

1. Project Information

Program	Microbial/CSP 2012
PMO Project	0
Seq Proj ID	1027166
Sequencing Project Name	Roseibacterium sp. HL7711_P4H5 JGI 000149CP-E10
JGI Project ID	0

2. Read Statistics

Illumina Std PE Statistics

File name	7667.7.80862.TCATTC.fastq
Library	TNGY
Number of reads	32,440,782
Sequencing depth [†]	973X
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	32,440,782	100
Contam removed	122	0.0
Artifact removed	453,038	1.4
Total removed	12,440,782	38.3
Total remaining	20,000,000	61.7

List of Contaminants Removed

Description	Num Reads	Pct Reads
gi 357579577 Canis_lupus_familiaris_chr3	80	0.00
human_chr2	70	0.00
gi 357579535 Canis_lupus_familiaris_chr20	32	0.00
gi 357579571 Canis_lupus_familiaris_chr5	8	0.00
human_chr11	4	0.00
human_chr17	2	0.00
human_chr14	2	0.00

human_chr4	2	0.00
------------	---	------

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

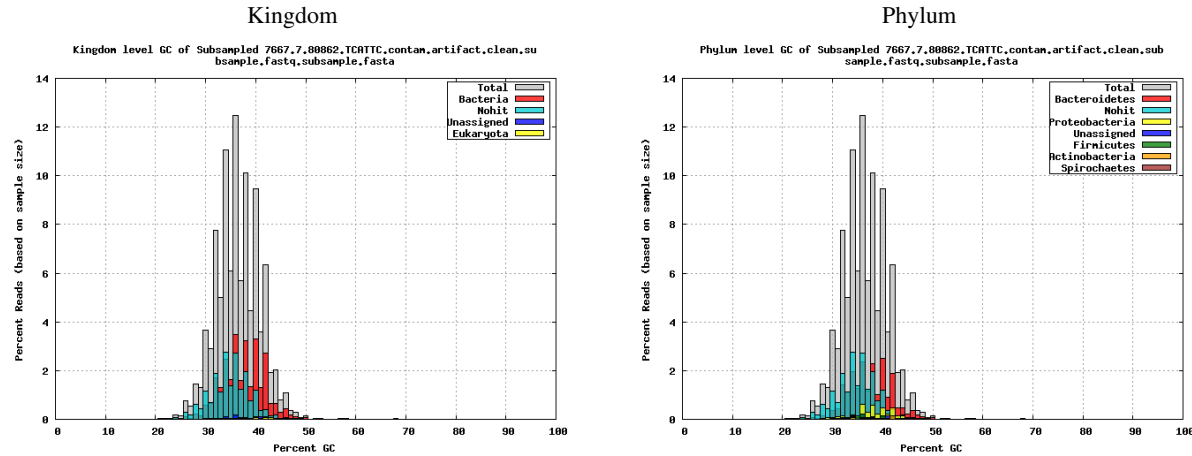
Illumina Std PE Contamination Identification Statistics

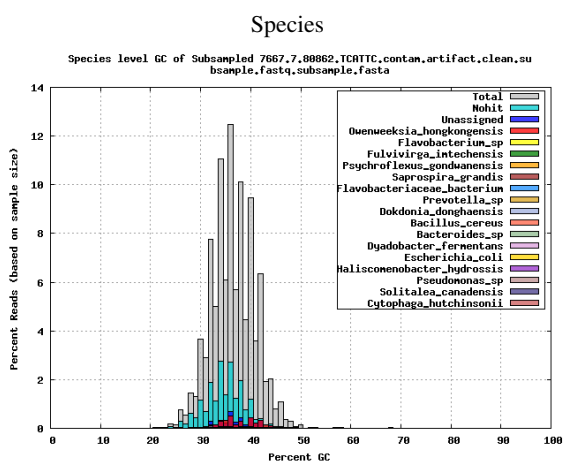
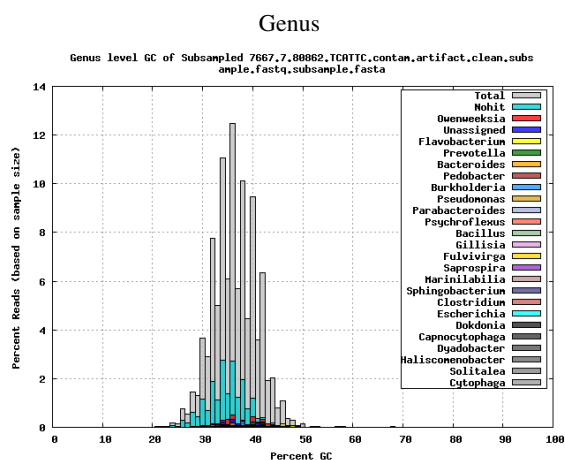
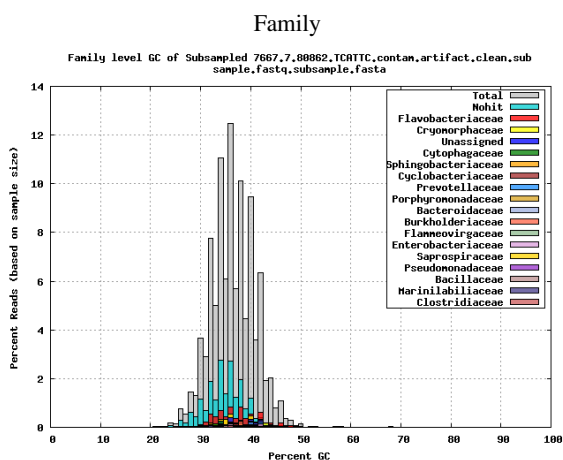
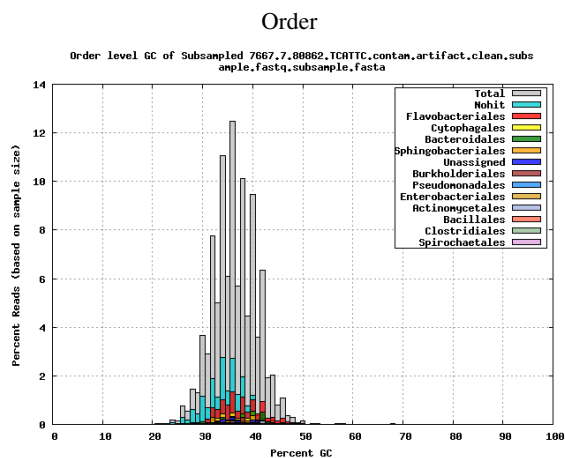
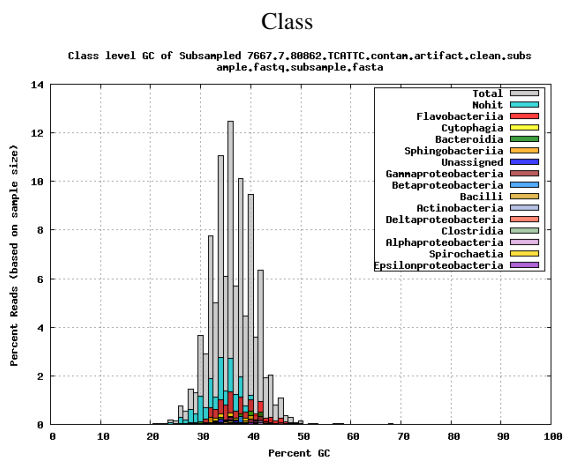
Description	Num Reads	Pct Reads
Input	32,440,782	100
Contam identified	4	0.0

List of Contaminants Identified

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



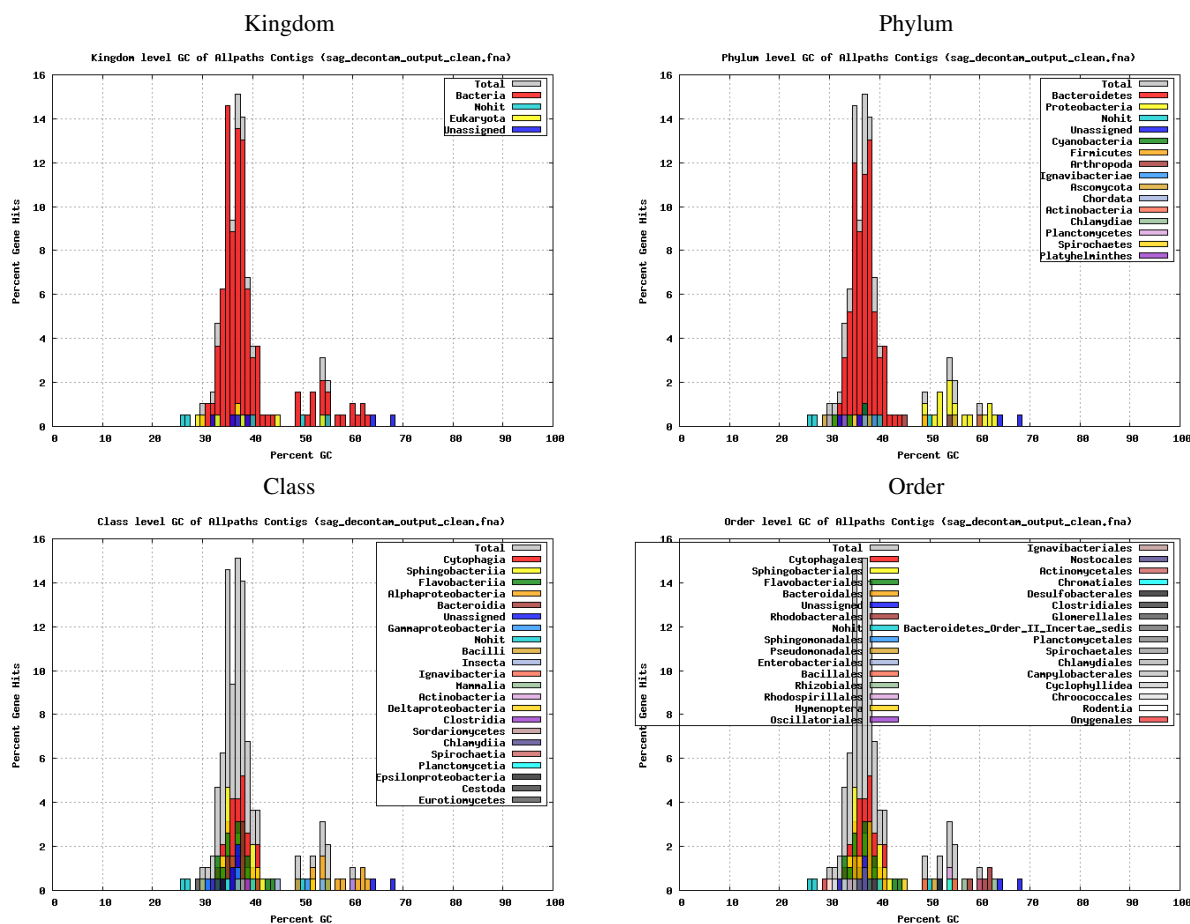


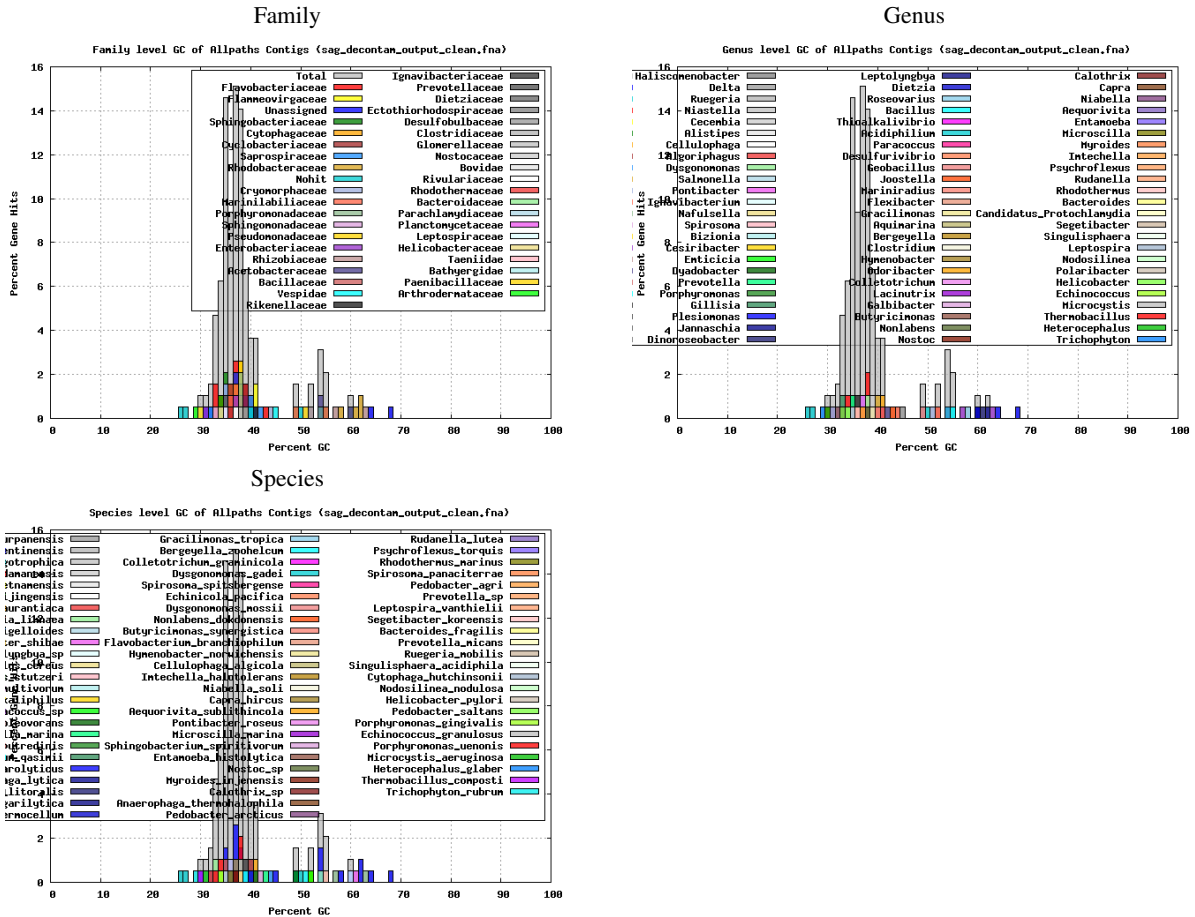
4. Assembly Statistics

Assembly method	SPAdes with auto decontamination
Scaffold total	29
Contig total	29
Scaffold sequence length	210.1 kb
Contig sequence length	210.1 kb (0.0% gap)
Scaffold N/L50	12/7.2 kb
Contig N/L50	12/7.2 kb
Largest Contig	14.6 kb
Number of scaffolds >50 kb	0
Pct of genome in scaffolds >50 kb	0.0
Pct of reads assembled (raw)	64.4
Pct of reads assembled (decontam)	11.8

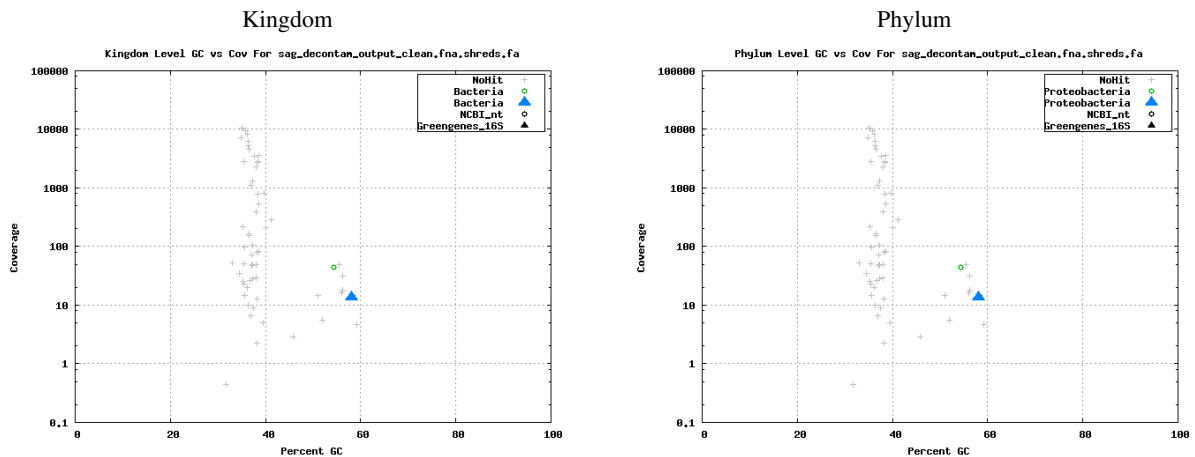
5. Assembly QC Results

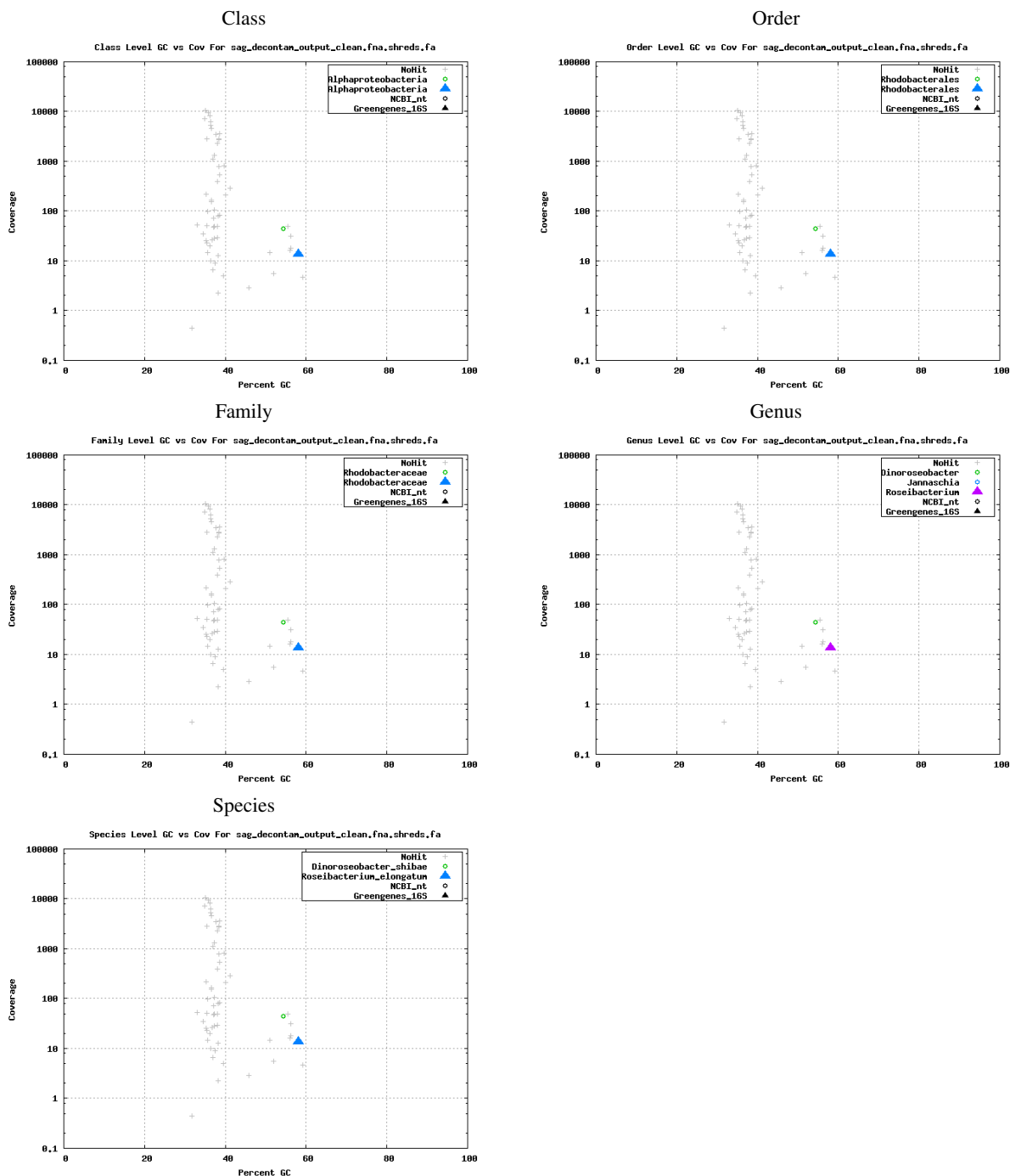
GC histogram of the predicted genes on each contig, overlaid with GC of hits based on BLASTP, shown for different taxonomic levels.



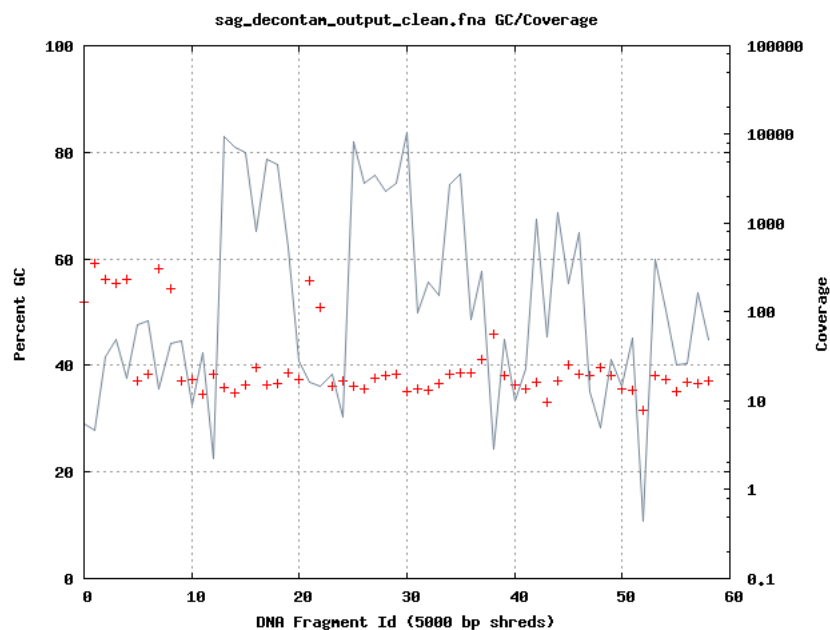


GC vs coverage based on GC of NCBI nt and Greengenes 16S rRNA gene hits to the assembly using megablast, shown for different taxonomic levels.

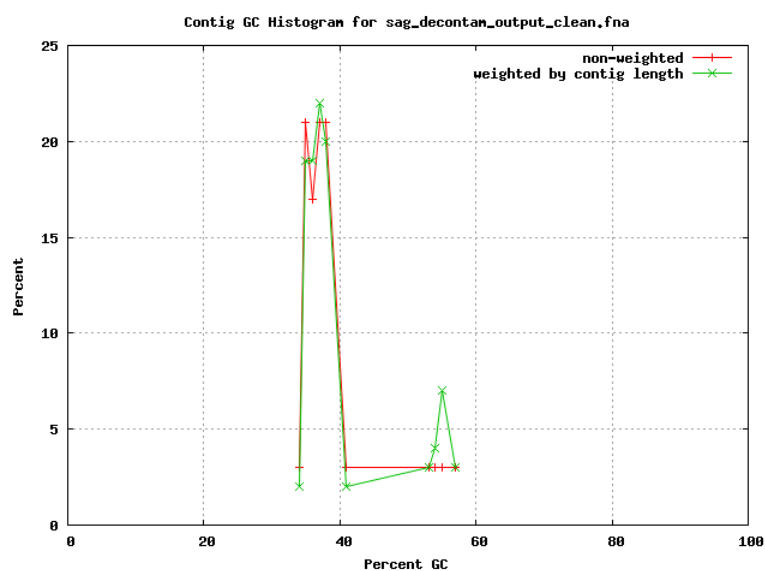




Coverage vs GC. Contigs were shredded into non-overlapping 5kbp and the GC of each shred was plotted as a point, colored by scaffold id. Coverage was calculated by mapping the fragment library to the final assembly and plotted as connected points.



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.25.length.5162.cov.22.9007.ID.53
35	NODE.2.length.11245.cov.1916.41.ID.5, NODE.3.length.10733.cov.4191.58.ID.7, NODE.4.length.9927.cov.3340.24.ID.9, NODE.5.length.9044.cov.76.2532.ID.11, NODE.6.length.8970.cov.2550.5.ID.13, NODE.7.length.8758.cov.443.871.ID.15, NODE.8.length.8145.cov.592.497.ID.17, NODE.10.length.7498.cov.22.9113.ID.21, NODE.11.length.7468.cov.5975.53.ID.23, NODE.12.length.7228.cov.18.6123.ID.25, NODE.13.length.6902.cov.2673.99.ID.27, NODE.14.length.6792.cov.17.6531.ID.33, NODE.15.length.6622.cov.91.2424.ID.35, NODE.16.length.6588.cov.139.91.ID.31,

	NODE_17.length_6482_cov_7.19387.ID_37, NODE_18.length_6284_cov_4815.56.ID_39, NODE_21.length_5900_cov_242.158.ID_45, NODE_22.length_5506_cov_12.8806.ID_47, NODE_23.length_5343_cov_31.6906.ID_49, NODE_24.length_5295_cov_332.598.ID_51, NODE_26.length_5120_cov_33.3931.ID_61, NODE_28.length_4929_cov_15.1143.ID_65, NODE_30.length_4824_cov_529.671.ID_69
40	NODE_27.length_5061_cov_189.413.ID_63
50	NODE_9.length_7622_cov_11.141.ID_19, NODE_19.length_6131_cov_3.65701.ID_41
55	NODE_1.length_14556_cov_22.7757.ID_3, NODE_20.length_5989_cov_12.8426.ID_43

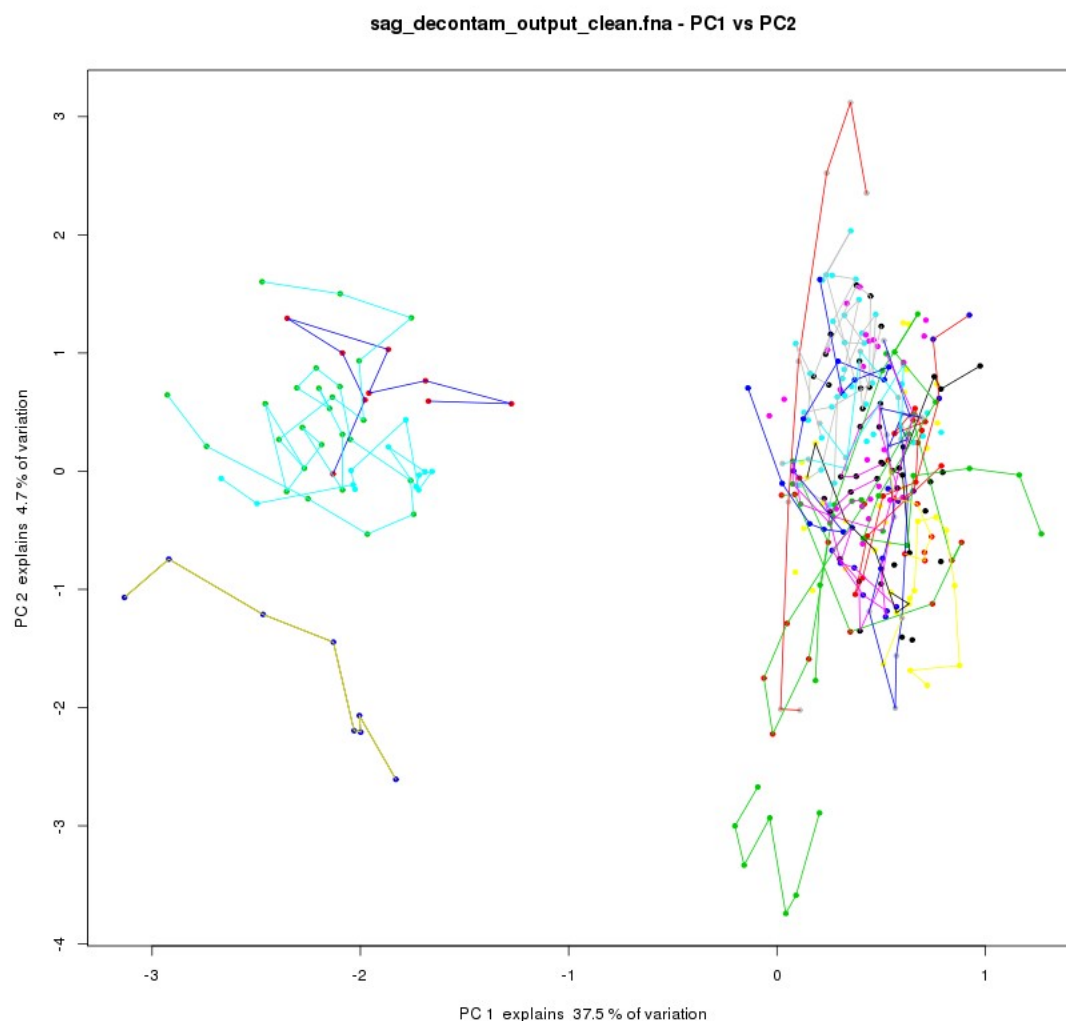
List of the top contig megablast hits against potential reagent and process contaminants.

Organism	Align Length (bp)	Pct Id	Contig Name
<i>Ralstonia solanacearum_str.PSI07_chromosome_comple</i> <i>te_genome</i>	198	91.92	NODE_20.length_5989_cov_12.8426.ID_43

List of the top contig megablast hits against 16S ribosomal RNA genes.

Organism	Align Length (bp)	Pct Id	Contig Name
<i>520995_Roseibacterium_elongatum_str.DSM.19469_FN66</i> <i>7962_1_1.1456</i>	1,459	98.15	NODE_20.length_5989_cov_12.8426.ID_43

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 %
archaea	1.37 %

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 32,440,782 reads totaling 4,866.1 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 29 contigs in 29 scaffolds, totalling 210.1 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.