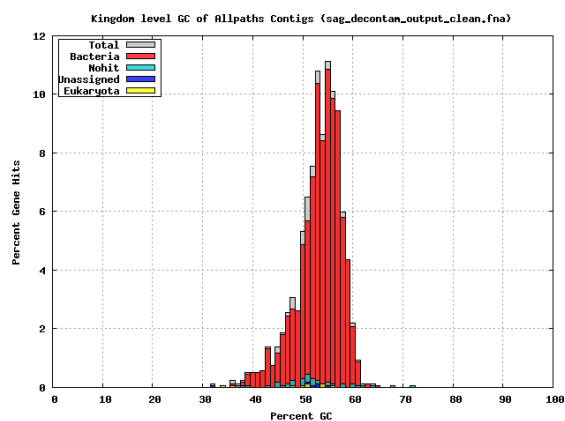
4. Assembly Statistics

Assembly method	SPAdes with auto decontamination
Scaffold total	72
Contig total	72
Scaffold sequence length	1.8 Mb
Contig sequence length	1.8 Mb (0.0% gap)
Scaffold N/L50	16/39.9 kb
Contig N/L50	16/39.9 kb
Largest Contig	92.1 kb
Number of scaffolds >50 kb	7
Pct of genome in scaffolds >50 kb	28.3
Pct of reads assembled (raw)	87.5
Pct of reads assembled (decontam)	29.6

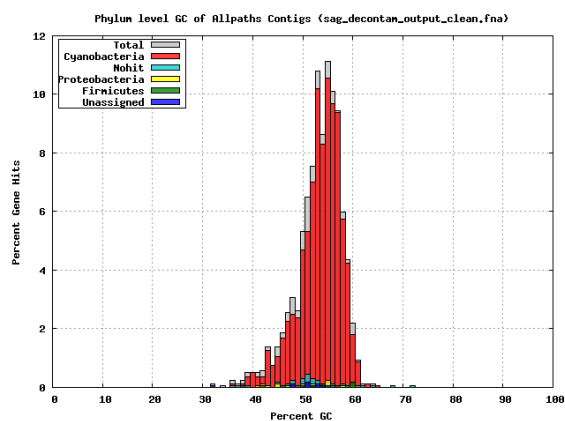
5. Assembly QC Results

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

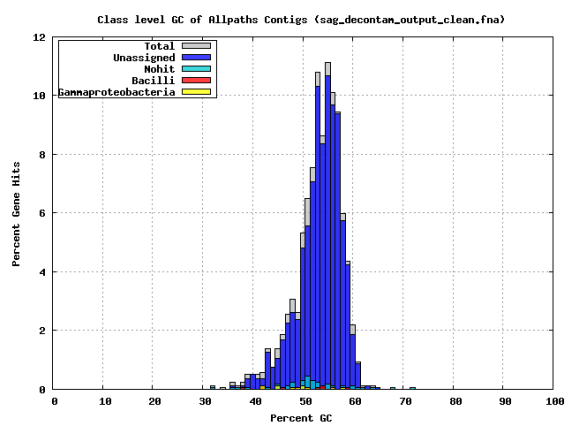
Kingdom



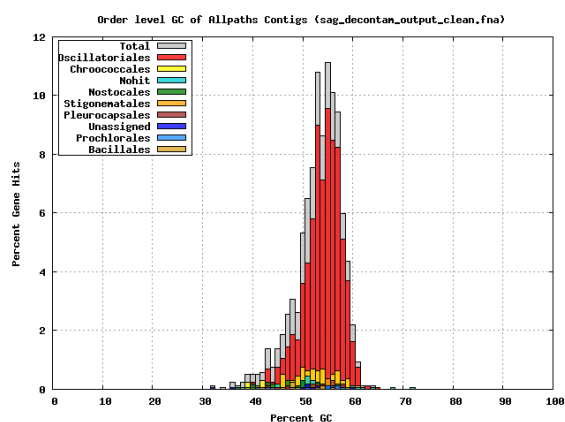
Phylum



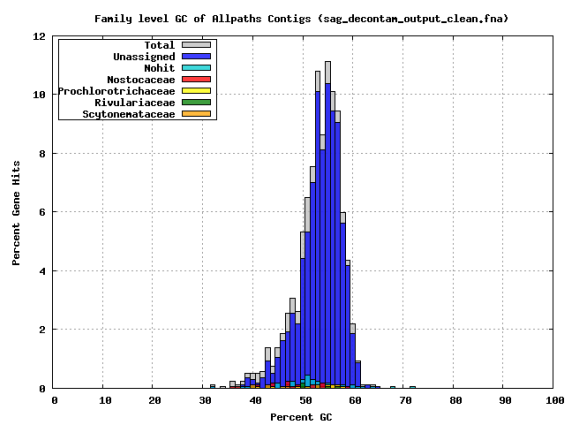
Class



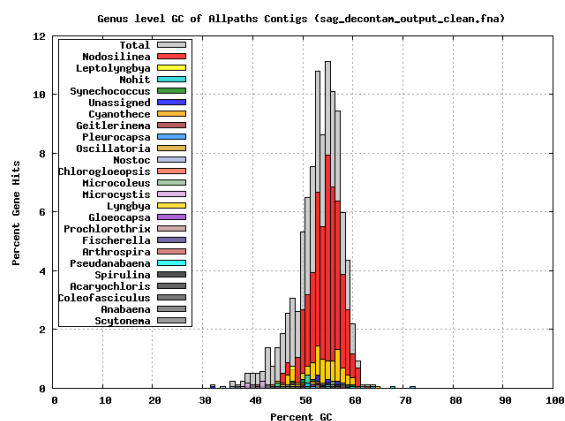
Order



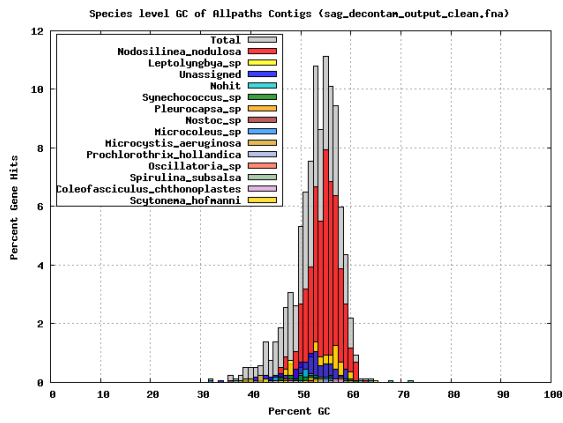
Family



Genus

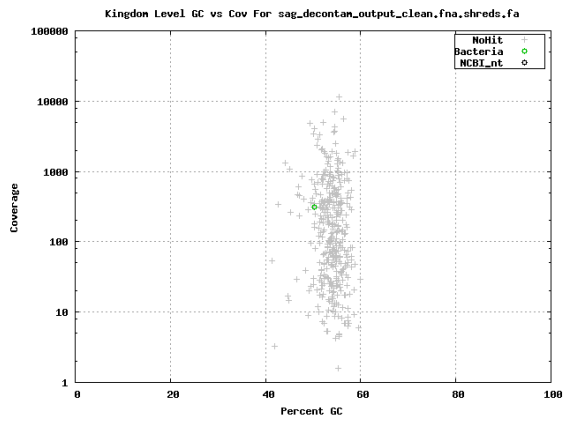


Species

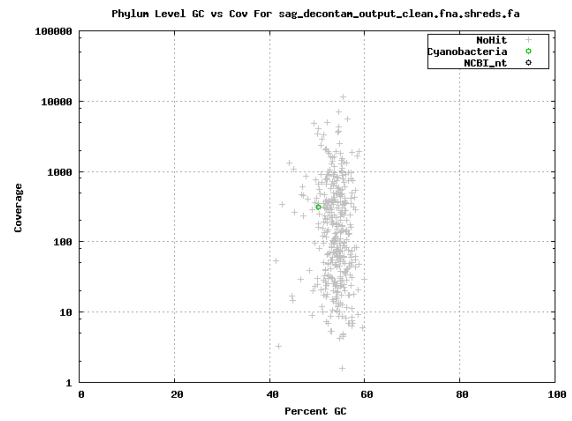


GC vs coverage based on GC of NCBI nt and Greengenes 16S rRNA gene hits to the assembly using megablast, shown for different taxonomic levels.

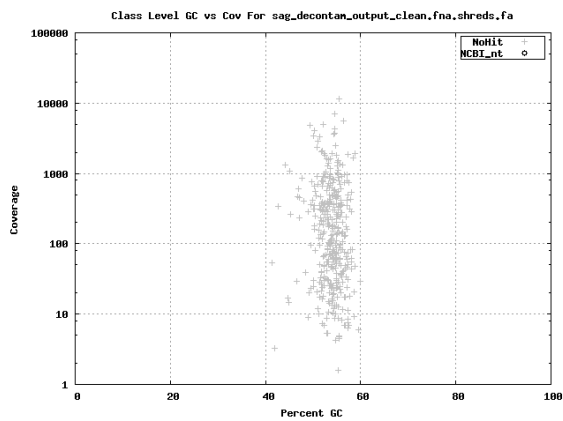
Kingdom



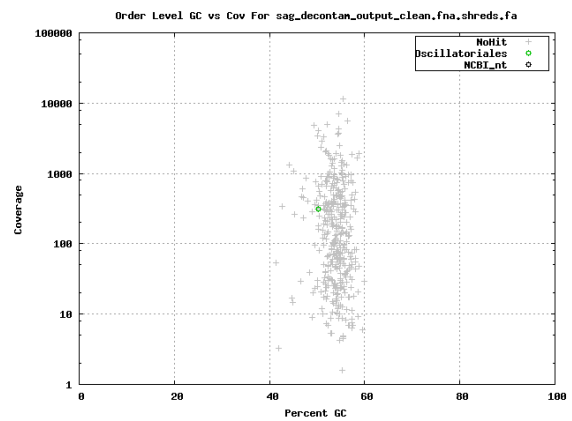
Phylum

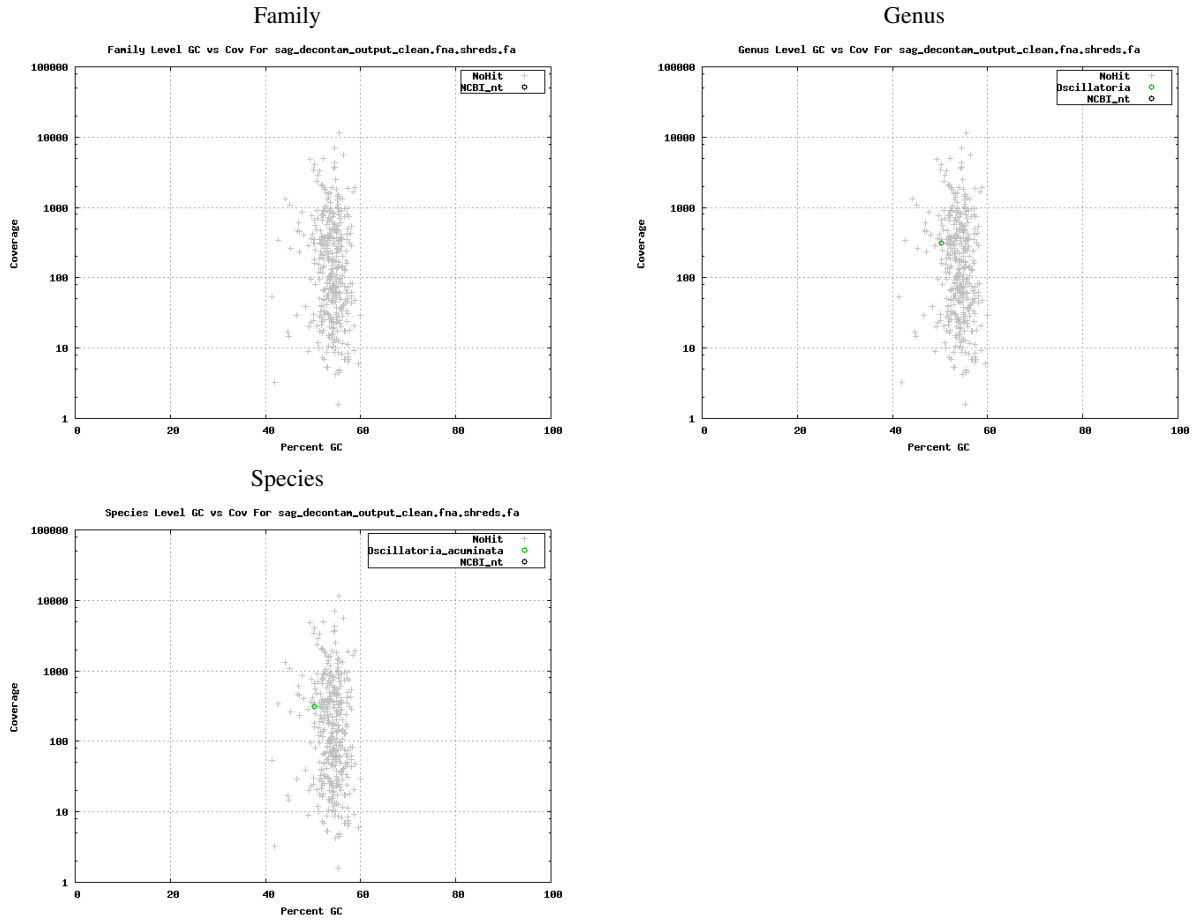


Class

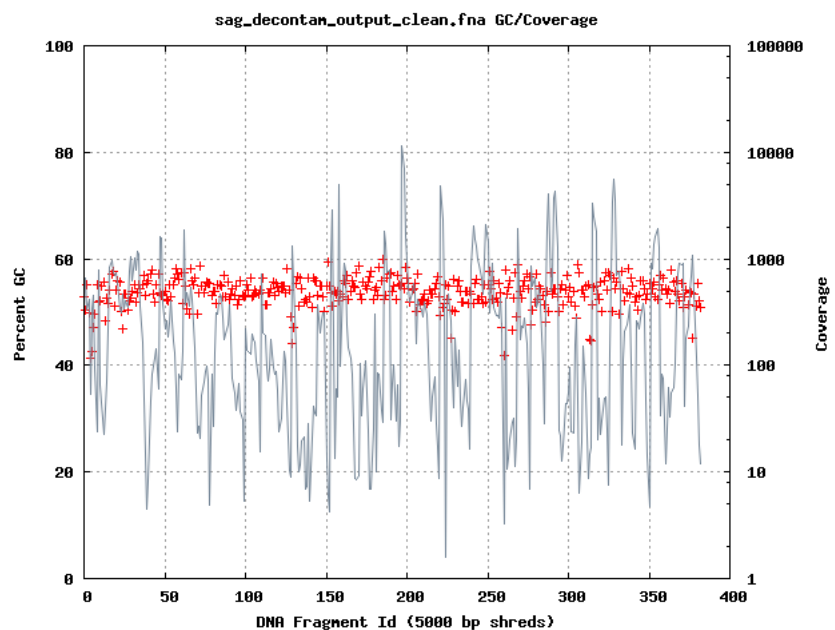


Order

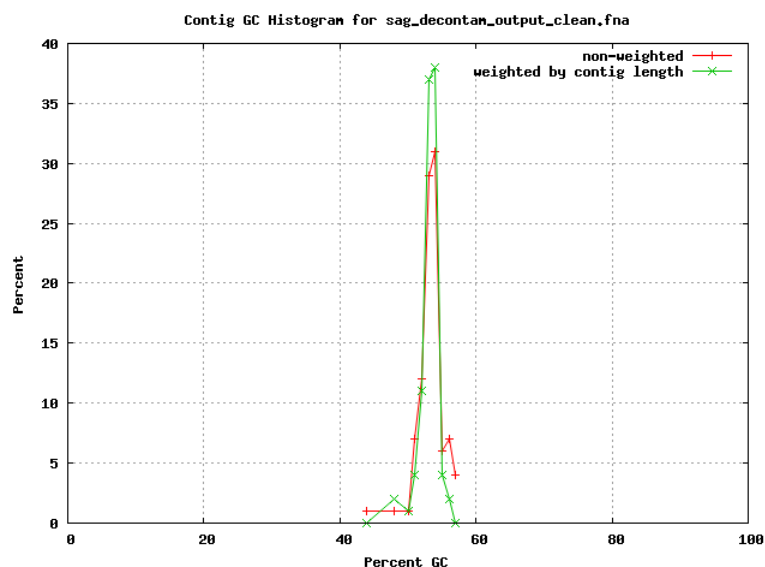




Coverage vs GC. Contigs were shredded into non-overlapping 5kbp and the GC of each shred was plotted as a point, colored by scaffold id. Coverage was calculated by mapping the fragment library to the final assembly and plotted as connected points.



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

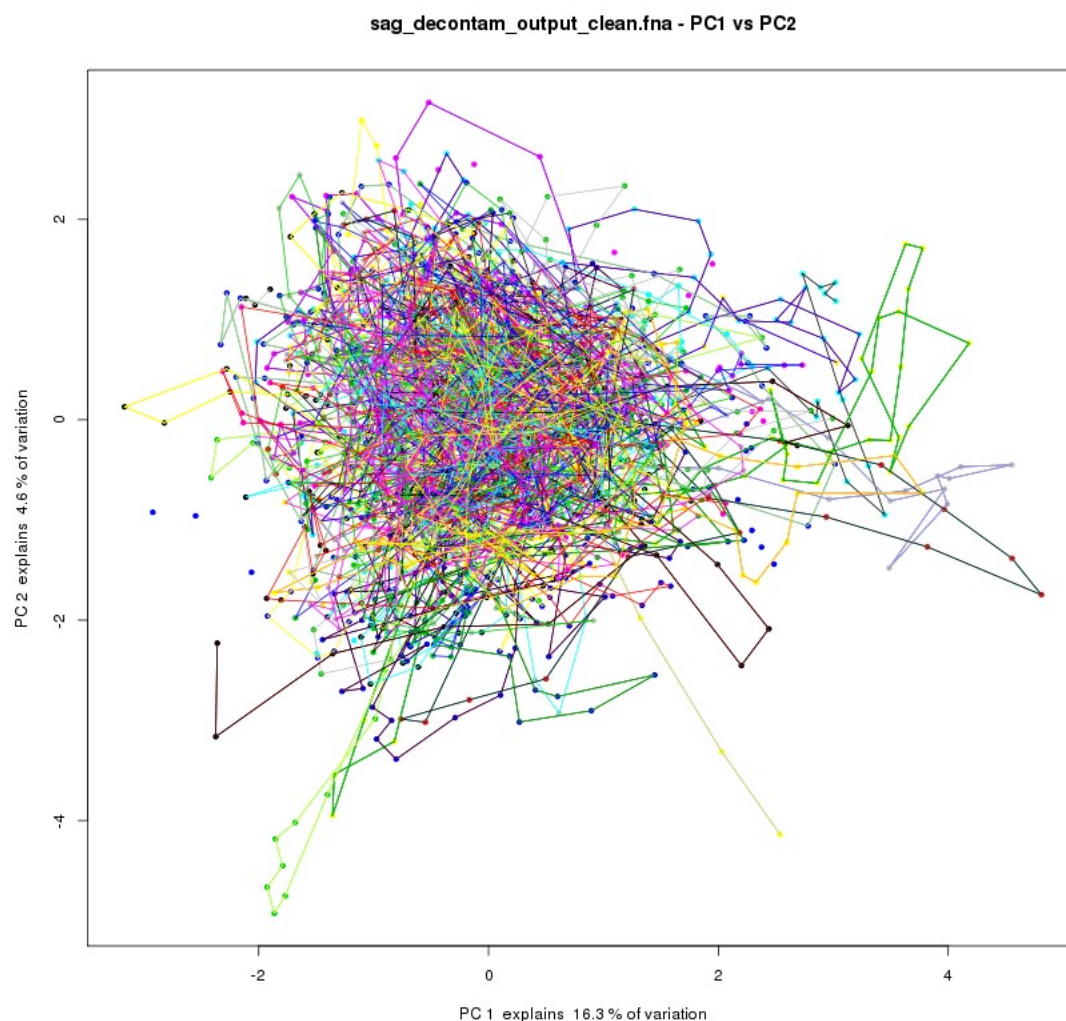


List of contigs and average percent GC, grouped in bins of 5:

Pct GC Bin	Contig Name
40	NODE_124.length_8358.cov_10.7063.ID_259
45	NODE_24.length_40479.cov_226.901.ID_47
50	NODE_2.length_92113.cov_209.769.ID_3, NODE_4.length_89700.cov_132.931.ID_7, NODE_6.length_72132.cov_390.16.ID_11, NODE_7.length_71107.cov_172.907.ID_13, NODE_8.length_67016.cov_127.74.ID_15, NODE_12.length_53043.cov_215.542.ID_23, NODE_14.length_51777.cov_799.922.ID_27, NODE_15.length_49192.cov_1616.76.ID_29, NODE_17.length_45699.cov_268.435.ID_33, NODE_18.length_45603.cov_389.236.ID_35, NODE_19.length_44335.cov_297.564.ID_37, NODE_20.length_44023.cov_41.6355.ID_39,

	<p> NODE.21.length.44008.cov.594.626.ID.41, NODE.23.length.42096.cov.37.8432.ID.45, NODE.25.length.39918.cov.158.679.ID.49, NODE.26.length.38643.cov.867.108.ID.51, NODE.28.length.37698.cov.384.578.ID.55, NODE.29.length.36535.cov.116.544.ID.57, NODE.30.length.36213.cov.125.651.ID.59, NODE.32.length.34242.cov.173.49.ID.63, NODE.36.length.32599.cov.42.2662.ID.71, NODE.38.length.32494.cov.1521.9.ID.75, NODE.40.length.31732.cov.65.9217.ID.79, NODE.41.length.30857.cov.158.625.ID.81, NODE.45.length.29440.cov.881.676.ID.89, NODE.47.length.28612.cov.168.357.ID.93, NODE.48.length.27881.cov.56.2302.ID.95, NODE.50.length.23915.cov.175.338.ID.99, NODE.54.length.21597.cov.80.5222.ID.107, NODE.55.length.21529.cov.24.5896.ID.109, NODE.63.length.19183.cov.1172.42.ID.125, NODE.64.length.18553.cov.38.2172.ID.127, NODE.65.length.18118.cov.21.9109.ID.129, NODE.66.length.17929.cov.824.738.ID.131, NODE.75.length.15736.cov.25.2833.ID.149, NODE.78.length.15058.cov.1670.05.ID.159, NODE.84.length.14270.cov.27.0028.ID.171, NODE.91.length.13283.cov.13.9537.ID.187, NODE.93.length.13109.cov.26.906.ID.191, NODE.96.length.12586.cov.1799.02.ID.197, NODE.99.length.12145.cov.19.2931.ID.203, NODE.100.length.11935.cov.22.0071.ID.205, NODE.104.length.11397.cov.58.2928.ID.211, NODE.114.length.9922.cov.127.37.ID.231, NODE.120.length.8752.cov.16.3855.ID.243, NODE.125.length.8316.cov.49.1684.ID.261, NODE.132.length.7818.cov.1132.6.ID.279, NODE.133.length.7773.cov.34.7356.ID.281, NODE.143.length.6960.cov.7.83664.ID.299, NODE.168.length.5606.cov.8.72275.ID.347, NODE.215.length.3917.cov.581.378.ID.453, NODE.216.length.3905.cov.4.70156.ID.455, NODE.220.length.3724.cov.4.82148.ID.463, NODE.223.length.3643.cov.2.97324.ID.469, NODE.235.length.3326.cov.162.369.ID.493, NODE.241.length.3253.cov.27.5804.ID.503, NODE.294.length.2485.cov.3.44733.ID.609, NODE.307.length.2338.cov.38.3355.ID.633 </p>
55	<p> NODE.31.length.35221.cov.298.507.ID.61, NODE.57.length.21063.cov.35.3396.ID.113, NODE.81.length.14499.cov.6.0432.ID.165, NODE.87.length.13835.cov.51.5428.ID.177, NODE.138.length.7406.cov.58.9336.ID.289, NODE.186.length.4825.cov.86.2273.ID.387, NODE.194.length.4567.cov.5.54078.ID.411, NODE.229.length.3496.cov.18.9326.ID.481, NODE.240.length.3255.cov.4.56438.ID.501, NODE.273.length.2762.cov.4.74067.ID.571, NODE.323.length.2068.cov.298.297.ID.661, NODE.335.length.2068.cov.10.1898.ID.685 </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	51.16 %
archaea	26.06 %

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 24,239,914 reads totaling 3,636.0 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 72 contigs in 72 scaffolds, totalling 1.8 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.