*uncultured archaeon*

Single–cell Assembly QC Report                                        04/07/2014

# 1. Project Information

| | |
|---|---|
| Program | Microbial/CSP 2012 |
| PMO Project | 0 |
| Seq Proj ID | 1027076 |
| Sequencing Project Name | Geoarchaeota archaeon JGI 000156CP–L10 |
| JGI Project ID | 0 |

# 2. Read Statistics

Illumina Std PE Statistics

| | |
|---|---|
| File name | 7666.6.80850.ATGTCA.fastq |
| Library | TGSS |
| Number of reads | 28,857,866 |
| Sequencing depth $^{\dagger}$ | 866X |
| Read type | 2x150 bp |

$^{\dagger}$ 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 | 28,857,866 | 100 |
| Contam removed | 86 | 0.0 |
| Artifact removed | 311,514 | 1.1 |
| Total removed | 8,857,866 | 30.7 |
| Total remaining | 20,000,000 | 69.3 |

List of Contaminants Removed

| Description | Num Reads | Pct Reads |
|---|---|---|
| gi\|357579535\|Canis_lupus_familiaris_chr20 | 44 | 0.00 |
| human_chr2 | 34 | 0.00 |
| gi\|357579577\|Canis_lupus_familiaris_chr3 | 32 | 0.00 |
| gi\|357579571\|Canis_lupus_familiaris_chr5 | 24 | 0.00 |
| human_chr18 | 2 | 0.00 |
| human_chr5 | 2 | 0.00 |
| human_chr3 | 2 | 0.00 |

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

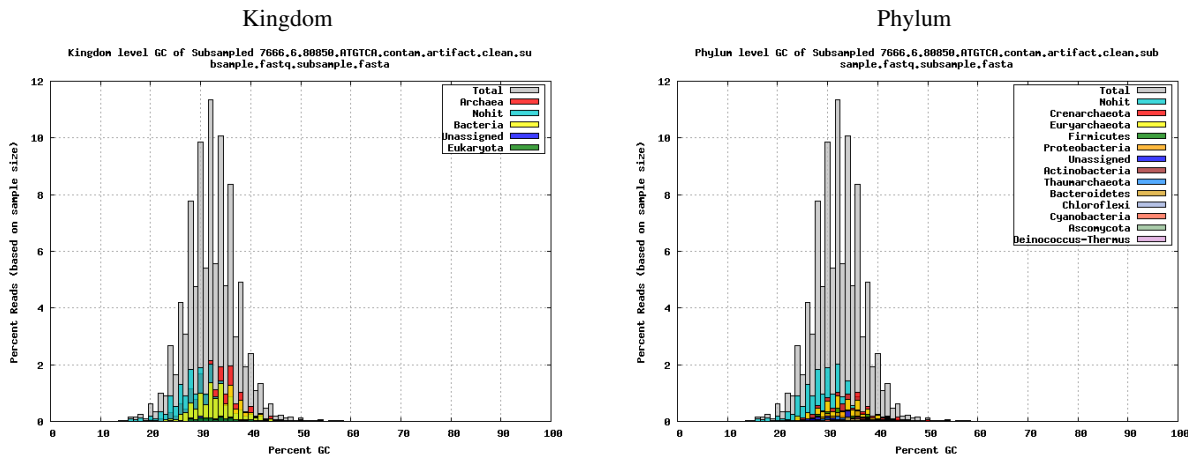The following are the results of reads screened against potential reagent and process contaminants but were not removed from the dataset.

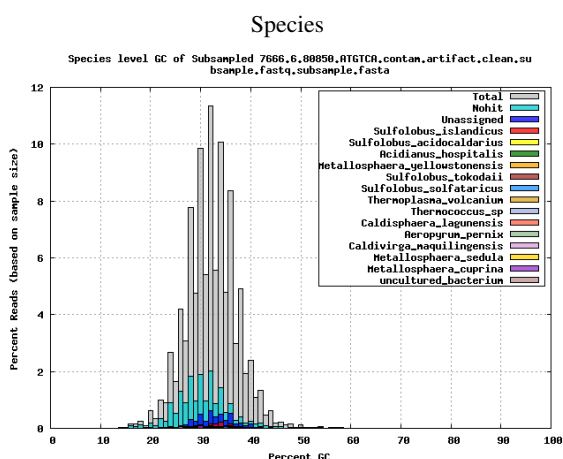Illumina Std PE Contamination Identification Statistics

| Description | Num Reads | Pct Reads |
| --- | --- | --- |
| Input | 28,857,866 | 100 |
| Contam identified | 10 | 0.0 |

List of Contaminants Identified

| Description | Num Reads | Pct Reads |
| --- | --- | --- |
| *Escherichia* | 2 | 0.00 |
| *Delftia* | 2 | 0.00 |
| *Pseudomonas* | 2 | 0.00 |
| *Shigella* | 2 | 0.00 |
| *Ralstonia* | 2 | 0.00 |

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

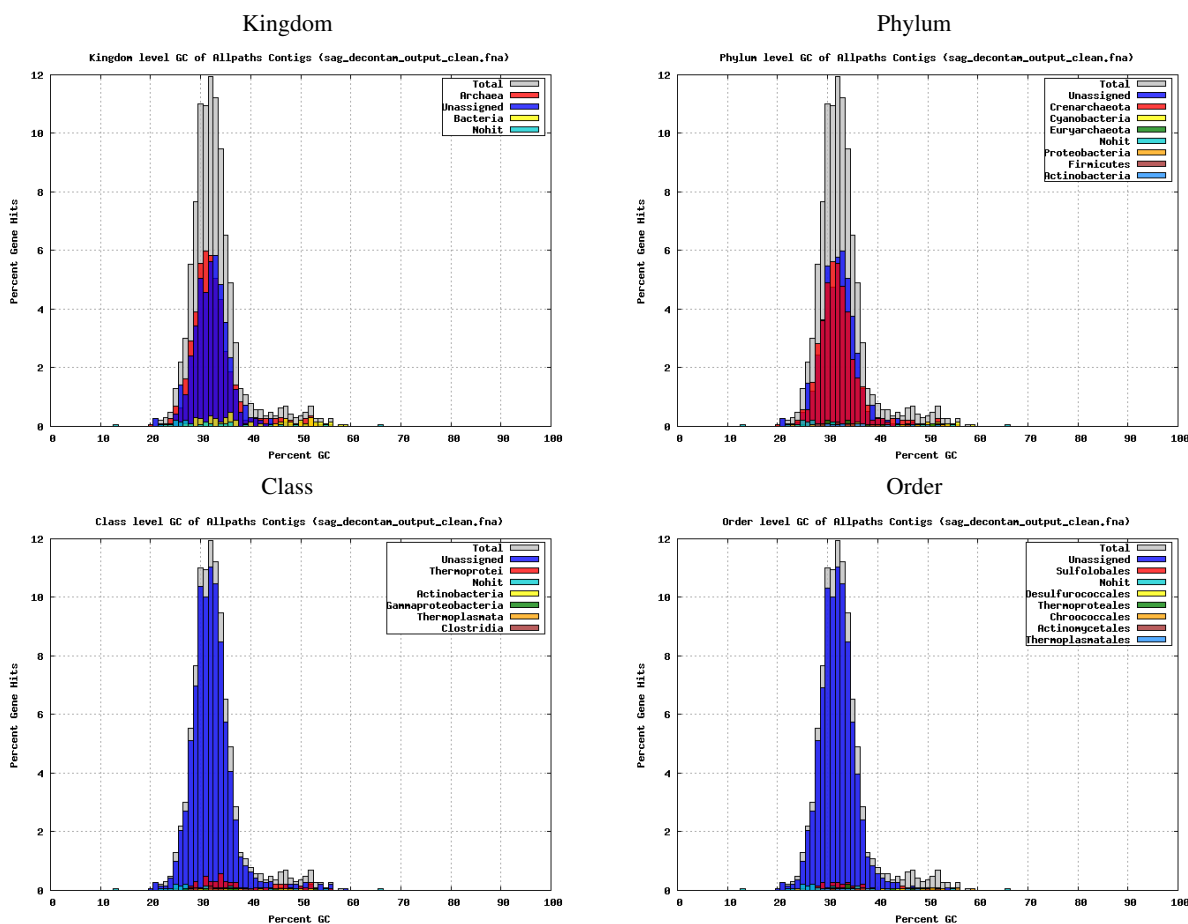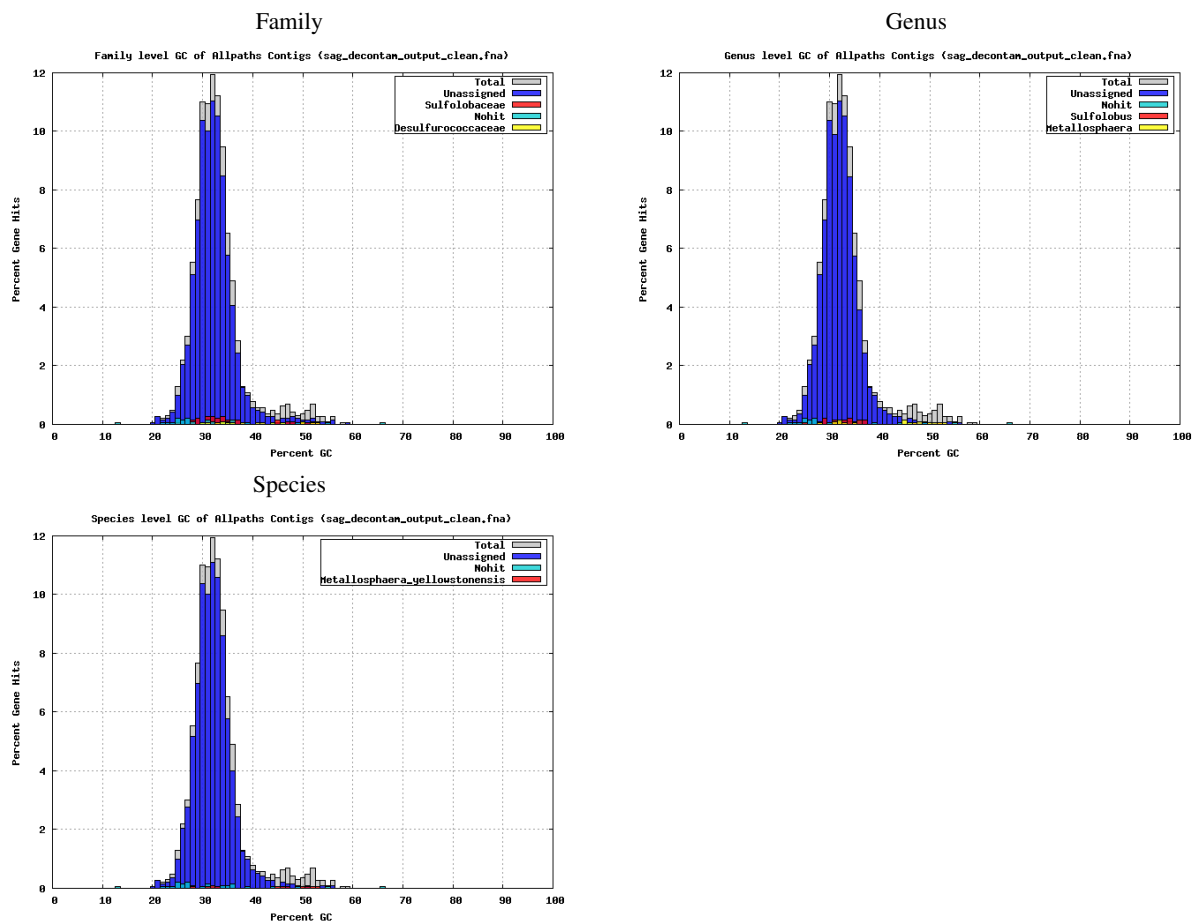## Class

Class level GC of Subsampled 7666.6.80850.ATGTCA.contam.artifact.clean.subsample.fastq.subsample.fasta

Legend: Total, Nohit, Thermoprotei, Unassigned, Bacilli, Thermococci, Clostridia, Actinobacteria, Halobacteria, Thermoplasmata, Alphaproteobacteria, Gammaproteobacteria, Methanomicrobia, Methanococci, Betaproteobacteria, Archaeoglobi, Methanobacteria, Deltaproteobacteria, Flavobacteriia, Deinococci

Y-axis: Percent Reads (based on sample size)
X-axis: Percent GC

## Order

Order level GC of Subsampled 7666.6.80850.ATGTCA.contam.artifact.clean.subsample.fastq.subsample.fasta

Legend: Total, Nohit, Sulfolobales, Unassigned, Thermoproteales, Desulfurococcales, Thermococcales, Bacillales, Actinomycetales, Halobacteriales, Clostridiales, Thermoplasmatales, Methanococcales, Archaeoglobales, Methanobacteriales, Rhizobiales, Acidilobales, Burkholderiales, Thermoanaerobacterales, Methanosarcinales, Flavobacteriales, Nitrosopumilales, Lactobacillales, Pseudomonadales, Methanomicrobiales

Y-axis: Percent Reads (based on sample size)
X-axis: Percent GC

## Family

Family level GC of Subsampled 7666.6.80850.ATGTCA.contam.artifact.clean.subsample.fastq.subsample.fasta

Legend: Total, Nohit, Sulfolobaceae, Unassigned, Thermoproteaceae, Thermococcaceae, Desulfurococcaceae, Halobacteriaceae, Bacillaceae, Thermoplasmataceae, Archaeoglobaceae, Methanobacteriaceae, Methanocaldococcaceae, Clostridiaceae, Nitrosopumilaceae, Flavobacteriaceae, Caldisphaeraceae, Mycobacteriaceae, Burkholderiaceae, Methanococcaceae, Methanosarcinaceae, Paenibacillaceae, Peptococcaceae

Y-axis: Percent Reads (based on sample size)
X-axis: Percent GC

## Genus

Genus level GC of Subsampled 7666.6.80850.ATGTCA.contam.artifact.clean.subsample.fastq.subsample.fasta

Legend: Total, Nohit, Sulfolobus, Unassigned, Pyrobaculum, Metallosphaera, Thermococcus, Candidatus_Caldiarchaeum, Pyrococcus, Bacillus, Thermoplasma, Vulcanisaeta, Acidianus, Archaeoglobus, Thermoproteus, Methanocaldococcus, Clostridium, Caldisphaera, Aeropyrum, Desulfurococcus, Caldivirga, Staphylothermus, Mycobacterium, Methanococcus, Methanobacterium

Y-axis: Percent Reads (based on sample size)
X-axis: Percent GC

## Species

Species level GC of Subsampled 7666.6.80850.ATGTCA.contam.artifact.clean.subsample.fastq.subsample.fasta

Legend: Total, Nohit, Unassigned, Sulfolobus_islandicus, Sulfolobus_acidocaldarius, Acidianus_hospitalis, Metallosphaera_yellowstonensis, Sulfolobus_tokodaii, Sulfolobus_solfataricus, Thermoplasma_volcanium, Thermococcus_sp, Caldisphaera_lagunensis, Aeropyrum_pernix, Caldivirga_maquilingensis, Metallosphaera_sedula, Metallosphaera_cuprina, uncultured_bacterium

Y-axis: Percent Reads (based on sample size)
X-axis: Percent GC

3

# 4. Assembly Statistics

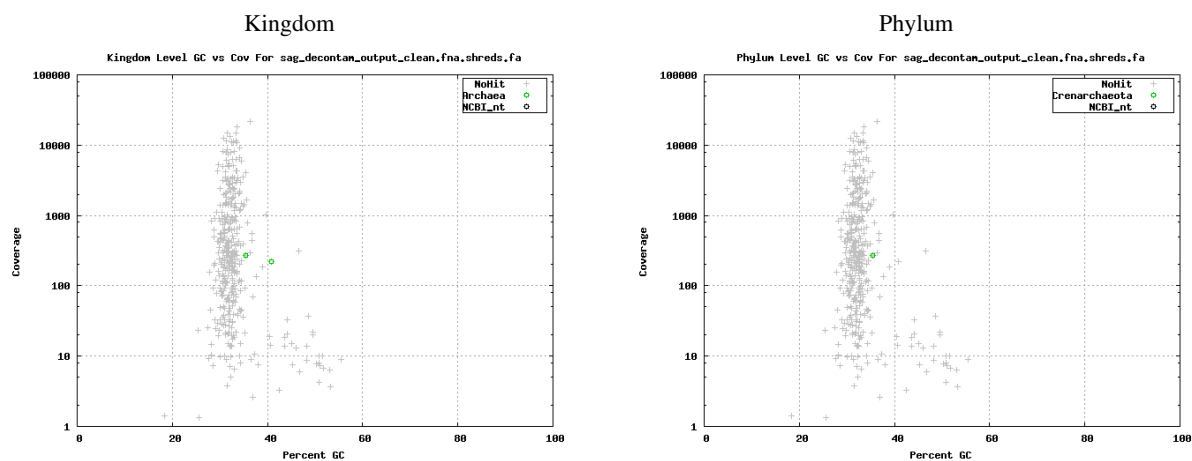| Assembly method | SPAdes with auto decontamination |
|---|---|
| Scaffold total | 89 |
| Contig total | 89 |
| Scaffold sequence length | 1.7 Mb |
| Contig sequence length | 1.7 Mb ( 0.0% gap) |
| Scaffold N/L50 | 18/32.6 kb |
| Contig N/L50 | 18/32.6 kb |
| Largest Contig | 75.9 kb |
| Number of scaffolds >50 kb | 6 |
| Pct of genome in scaffolds >50 kb | 22.7 |
| Pct of reads asssembled (raw) | 83.1 |
| Pct of reads asssembled (decontam) | 77.4 |

# 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.
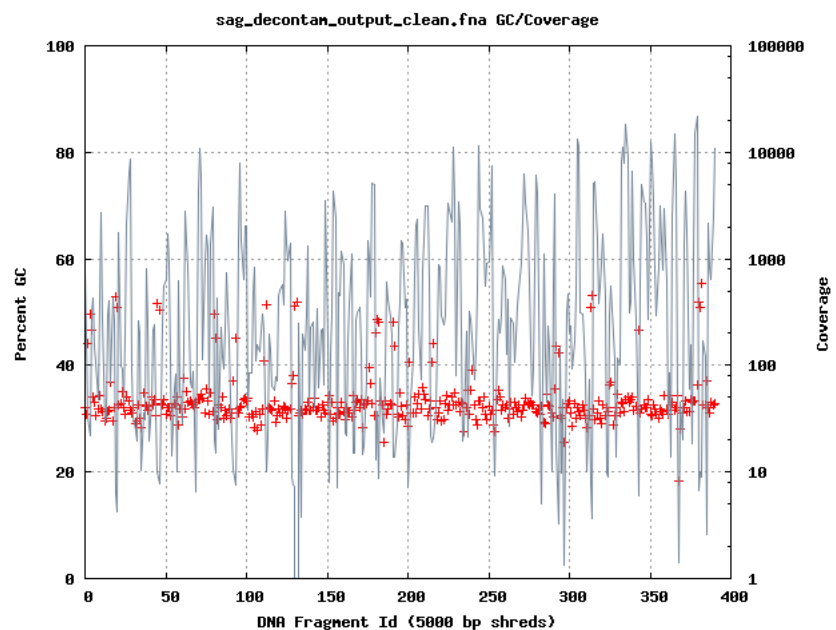
Kingdom

Phylum

Class

Order

## Family

Family level GC of Allpaths Contigs (sag_decontam_output_clean.fna)



## Genus

Genus level GC of Allpaths Contigs (sag_decontam_output_clean.fna)



## Species

Species level GC of Allpaths Contigs (sag_decontam_output_clean.fna)



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.
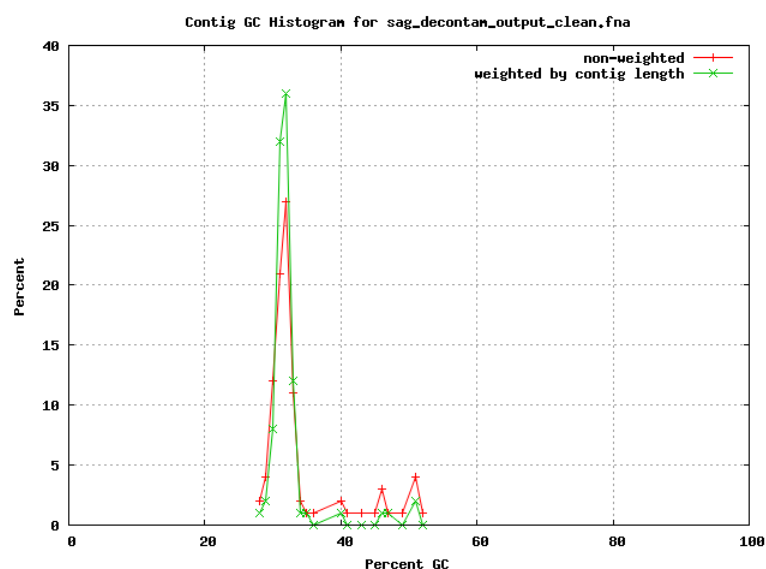
## Kingdom

Kingdom Level GC vs Cov For sag_decontam_output_clean.fna.shreds.fa



## Phylum

Phylum Level GC vs Cov For sag_decontam_output_clean.fna.shreds.fa

## Class



Class Level GC vs Cov For sag_decontam_output_clean.fna.shreds.fa

## Order



Order Level GC vs Cov For sag_decontam_output_clean.fna.shreds.fa

## Family



Family Level GC vs Cov For sag_decontam_output_clean.fna.shreds.fa

## Genus



Genus Level GC vs Cov For sag_decontam_output_clean.fna.shreds.fa

## Species



Species Level GC vs Cov For sag_decontam_output_clean.fna.shreds.fa

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 asssembly and plotted as connected points.

**sag_decontam_output_clean.fna GC/Coverage**

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



**Contig GC Histogram for sag_decontam_output_clean.fna**

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

| Pct GC Bin | Contig Name |
|---|---|
| 25 | NODE_50_length_11236_cov_95.8108_ID_103, NODE_71_length_6503_cov_192.994_ID_149, NODE_73_length_6314_cov_26.4531_ID_153, NODE_76_length_6056_cov_18.8212_ID_159, NODE_80_length_5604_cov_281.579_ID_165, NODE_96_length_4878_cov_281.327_ID_211 |
| 30 | NODE_1_length_75851_cov_758.926_ID_1, NODE_2_length_73144_cov_398.274_ID_3, NODE_3_length_66427_cov_1778.35_ID_5, NODE_4_length_61123_cov_824.336_ID_7, NODE_5_length_55719_cov_1079.9_ID_9, NODE_6_length_52032_cov_282.999_ID_11, NODE_7_length_48417_cov_185.684_ID_13, NODE_8_length_44711_cov_2190.23_ID_15, |

7

|    | |
|----|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
|    | NODE_9_length_41378_cov_144.777_ID_17, NODE_10_length_41321_cov_2356.69_ID_19, NODE_11_length_38327_cov_1375.49_ID_21, NODE_12_length_38195_cov_2169.82_ID_23, NODE_13_length_37190_cov_490.381_ID_25, NODE_14_length_36550_cov_255.024_ID_27, NODE_15_length_36476_cov_407.507_ID_29, NODE_16_length_35874_cov_346.141_ID_31, NODE_17_length_34702_cov_1520.27_ID_33, NODE_18_length_32580_cov_2522.65_ID_35, NODE_19_length_32126_cov_5783.81_ID_37, NODE_20_length_31303_cov_124.723_ID_39, NODE_21_length_28042_cov_1703.95_ID_41, NODE_22_length_25252_cov_163.011_ID_43, NODE_24_length_23569_cov_218.285_ID_47, NODE_25_length_22378_cov_163.706_ID_49, NODE_26_length_21957_cov_136.72_ID_51, NODE_27_length_21189_cov_1689.07_ID_53, NODE_28_length_20919_cov_1375.65_ID_55, NODE_29_length_20293_cov_249.362_ID_57, NODE_30_length_18835_cov_153.268_ID_59, NODE_31_length_18776_cov_294.389_ID_61, NODE_33_length_17009_cov_1556.15_ID_65, NODE_34_length_16095_cov_261.278_ID_67, NODE_35_length_15495_cov_85.551_ID_69, NODE_36_length_15479_cov_111.284_ID_71, NODE_37_length_15322_cov_156.297_ID_73, NODE_38_length_15090_cov_374.476_ID_75, NODE_39_length_14457_cov_1554.54_ID_77, NODE_40_length_14420_cov_225.012_ID_79, NODE_41_length_14090_cov_5240.78_ID_81, NODE_42_length_13252_cov_137.758_ID_83, NODE_43_length_13007_cov_940.997_ID_89, NODE_44_length_12715_cov_836.495_ID_91, NODE_45_length_12273_cov_292.523_ID_93, NODE_46_length_12220_cov_732.621_ID_95, NODE_47_length_12113_cov_502.159_ID_97, NODE_48_length_12048_cov_219.225_ID_99, NODE_51_length_10961_cov_88.4575_ID_105, NODE_53_length_10542_cov_35.6054_ID_109, NODE_54_length_10338_cov_1806.05_ID_111, NODE_55_length_10227_cov_102.843_ID_85, NODE_56_length_10187_cov_232.742_ID_113, NODE_57_length_9837_cov_145.557_ID_115, NODE_59_length_9010_cov_3456.74_ID_119, NODE_61_length_8339_cov_26.154_ID_131, NODE_63_length_8142_cov_63.013_ID_135, NODE_64_length_7869_cov_378.122_ID_137, NODE_66_length_7395_cov_103.374_ID_141, NODE_70_length_6527_cov_322.241_ID_147, NODE_74_length_6248_cov_761.281_ID_155, NODE_75_length_6082_cov_18.5444_ID_157, NODE_78_length_5704_cov_8931.69_ID_163, NODE_79_length_5608_cov_15.6231_ID_123, NODE_83_length_5497_cov_281.99_ID_171, NODE_87_length_5223_cov_30.1418_ID_187, NODE_89_length_5142_cov_29.3668_ID_195, NODE_92_length_4960_cov_48.8377_ID_203 |
| 35 | NODE_23_length_24734_cov_413.486_ID_45, NODE_65_length_7397_cov_5.7088_ID_139 |
| 40 | NODE_60_length_8601_cov_6.30084_ID_129, NODE_68_length_6826_cov_14.0418_ID_145, NODE_85_length_5278_cov_11.7722_ID_179, NODE_90_length_4979_cov_9.48457_ID_185 |
| 45 | NODE_49_length_11800_cov_15.229_ID_101, NODE_58_length_9826_cov_18.4541_ID_117, NODE_62_length_8262_cov_9.36847_ID_133, NODE_77_length_5722_cov_14.1288_ID_161, NODE_81_length_5574_cov_229.352_ID_167, NODE_91_length_4971_cov_3.97803_ID_201 |
| 50 | NODE_52_length_10744_cov_5.75208_ID_107, NODE_69_length_6798_cov_6.3365_ID_127, NODE_72_length_6418_cov_4.84441_ID_151, NODE_82_length_5527_cov_4.11056_ID_169 NODE_86_length_5258_cov_4.90698_ID_181 |

Principal component analysis of tetramer frequencies of contigs. Detectable variations are highlighted in color.

sag_decontam_output_clean.fna - PC1 vs PC2

Estimated genome recovery derived from analysis of universal single-copy genes detected in final assembly.

| HMM | Pct Recovered |
|---|---|
| bacteria | 49.56 % |
| archaea | 100 % |

# 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 Availiability

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 28,857,866 reads totaling 4,328.7 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 89 contigs in 89 scaffolds, totalling 1.7 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.