

1. Project Information

Program	Microbial/CSP 2012
PMO Project	0
Seq Proj ID	1027133
Sequencing Project Name	Thioalkalivibrio sp. HL7711.P3F6 JGI 000155CP-M11
JGI Project ID	0

2. Read Statistics

Illumina Std PE Statistics

File name	7667.6.80858.CATGGC.fastq
Library	TNGC
Number of reads	25,656,480
Sequencing depth [†]	770X
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	25,656,480	100
Contam removed	1128	0.0
Artifact removed	368,132	1.4
Total removed	5,656,480	22.0
Total remaining	20,000,000	78.0

List of Contaminants Removed

Description	Num Reads	Pct Reads
human_chr11	490	0.00
human_chr6	338	0.00
gi 357579577 Canis_lupus_familiaris_chr3	262	0.00
human_chr2	256	0.00
gi 357579535 Canis_lupus_familiaris_chr20	22	0.00
gi 357579571 Canis_lupus_familiaris_chr5	18	0.00
human_chr1	2	0.00

human_chr19	2	0.00
human_chr3	2	0.00
human_chr8	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.

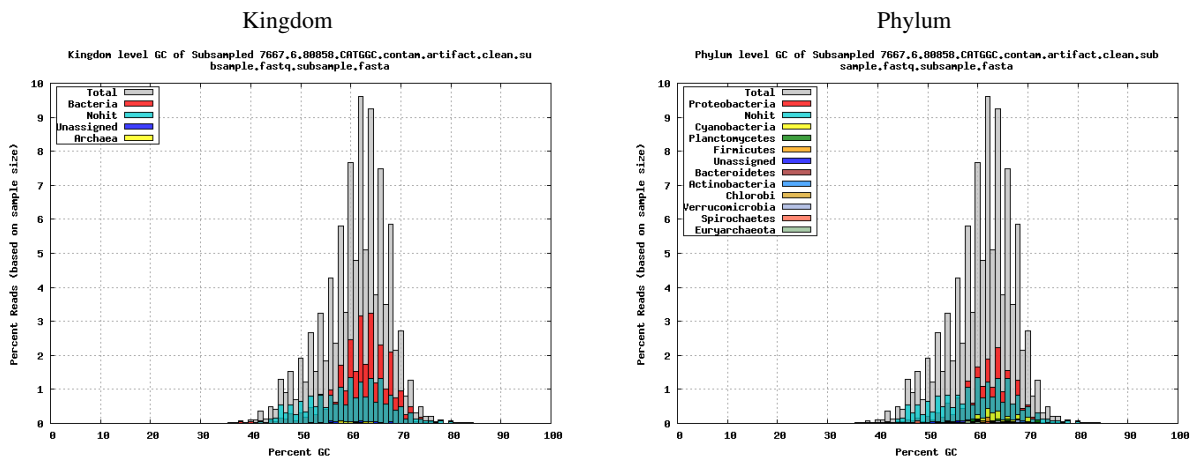
Illumina Std PE Contamination Identification Statistics

Description	Num Reads	Pct Reads
Input	25,656,480	100
Contam identified	16	0.0

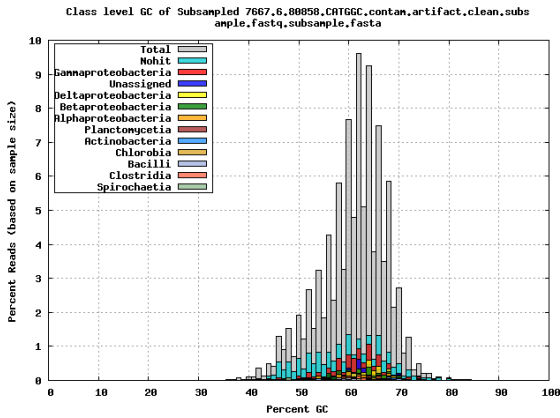
List of Contaminants Identified

Description	Num Reads	Pct Reads
<i>Ralstonia</i>	14	0.00
<i>Delftia</i>	2	0.00

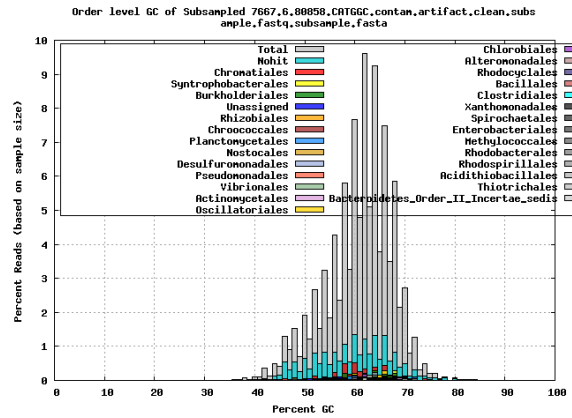
GC histogram of the reads subsampled to 10k, overlaid with GC of hits based on BLASTX, shown for different taxonomic levels.



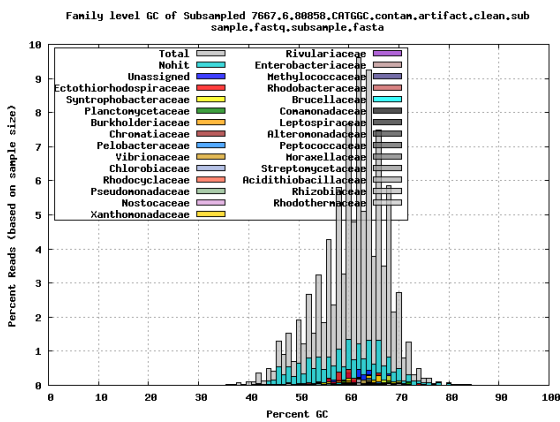
Class



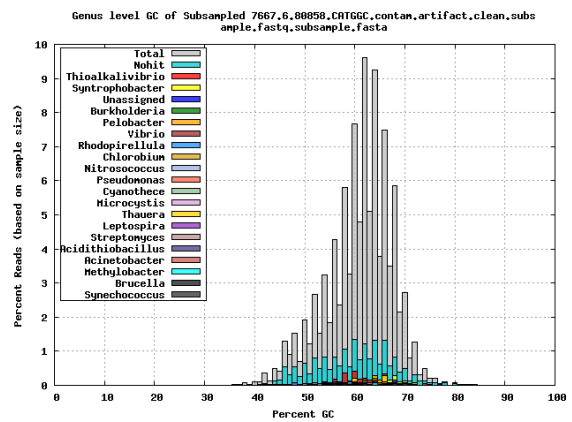
Order



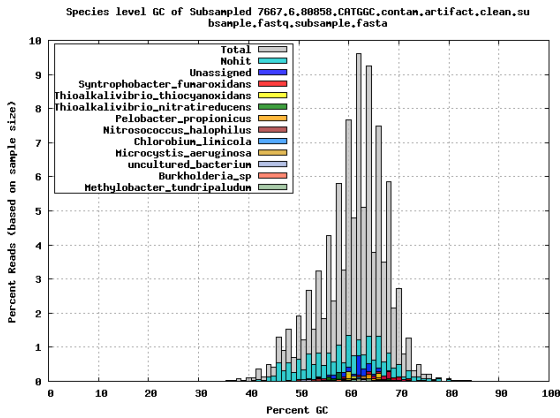
Family



Genus



Species

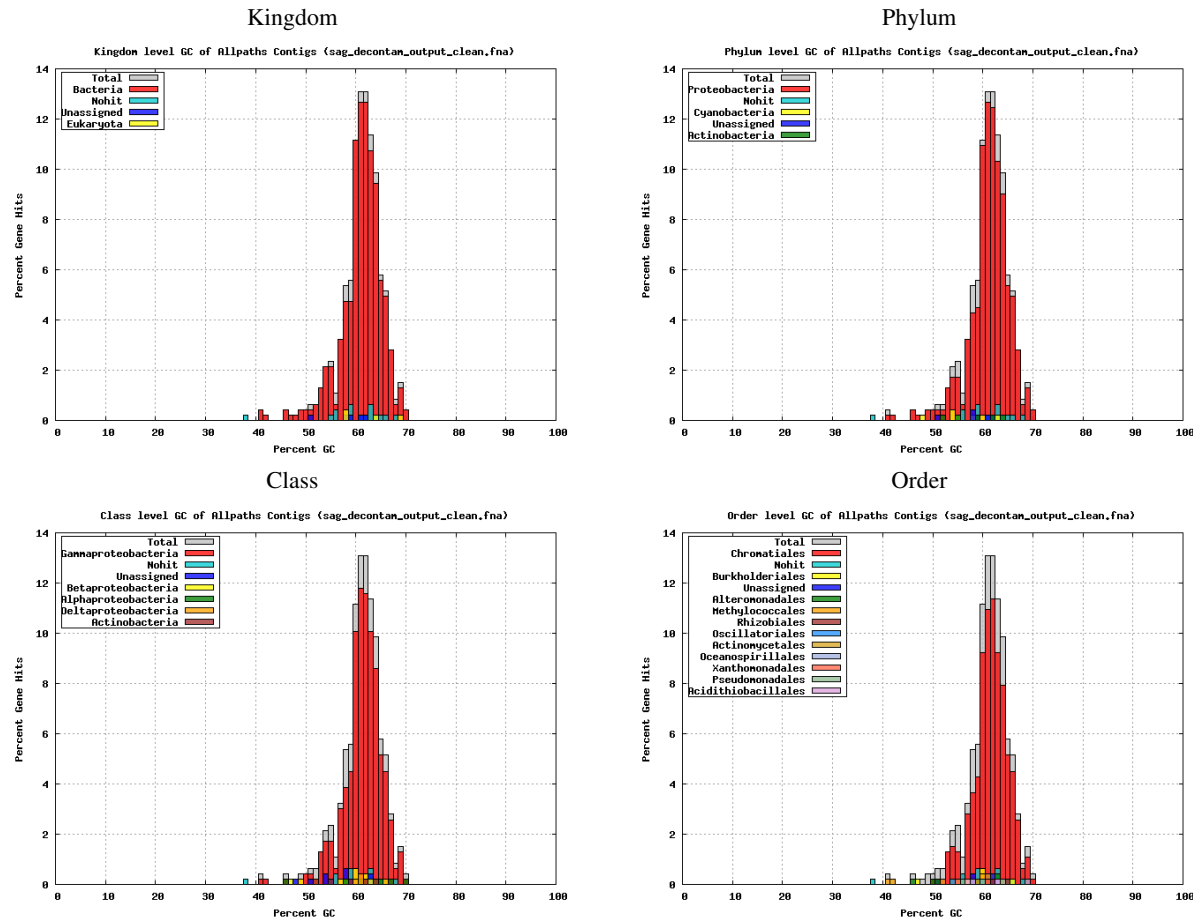


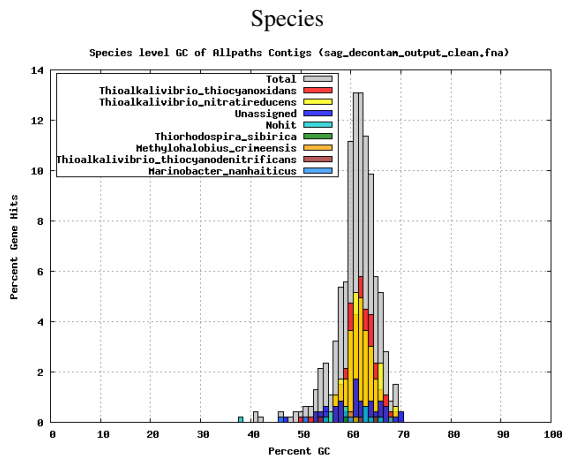
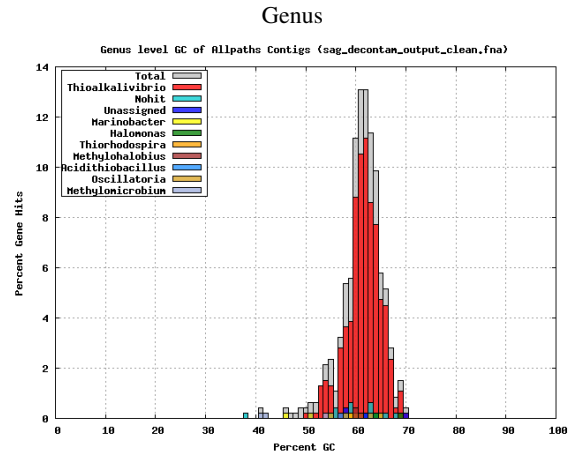
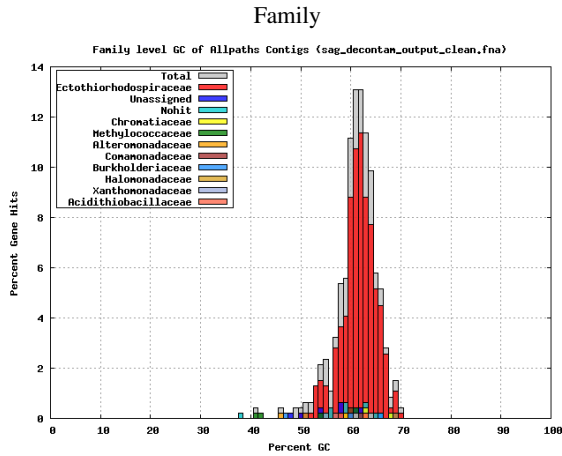
4. Assembly Statistics

Assembly method	SPAdes with auto decontamination
Scaffold total	48
Contig total	48
Scaffold sequence length	440.9 kb
Contig sequence length	440.9 kb (0.0% gap)
Scaffold N/L50	11/12.5 kb
Contig N/L50	11/12.5 kb
Largest Contig	35.1 kb
Number of scaffolds >50 kb	0
Pct of genome in scaffolds >50 kb	0.0
Pct of reads assembled (raw)	52.2
Pct of reads assembled (decontam)	18.0

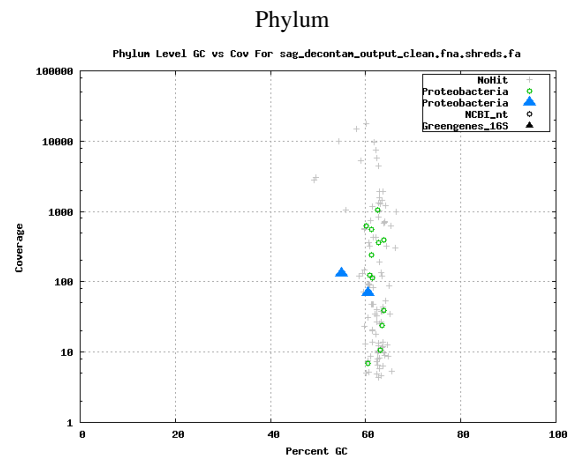
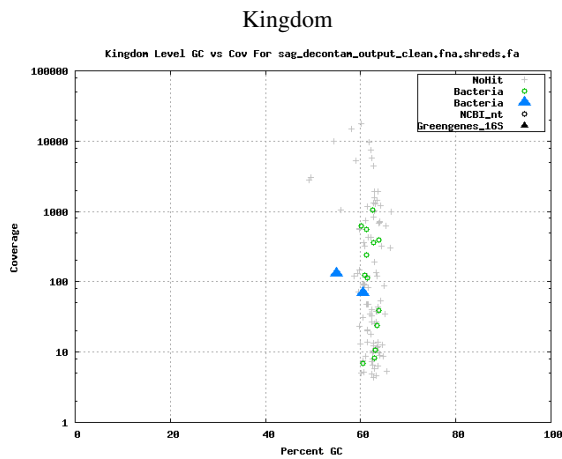
5. Assembly QC Results

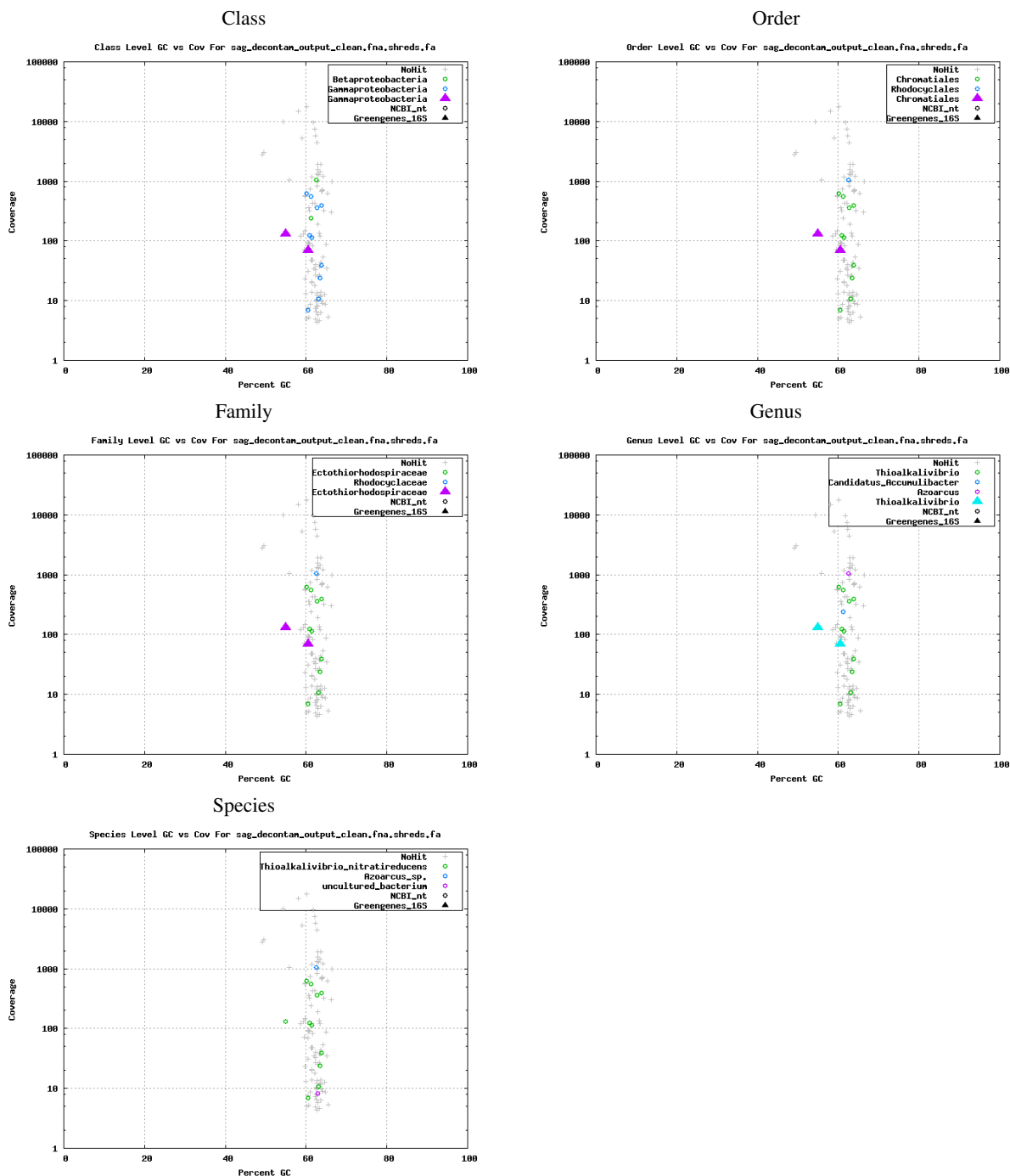
GC histogram of the predicted genes on each contig, overlaid with GC of hits based on BLASTP, shown for different taxonomic levels.



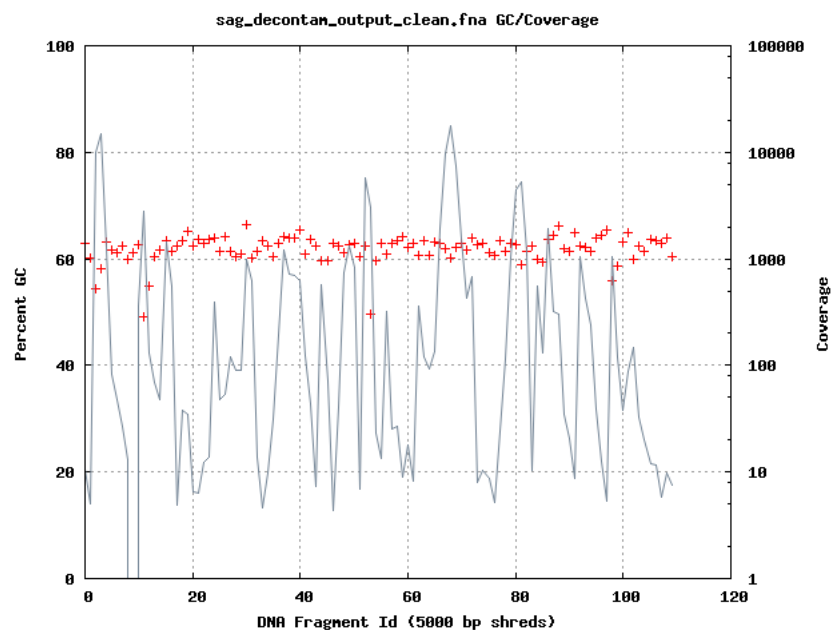


GC vs coverage based on GC of NCBI nt and Greengenes 16S rRNA gene hits to the assembly using megablast, shown for different taxonomic levels.

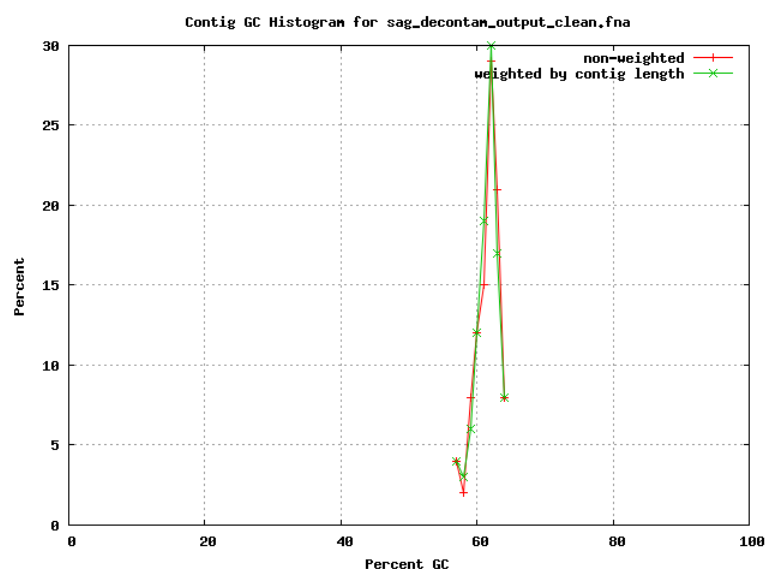




Coverage vs GC. Contigs were shredded into non-overlapping 5kbp and the GC of each shred was plotted as a point, colored by scaffold id. Coverage was calculated by mapping the fragment library to the final assembly and plotted as connected points.



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.9.length.14657.cov.55.7333.ID.17, NODE.18.length.10850.cov.357.475.ID.35, NODE.24.length.9481.cov.217.496.ID.51, NODE.30.length.8142.cov.893.842.ID.63, NODE.35.length.7613.cov.268.341.ID.73, NODE.43.length.6494.cov.3444.ID.93 NODE.51.length.4278.cov.97.276.ID.107
60	NODE.1.length.35134.cov.2500.51.ID.1, NODE.2.length.30992.cov.4124.12.ID.3, NODE.3.length.27534.cov.387.07.ID.5, NODE.4.length.26544.cov.1446.42.ID.7, NODE.6.length.19205.cov.489.073.ID.11, NODE.7.length.19200.cov.12.9291.ID.13,

NODE.8.length.19031.cov.58.4558.ID.15, NODE.11.length.13831.cov.403.988.ID.21, NODE.13.length.13083.cov.5.0895.ID.25, NODE.15.length.12482.cov.622.018.ID.29, NODE.17.length.11621.cov.151.321.ID.33, NODE.19.length.10859.cov.22.3055.ID.37, NODE.20.length.10810.cov.17.1638.ID.39, NODE.21.length.10721.cov.15.8692.ID.41, NODE.25.length.8899.cov.8.47286.ID.53, NODE.26.length.8548.cov.552.065.ID.55, NODE.27.length.8540.cov.162.486.ID.57, NODE.28.length.8203.cov.17.6103.ID.59, NODE.34.length.7658.cov.12.9444.ID.71, NODE.36.length.7397.cov.8.76246.ID.75, NODE.42.length.6692.cov.5.12566.ID.87, NODE.44.length.6456.cov.7.67318.ID.47, NODE.47.length.4965.cov.5.17943.ID.99, NODE.49.length.4747.cov.28.9386.ID.103, NODE.52.length.4294.cov.5.14485.ID.109, NODE.58.length.3768.cov.8.32534.ID.119, NODE.74.length.3052.cov.456.668.ID.149, NODE.76.length.3008.cov.428.988.ID.153, NODE.78.length.2942.cov.59.2082.ID.157, NODE.79.length.2937.cov.3.22623.ID.159, NODE.81.length.2769.cov.3.85741.ID.163, NODE.83.length.2696.cov.6.51571.ID.167, NODE.85.length.2560.cov.3.12335.ID.171, NODE.86.length.2517.cov.218.8.ID.173, NODE.87.length.2501.cov.4.51431.ID.175, NODE.91.length.2430.cov.4.52505.ID.181, NODE.100.length.2222.cov.7.22981.ID.203, NODE.102.length.2199.cov.87.1674.ID.207, NODE.103.length.2198.cov.963.294.ID.209, NODE.108.length.2075.cov.4.12871.ID.219, NODE.113.length.2006.cov.2.95079.ID.229

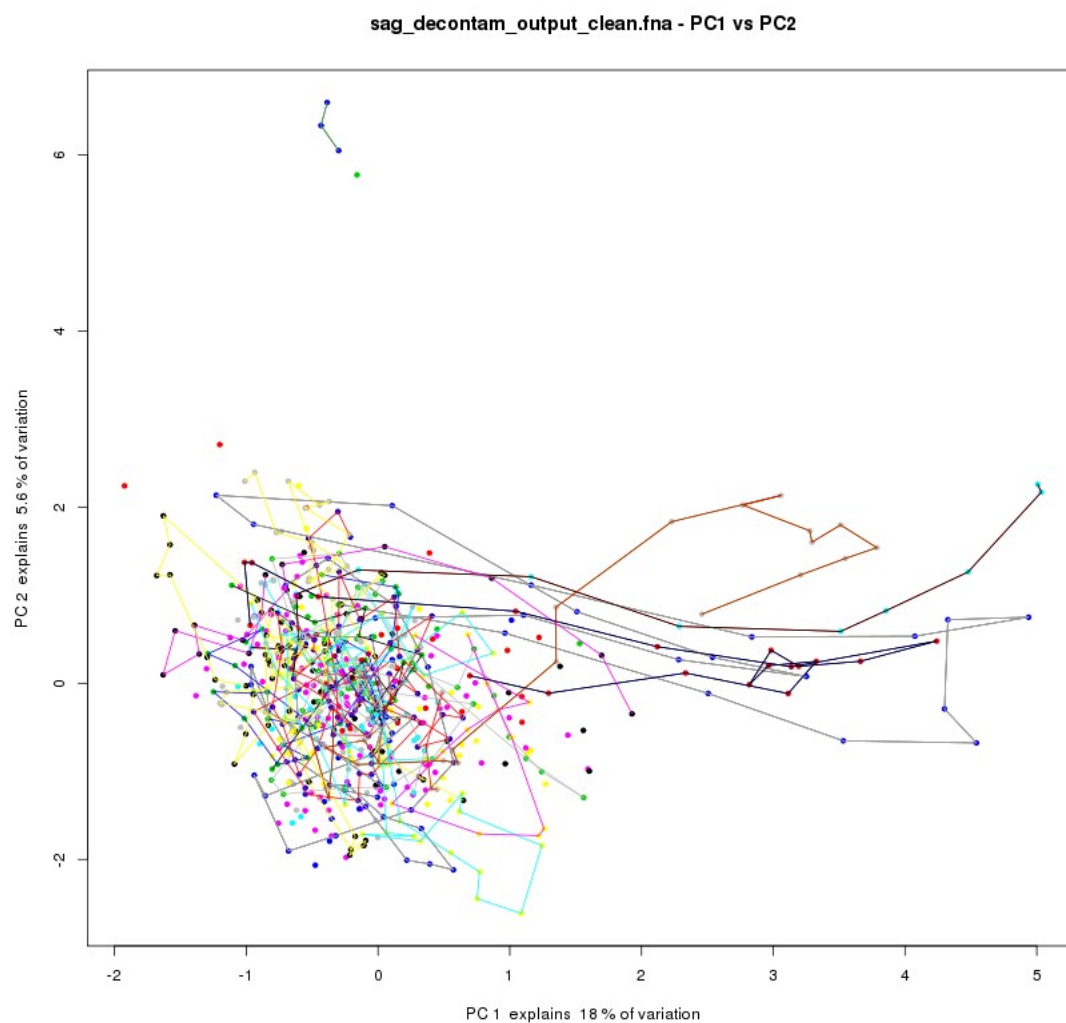
List of the top contig megablast hits against potential reagent and process contaminants.

Organism	Align Length (bp)	Pct Id	Contig Name
<i>Ralstonia solanacearum_str.CFBP2957_chromosome.com</i> <i>plete_genome</i>	84	97.62	NODE.11.length.13831.cov.403.988.ID.21

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

Organism	Align Length (bp)	Pct Id	Contig Name
<i>517623_Thioalkalivibrio_sp_str.K90mix.CP001905.1.4</i> <i>23224_424763</i>	1,542	95.78	NODE.9.length.14657.cov.55.7333.ID.17

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	30.38 %
archaea	13.72 %

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 25,656,480 reads totaling 3,848.5 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 48 contigs in 48 scaffolds, totalling 440.9 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.