

1. Project Information

Program	Microbial/CSP 2012
PMO Project	0
Seq Proj ID	1027118
Sequencing Project Name	Alphaproteobacteria bacterium HL7711_P5A1 JGI 000151CP-C17
JGI Project ID	0

2. Read Statistics

Illumina Std PE Statistics

File name	7667.6.80858.CAACTA.fastq
Library	TGTP
Number of reads	24,207,654
Sequencing depth [†]	726X
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	24,207,654	100
Contam removed	10056	0.0
Artifact removed	357,102	1.5
Total removed	4,207,654	17.4
Total remaining	20,000,000	82.6

List of Contaminants Removed

Description	Num Reads	Pct Reads
human_chr11	9,056	0.04
human_chr6	488	0.00
gi 357579577 Canis_lupus_familiaris_chr3	346	0.00
human_chr2	304	0.00
gi 357579507 Canis_lupus_familiaris_chr38	130	0.00
gi 357579535 Canis_lupus_familiaris_chr20	18	0.00
gi 357579571 Canis_lupus_familiaris_chr5	14	0.00

human_chr8	6	0.00
human_chr14	4	0.00
human_chr18	2	0.00
human_chr7	2	0.00
human_chr9	2	0.00
gi 357579523 Canis_lupus_familiaris_chr27	2	0.00

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

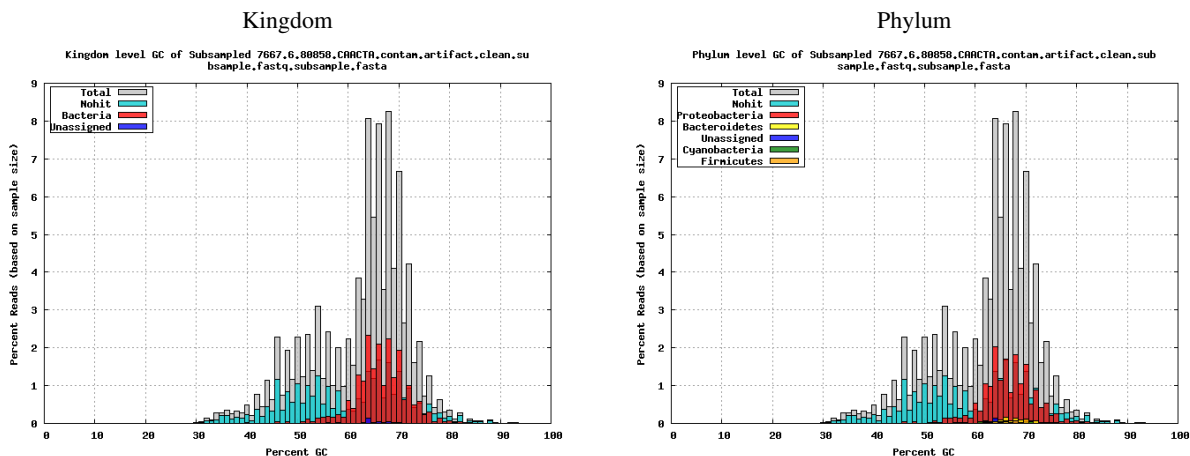
Illumina Std PE Contamination Identification Statistics

Description	Num Reads	Pct Reads
Input	24,207,654	100
Contam identified	2,212	0.0

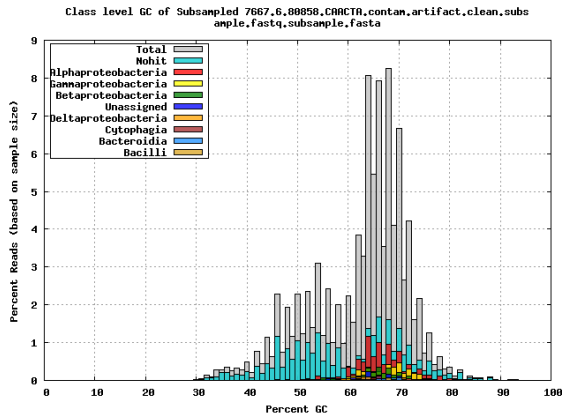
List of Contaminants Identified

Description	Num Reads	Pct Reads
<i>Delftia</i>	2,202	0.01
<i>Escherichia</i>	4	0.00
<i>Pseudomonas</i>	2	0.00
<i>Shigella</i>	2	0.00
<i>Ralstonia</i>	2	0.00

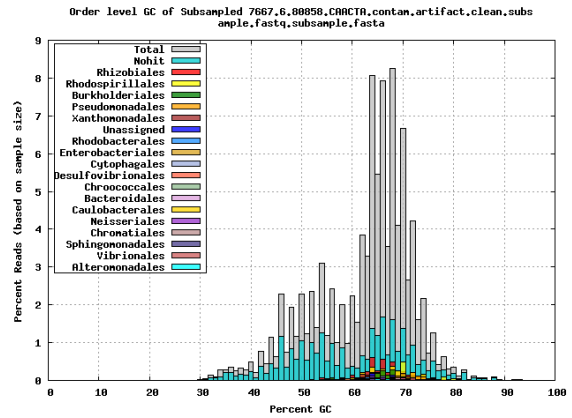
GC histogram of the reads subsampled to 10k, overlaid with GC of hits based on BLASTX, shown for different taxonomic levels.



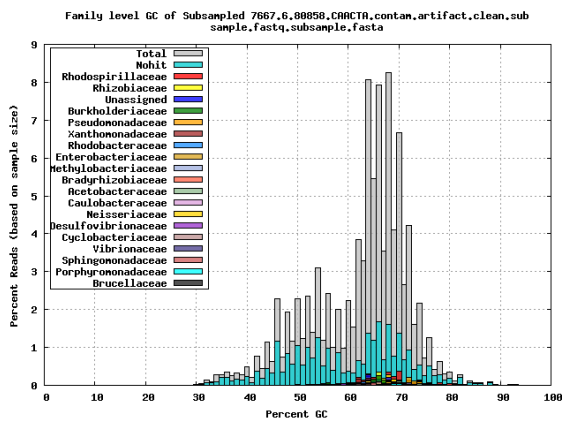
Class



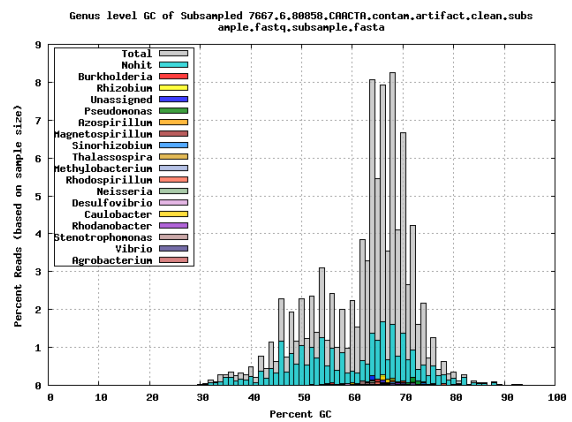
Order



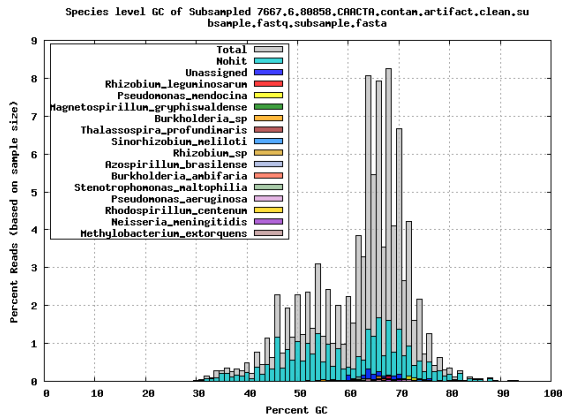
Family



Genus



Species

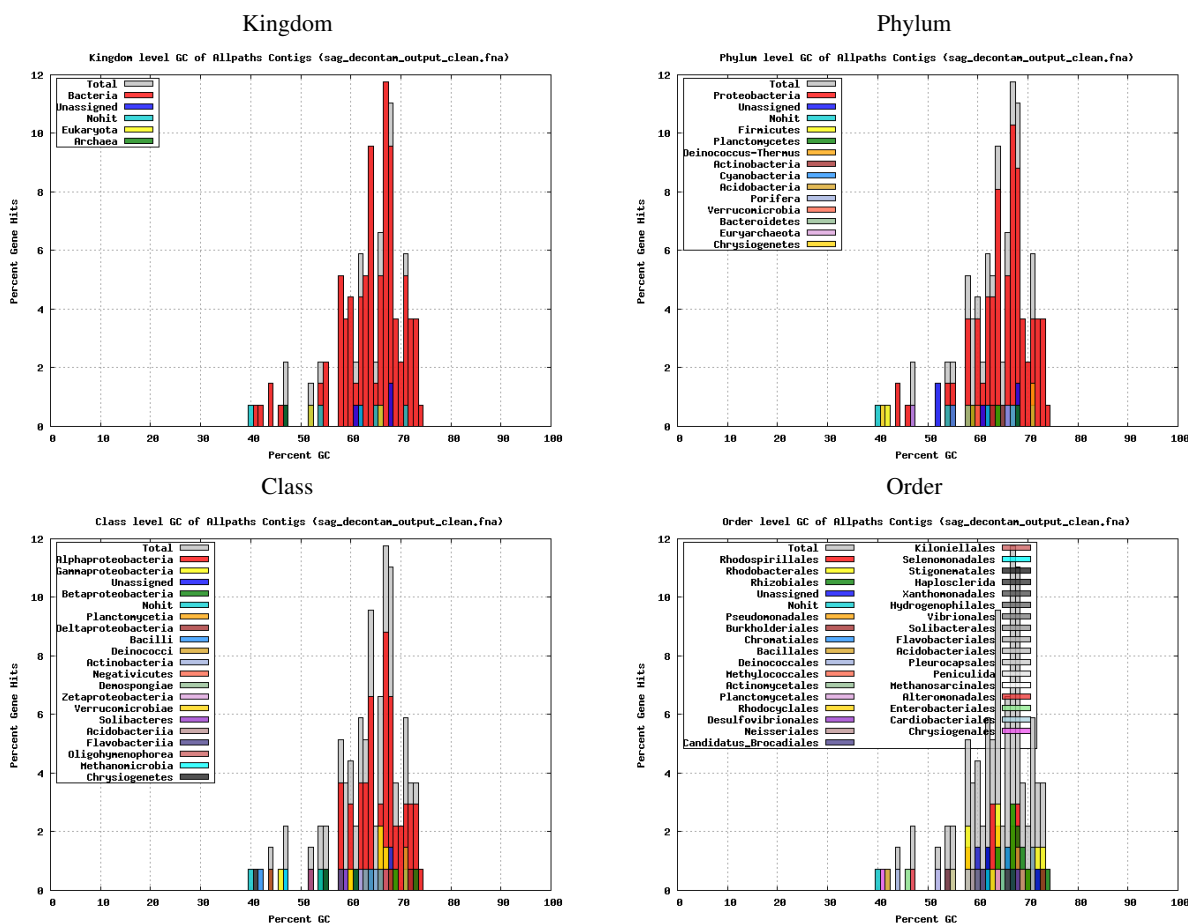


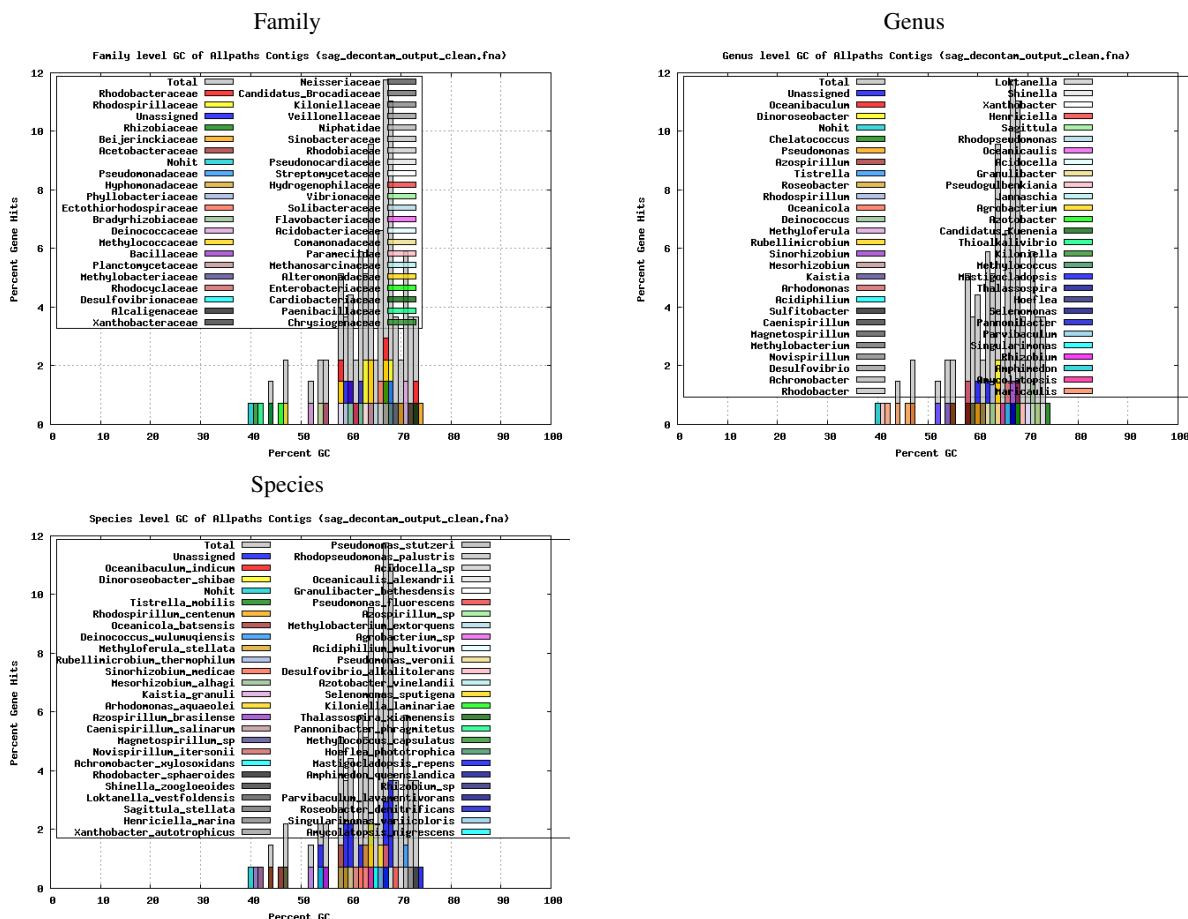
4. Assembly Statistics

Assembly method	SPAdes with auto decontamination
Scaffold total	17
Contig total	17
Scaffold sequence length	134.6 kb
Contig sequence length	134.6 kb (0.0% gap)
Scaffold N/L50	6/8.3 kb
Contig N/L50	6/8.3 kb
Largest Contig	19.8 kb
Number of scaffolds >50 kb	0
Pct of genome in scaffolds >50 kb	0.0
Pct of reads assembled (raw)	24.3
Pct of reads assembled (decontam)	2.4

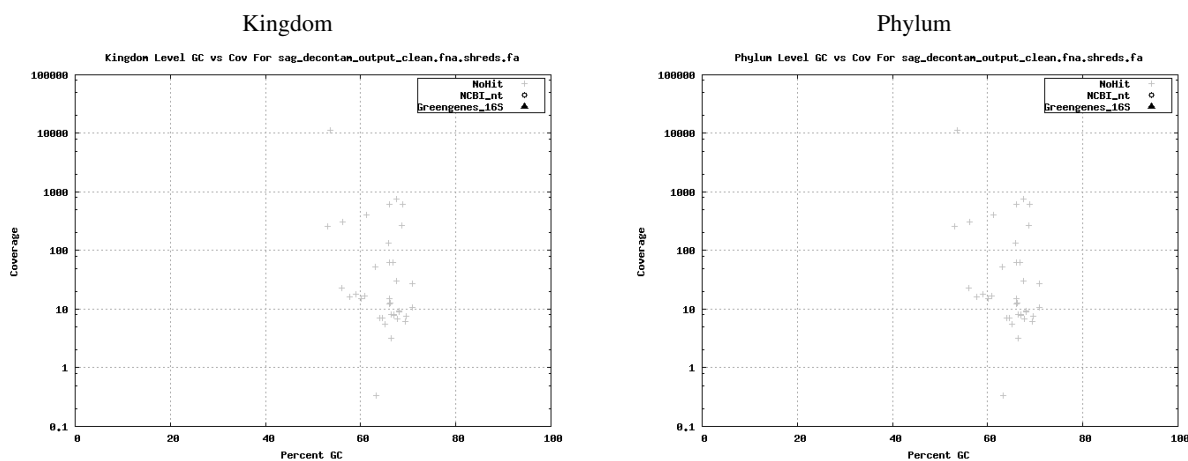
5. Assembly QC Results

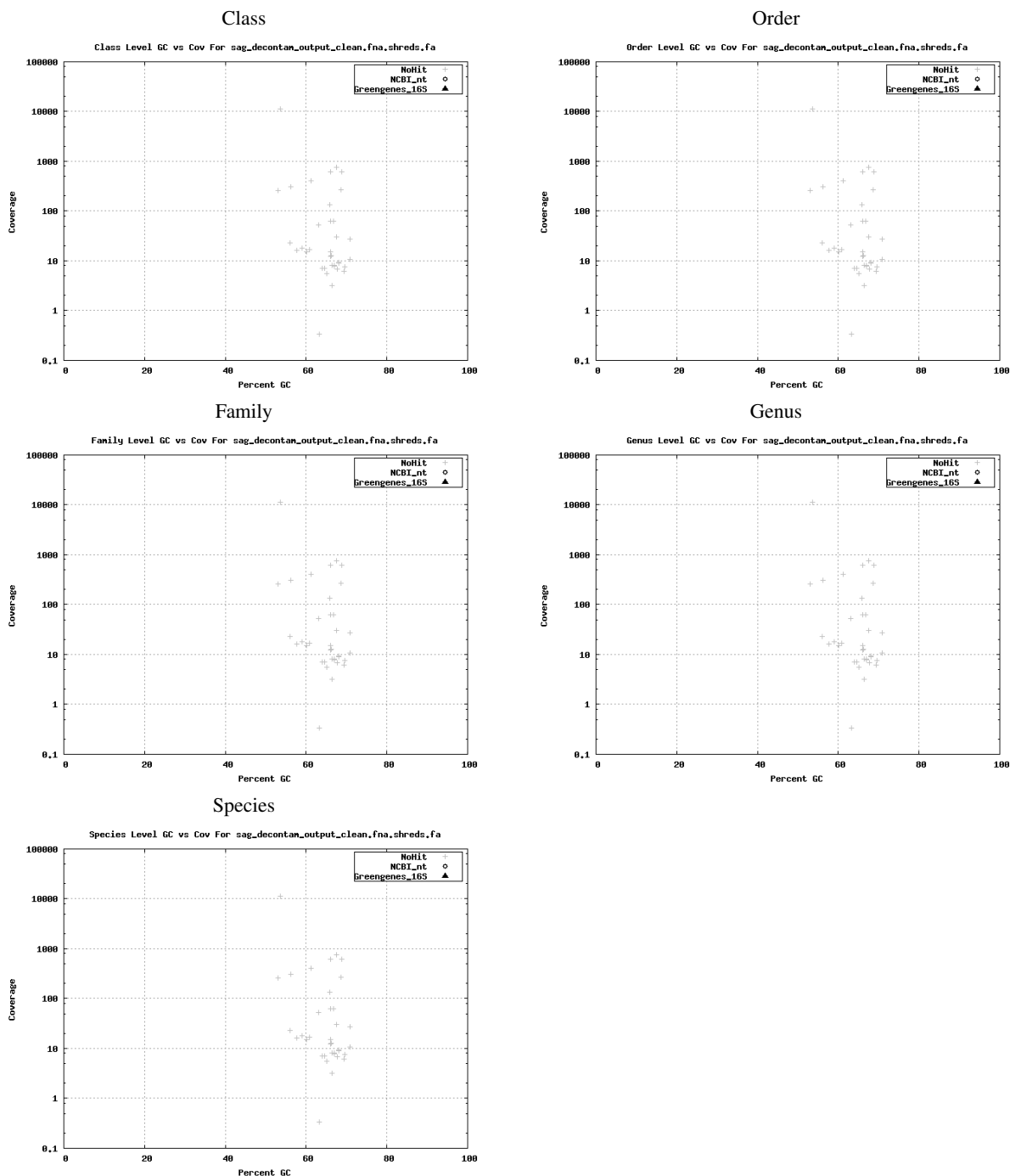
GC histogram of the predicted genes on each contig, overlaid with GC of hits based on BLASTP, shown for different taxonomic levels.



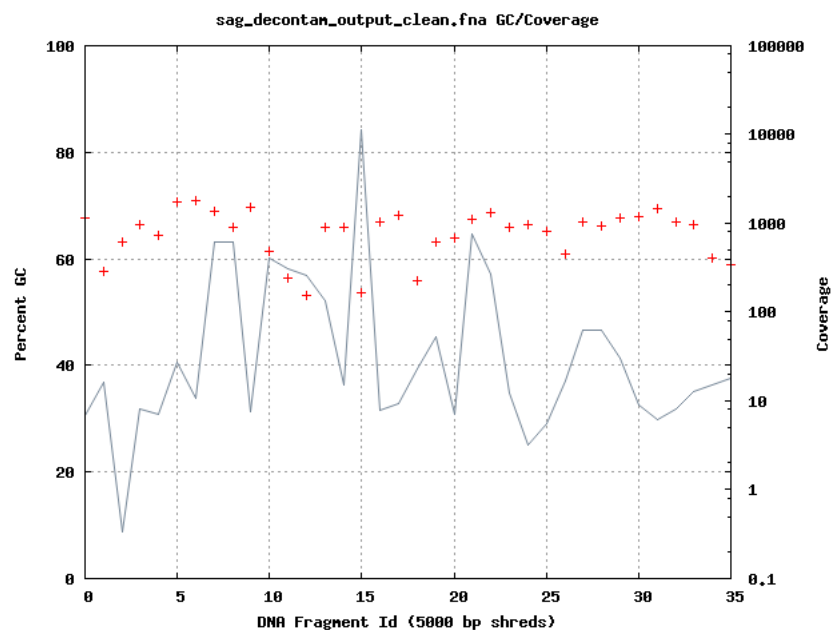


GC vs coverage based on GC of NCBI nt and Greengenes 16S rRNA gene hits to the assembly using megablast, shown for different taxonomic levels.

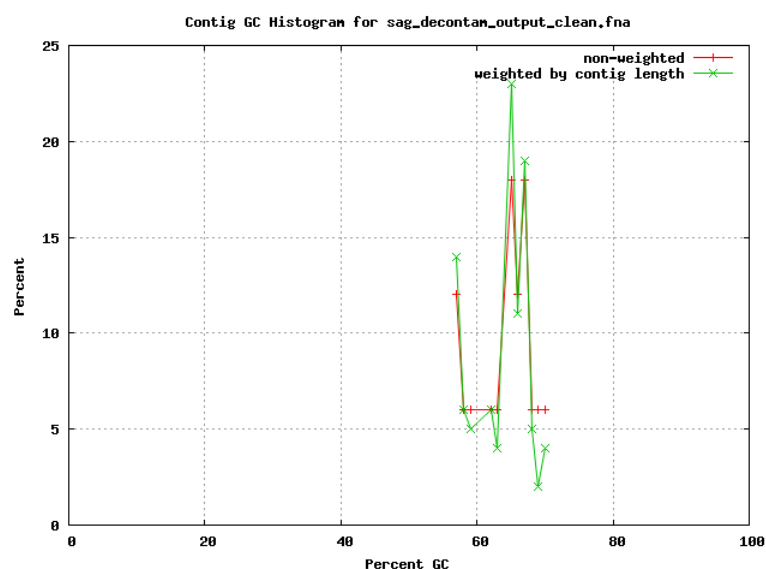




Coverage vs GC. Contigs were shredded into non-overlapping 5kbp and the GC of each shred was plotted as a point, colored by scaffold id. Coverage was calculated by mapping the fragment library to the final assembly and plotted as connected points.



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

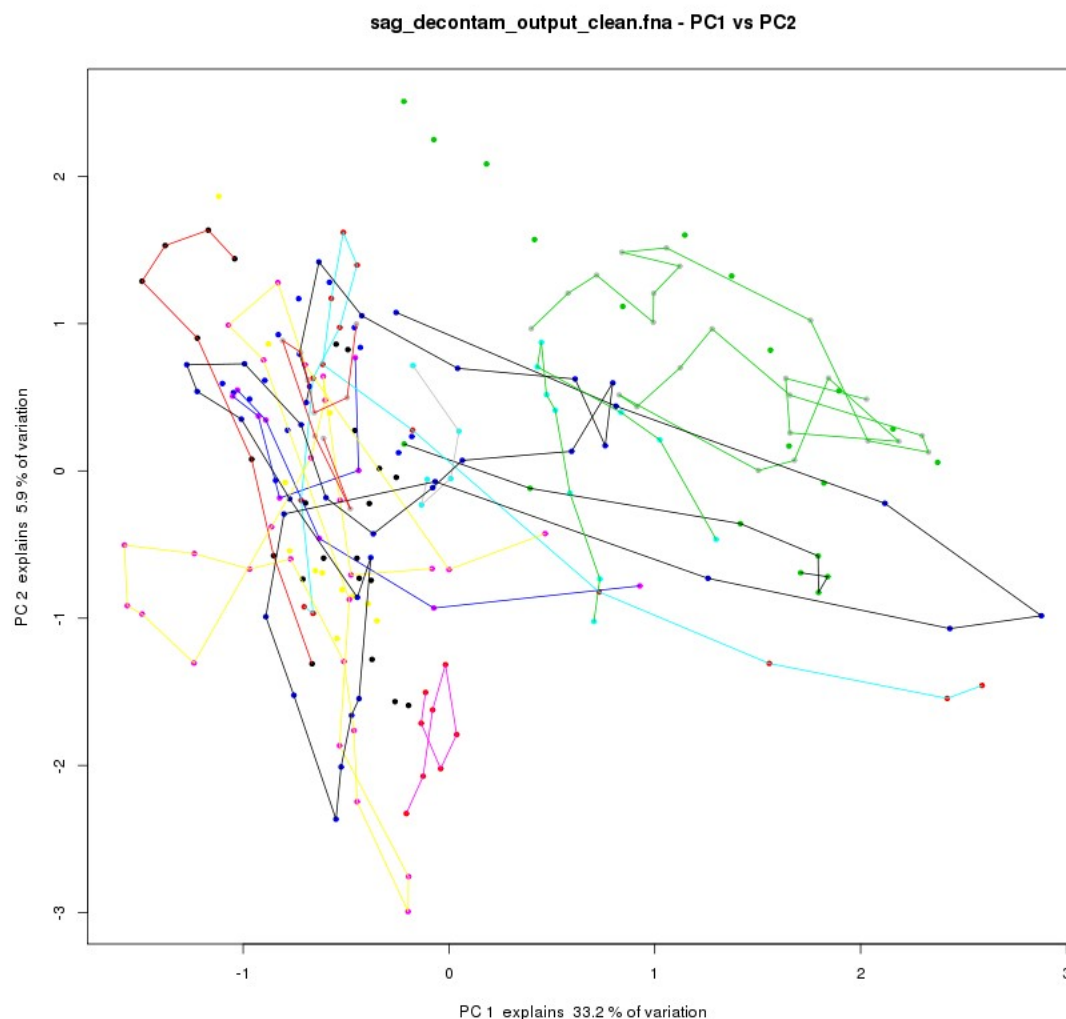


List of contigs and average percent GC, grouped in bins of 5:

Pct GC Bin	Contig Name
55	NODE.6.length.13977.cov.220.579.ID.11, NODE.17.length.8256.cov.5227.94.ID.37, NODE.25.length.6735.cov.10.8698.ID.53, NODE.39.length.5076.cov.10.2633.ID.81
60	NODE.15.length.8710.cov.9.93507.ID.33, NODE.35.length.5623.cov.31.2552.ID.73
65	NODE.3.length.19775.cov.27.968.ID.5, NODE.4.length.15333.cov.218.005.ID.7, NODE.14.length.8740.cov.6.44134.ID.31, NODE.18.length.8134.cov.593.378.ID.39, NODE.26.length.6676.cov.70.5268.ID.55, NODE.30.length.6258.cov.5.18926.ID.63,

	NODE_32_length.6125_cov.5.43591_ID.67, NODE_48_length.4080_cov.3.50534_ID.103, NODE_57_length.3125_cov.4.64495_ID.121, NODE_71_length.2370_cov.4.31965_ID.151
70	NODE_34_length.5635_cov.17.0763_ID.71

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	4.8 %
archaea	3.43 %

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 24,207,654 reads totaling 3,631.1 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 17 contigs in 17 scaffolds, totalling 134.6 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.