

1. Project Information

Program	Microbial/CSP 2012
PMO Project	0
Seq Proj ID	1027115
Sequencing Project Name	Alphaproteobacteria bacterium HL7711_P5A1 JGI 000149CP-E07
JGI Project ID	0

2. Read Statistics

Illumina Std PE Statistics

File name	7667.6.80858.CAAAAG.fastq
Library	TGTO
Number of reads	28,183,458
Sequencing depth [†]	846X
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	28,183,458	100
Contam removed	13845384	49.1
Artifact removed	1,503,274	5.3
Total removed	15,348,658	54.5
Total remaining	12,834,800	45.5

List of Contaminants Removed

Description	Num Reads	Pct Reads
human_chr6	13,841,730	49.11
human_chr11	3,316	0.01
gi 357579577 Canis_lupus_familiaris_chr3	270	0.00
human_chr2	266	0.00
human_chr8	12	0.00
human_chr4	12	0.00
human_chr13	10	0.00

human_chr7	10	0.00
gi 357579535 Canis_lupus_familiaris_chr20	10	0.00
gi 357579571 Canis_lupus_familiaris_chr5	10	0.00
human_chr5	6	0.00
human_chr16	6	0.00
human_chr14	4	0.00
human_chr3	4	0.00
human_chr9	4	0.00
human_chr20	4	0.00
human_chr1	2	0.00
human_chr18	2	0.00
human_chr22	2	0.00
human_chr15	2	0.00
gi 357579523 Canis_lupus_familiaris_chr27	2	0.00
human_chr17	2	0.00
human_chr21	2	0.00
human_chrX	2	0.00
gi 357579551 Canis_lupus_familiaris_chr11	2	0.00
human_chr12	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

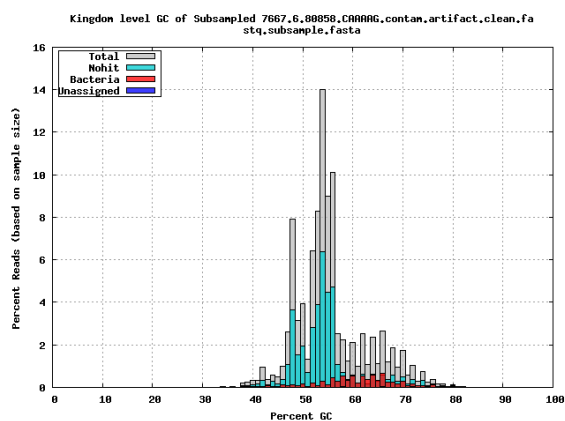
Description	Num Reads	Pct Reads
Input	28,183,458	100
Contam identified	8	0.0

List of Contaminants Identified

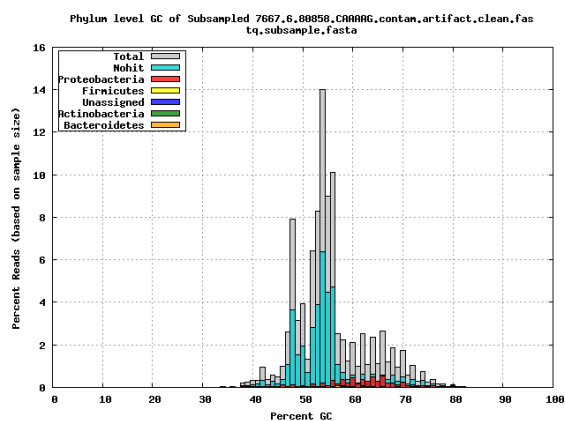
Description	Num Reads	Pct Reads
<i>Delftia</i>	4	0.00
<i>Escherichia</i>	2	0.00
<i>Klebsiella</i>	2	0.00

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

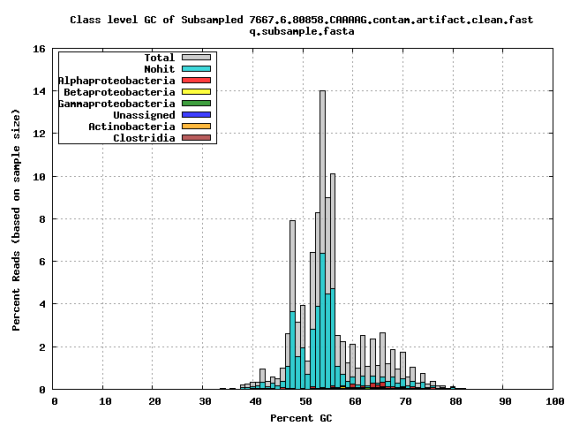
Kingdom



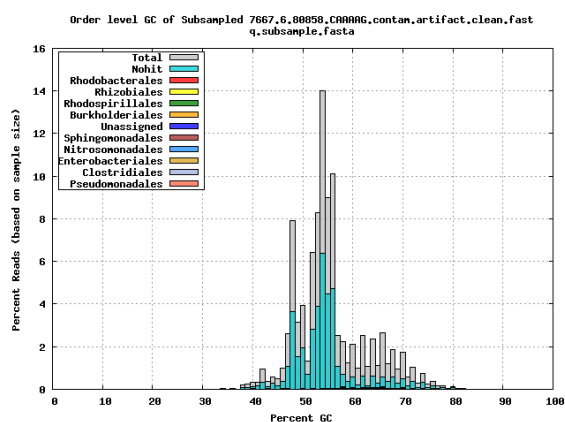
Phylum



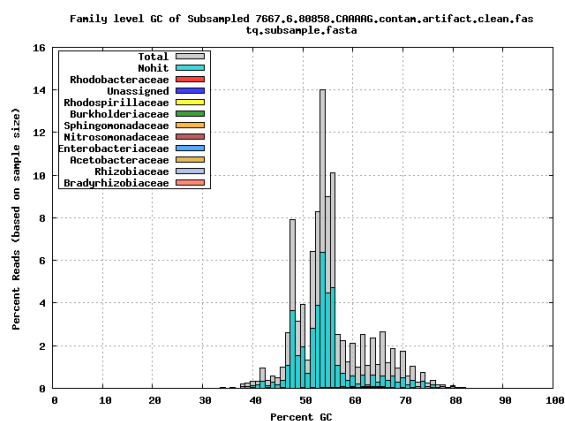
Class



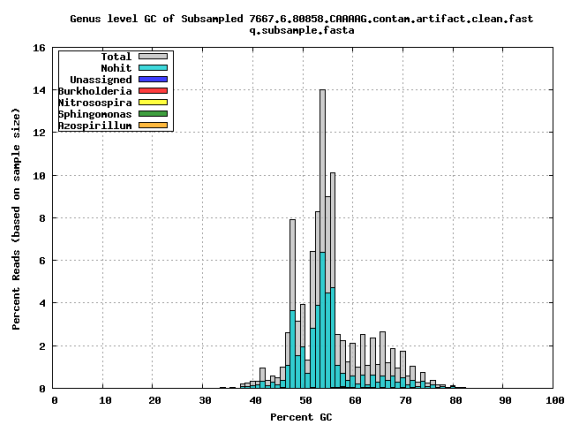
Order

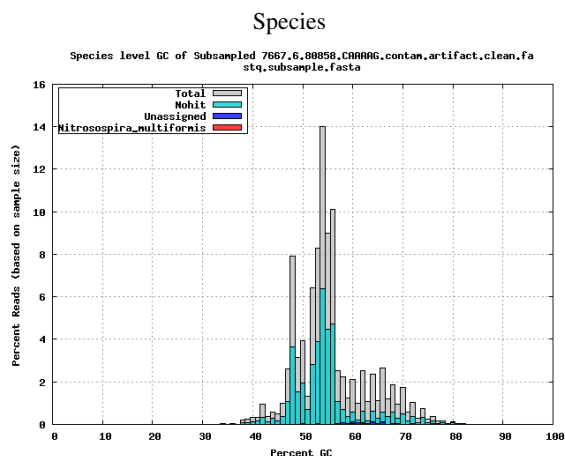


Family



Genus



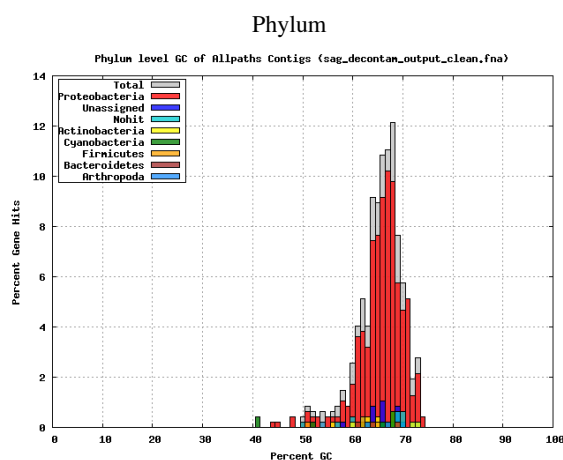
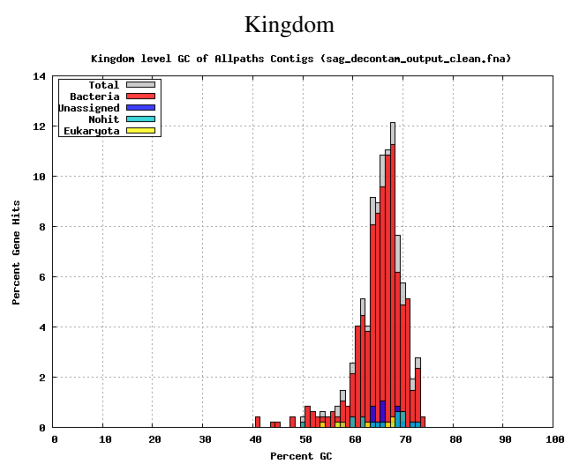


4. Assembly Statistics

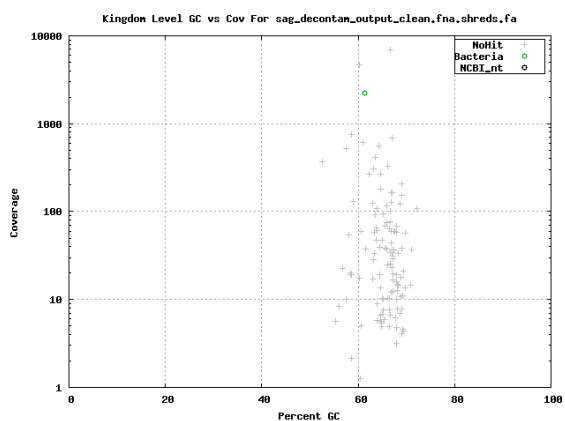
Assembly method	SPAdes with auto decontamination
Scaffold total	50
Contig total	50
Scaffold sequence length	494.6 kb
Contig sequence length	494.6 kb (0.0% gap)
Scaffold N/L50	10/18.0 kb
Contig N/L50	10/18.0 kb
Largest Contig	53.0 kb
Number of scaffolds >50 kb	1
Pct of genome in scaffolds >50 kb	10.7
Pct of reads assembled (raw)	27.3
Pct of reads assembled (decontam)	4.6

5. Assembly QC Results

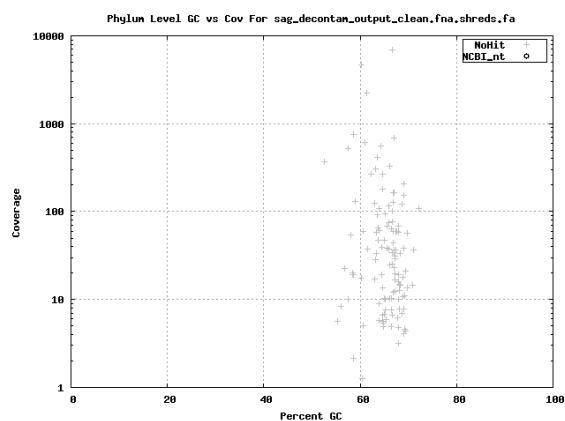
GC histogram of the predicted genes on each contig, overlaid with GC of hits based on BLASTP, shown for different taxonomic levels.



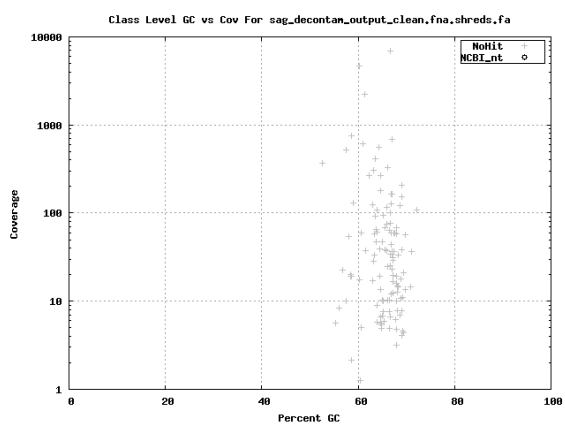
Kingdom



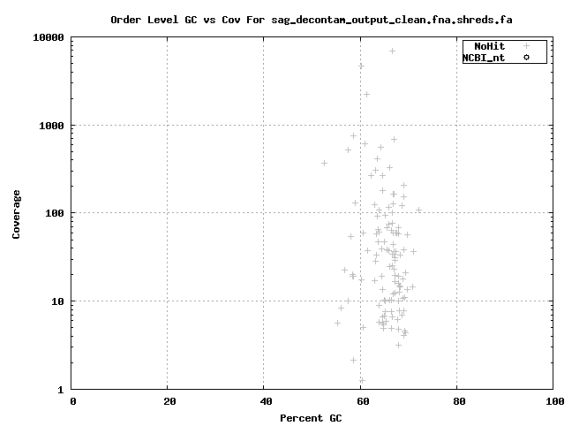
Phylum



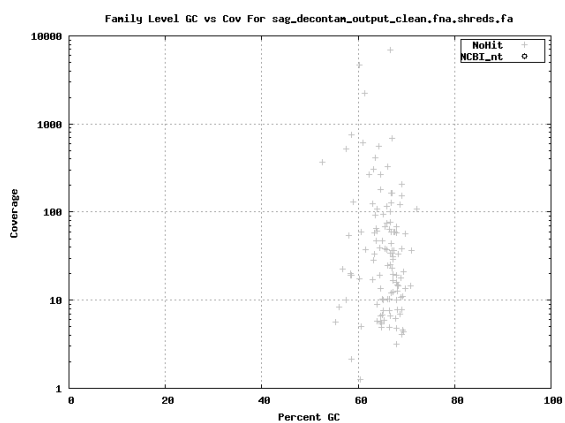
Class



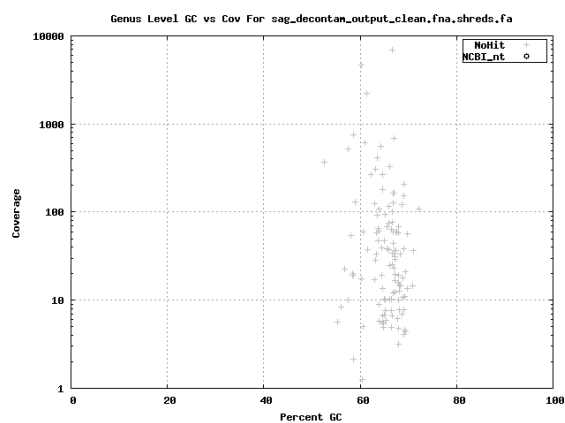
Order

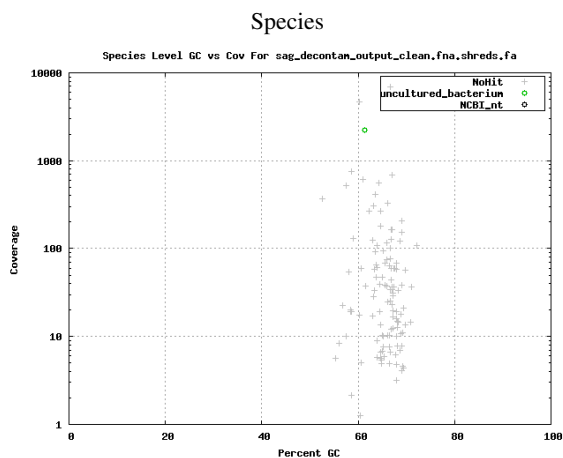


Family

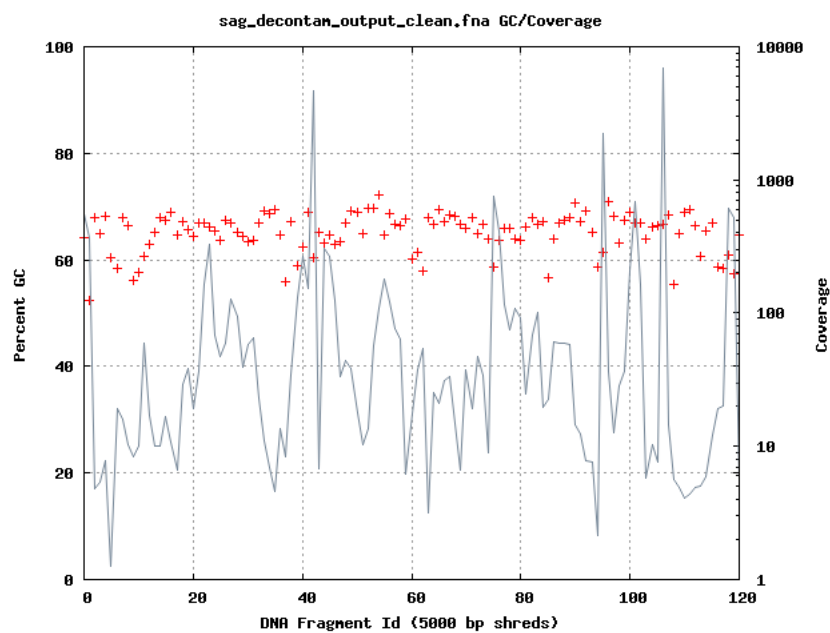


Genus

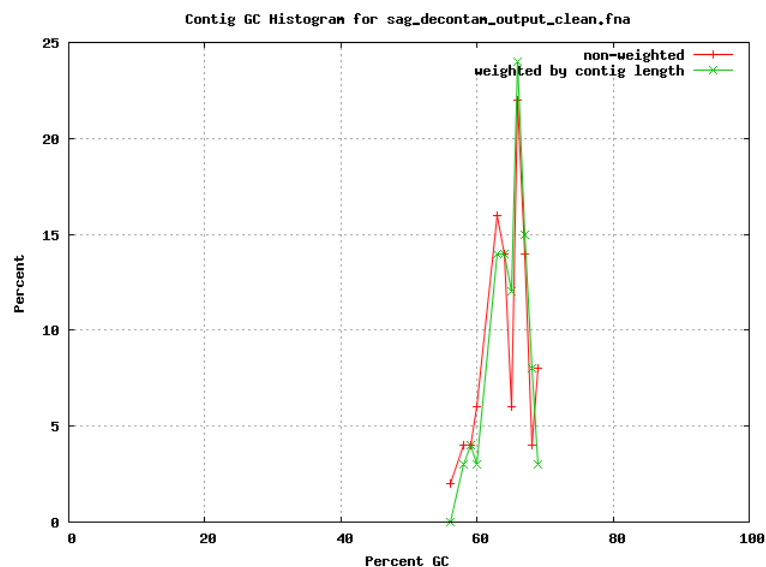




Coverage vs GC. Contigs were shredded into non-overlapping 5kbp and the GC of each shred was plotted as a point, colored by scaffold id. Coverage was calculated by mapping the fragment library to the final assembly and plotted as connected points.



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



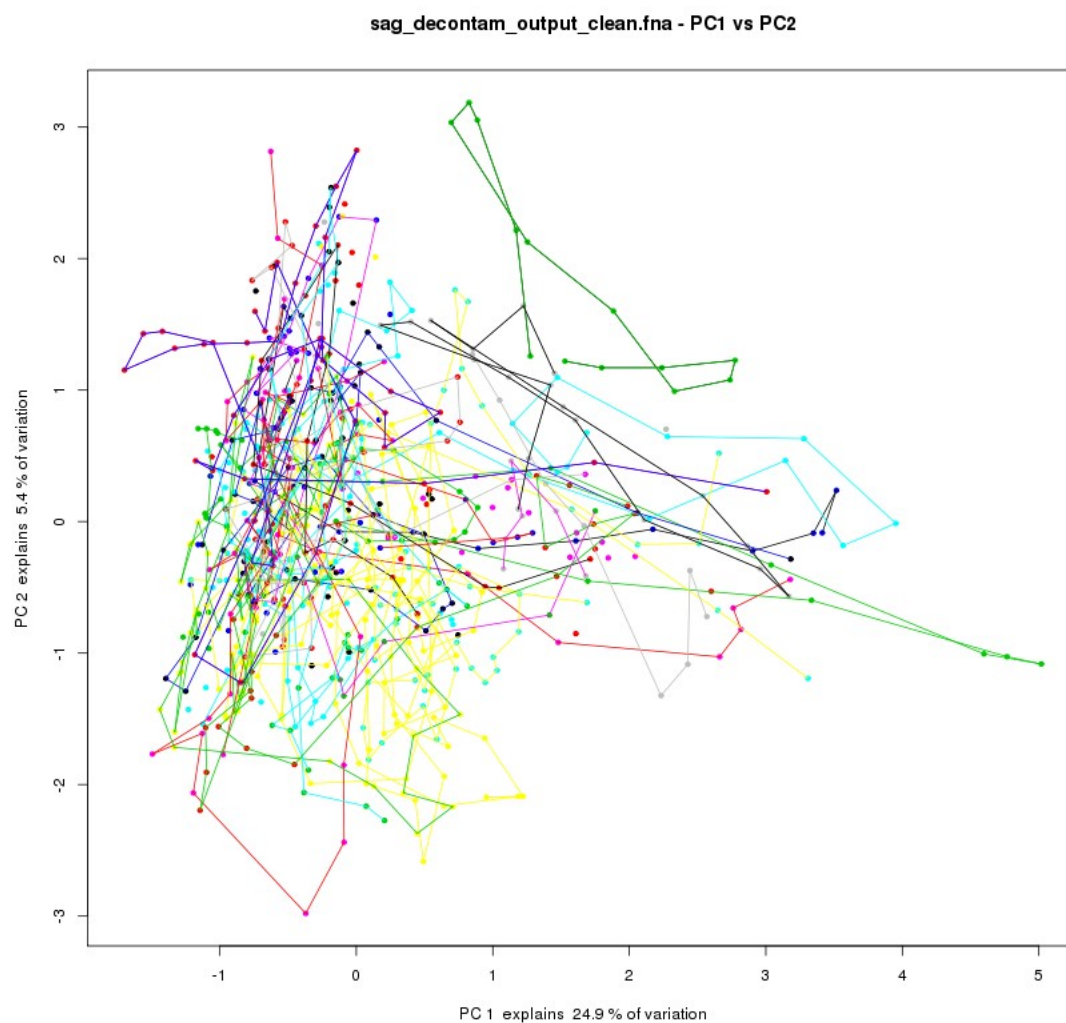
List of contigs and average percent GC, grouped in bins of 5:

Pct GC Bin	Contig Name
55	NODE.21.length.10836.cov.347.856.ID.41, NODE.26.length.8172.cov.18.9261.ID.55, NODE.27.length.7999.cov.323.198.ID.57, NODE.32.length.5986.cov.13.1074.ID.69 NODE.92.length.2042.cov.5.57272.ID.185
60	NODE.3.length.26477.cov.33.0196.ID.5, NODE.4.length.25374.cov.193.264.ID.7, NODE.7.length.20993.cov.318.772.ID.13, NODE.19.length.12069.cov.10.7419.ID.37, NODE.24.length.9379.cov.131.196.ID.51, NODE.28.length.7476.cov.7.31909.ID.63, NODE.29.length.7312.cov.40.964.ID.65, NODE.34.length.5638.cov.90.8506.ID.71, NODE.35.length.5633.cov.11.0477.ID.73, NODE.38.length.5532.cov.4.50685.ID.79, NODE.44.length.4590.cov.3.34443.ID.91, NODE.45.length.4554.cov.11.377.ID.93, NODE.46.length.4542.cov.21.7631.ID.95, NODE.47.length.4530.cov.5.7048.ID.97, NODE.50.length.4249.cov.204.478.ID.103, NODE.57.length.3619.cov.3.47475.ID.115, NODE.75.length.2582.cov.3.0922.ID.151, NODE.78.length.2361.cov.3.78925.ID.157
65	NODE.1.length.53014.cov.36.6049.ID.1, NODE.2.length.33666.cov.57.6013.ID.3, NODE.5.length.23884.cov.146.627.ID.9, NODE.6.length.23261.cov.16.8127.ID.11, NODE.8.length.19964.cov.17.718.ID.15, NODE.9.length.18496.cov.36.6417.ID.17, NODE.10.length.18035.cov.7.36085.ID.19, NODE.13.length.13919.cov.534.599.ID.25, NODE.14.length.13560.cov.21.276.ID.27, NODE.16.length.12509.cov.98.2824.ID.31, NODE.18.length.12153.cov.20.2937.ID.35, NODE.23.length.9640.cov.5.83944.ID.49, NODE.37.length.5555.cov.2.81618.ID.77, NODE.40.length.5451.cov.4.67661.ID.83, NODE.42.length.4826.cov.23.813.ID.87, NODE.48.length.4512.cov.4.79403.ID.99, NODE.52.length.3924.cov.4505.62.ID.107, NODE.56.length.3668.cov.3.05536.ID.113, NODE.64.length.3126.cov.7.65549.ID.129, NODE.66.length.3022.cov.4.20998.ID.133, NODE.71.length.2720.cov.7.92083.ID.143, NODE.72.length.2677.cov.3.92105.ID.145, NODE.73.length.2630.cov.3.72738.ID.147, NODE.82.length.2298.cov.2.16139.ID.165, NODE.89.length.2072.cov.2.93009.ID.179, NODE.90.length.2062.cov.2.70055.ID.181 NODE.93.length.2039.cov.4.1124.ID.187

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

Organism	Align Length (bp)	Pct Id	Contig Name
<i>Pseudomonas aeruginosa</i> .UCBPP_PA14.complete.genome	202	90.59	NODE.13.length.13919.cov.534.599.ID.25

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	7.19 %
archaea	2.06 %

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 28,183,458 reads totaling 4,227.5 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 50 contigs in 50 scaffolds, totalling 494.6 Kb in size. The final assembly was based on 1,925.2 Mb of Illumina data. Based on a presumed genome size of 5.0 Mb, the average input read coverage used for the assembly was 385.0X.

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.