

SAG quality control

Current QC reporting standards

PERSPECTIVE

nature
biotechnology
OPEN

Minimum information about a single amplified genome (MISAG) and a metagenome-assembled genome (MIMAG) of bacteria and archaea

Table 1 Genome reporting standards

Criterion	Definition
Finished	Assembly quality ^a : >99% of the genome is covered by contigs with >95% sequence identity. Completion ^b : >99% of the genome is covered by contigs with >95% sequence identity. Contamination ^c : <5%.
High-quality draft (SAG/MAG)	Assembly quality ^a : Multiple fragments where gaps span repetitive regions. Presence of the 23S, 16S, and 5S rRNA genes and at least 18 tRNAs. Completion ^b : >90%. Contamination ^c : <5%.
Medium-quality draft (SAG/MAG)	Assembly quality ^a : Many fragments with little to no review or assembly other than reporting of standard assembly statistics. Completion ^b : <50%. Contamination ^c : <10%.
Low-quality draft (SAG/MAG)	This is a compressed representation of the data in Table 1. ^a Assembly quality: >99% of the genome is covered by contigs with >95% sequence identity. Completion: >99% of the genome is covered by contigs with >95% sequence identity. Contamination: <5% of the genome is covered by contigs with >95% sequence identity. Predicted genes per genome: >1000 predicted genes per genome. Genes in ≥2 copies to total genome: >10% genes in ≥2 copies to total genome.

Problem #1: Contamination limits too permissive

Problem #2: Misassemblies, indels, mismatches not covered

Problem #3: QC method validation not covered

Current QC reporting standards

Method

CheckM: assessing the quality of microbial genomes recovered from isolates, single cells, and metagenomes

Donovan H. Parks,¹ Michael Imelfort,¹ Connor T. Skennerton,¹ Philip Hugenholtz,^{1,2} and Gene W. Tyson^{1,3}

¹Australian Centre for Ecogenomics, School of Chemistry and Molecular Biosciences, The University of Queensland, St. Lucia, QLD 4072, Queensland, Australia; ²Institute for Molecular Bioscience, The University of Queensland, St. Lucia, QLD 4072, Queensland, Australia; ³Advanced Water Management Centre, The University of Queensland, St. Lucia, QLD 4072, Queensland, Australia

Contamination under-estimation by checkM

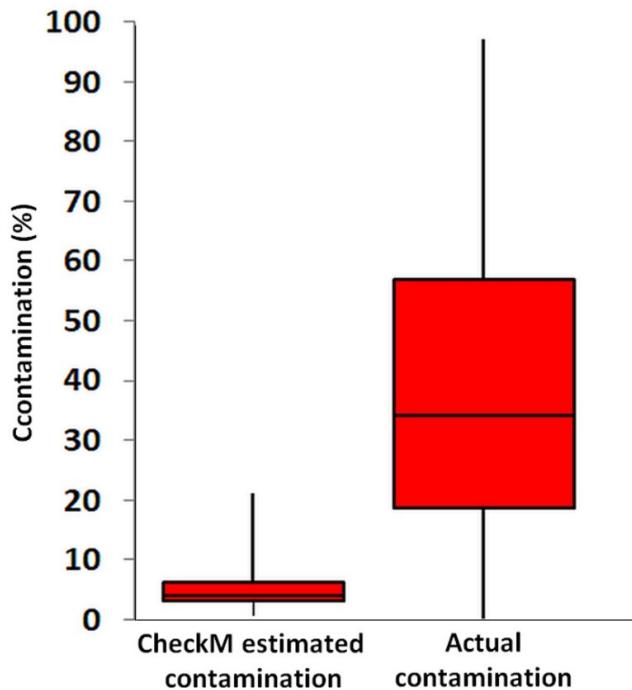


FIGURE 6 | Contamination predicted by CheckM software from all pairwise SAG combinations across the class-level lineages Rokumicrobia and Infratellusbacteria (15 Rokumicrobia and 4 Infratellusbacteria SAGs; 60 combinations total), compared to actual contamination calculated from all artificially combined SAGs. Two-sample *t*-test assuming equal variances was significant ($p = <000.1$).

From: Becroft et al. 2017

Reasons behind this discrepancy:

1. CheckM searches for multiple occurrences of same marker gene, not phylogeny conflicts.
2. Assembly incompleteness is not accounted for.
3. Reference database limitations may lead to detection failures of some marker genes .

Recent microbial genome quality re-evaluation

Orakov *et al.* *Genome Biology* (2021) 22:178
<https://doi.org/10.1186/s13059-021-02393-0>

Genome Biology

METHOD

Open Access

GUNC: detection of chimerism and contamination in prokaryotic genomes

Askarbek Orakov^{1†} , Anthony Fullam^{1†}, Luis Pedro Coelho^{2,3}, Supriya Khedkar¹, Damian Szkłarczyk^{4,5}, Daniel R. Mende⁶, Thomas S. B. Schmidt^{1*} and Peer Bork^{1,7,8,9*}



* Correspondence: sebastian.schmidt@embl.de; peer.bork@embl.org

[†]Askarbek Orakov and Anthony Fullam contributed equally to this work.

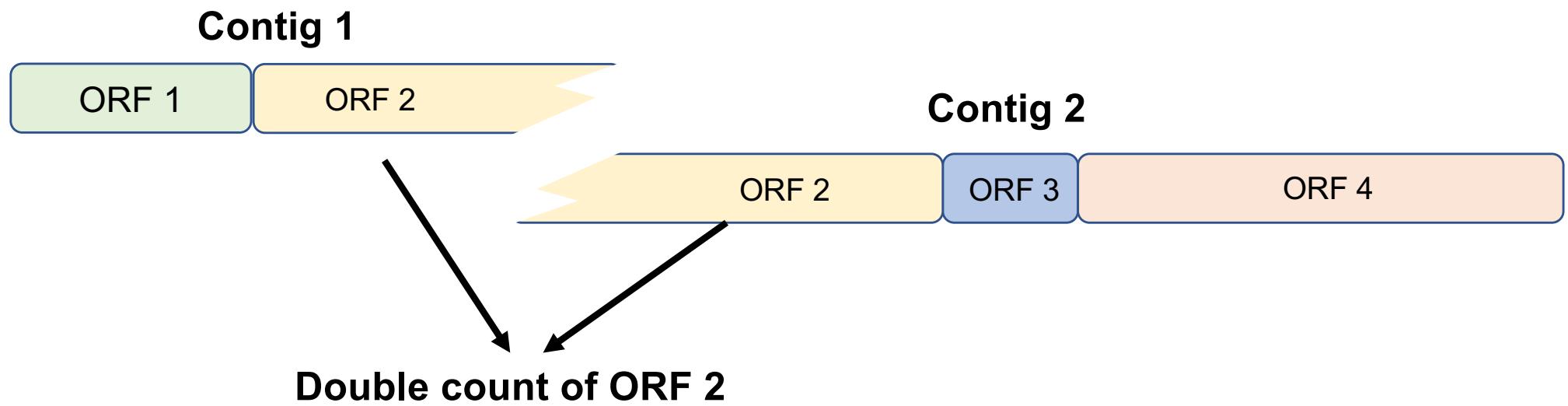
¹Structural and Computational Biology Unit, European Molecular Biology Laboratory, 69117 Heidelberg, Germany
Full list of author information is available at the end of the article

Abstract

Genomes are critical units in microbiology, yet ascertaining quality in prokaryotic genome assemblies remains a formidable challenge. We present GUNC (the Genome UNClutterer), a tool that accurately detects and quantifies genome chimerism based on the lineage homogeneity of individual contigs using a genome's full complement of genes. GUNC complements existing approaches by targeting previously underdetected types of contamination: we conservatively estimate that 5.7% of genomes in GenBank, 5.2% in RefSeq, and 15–30% of pre-filtered “high-quality” metagenome-assembled genomes in recent studies are undetected chimeras. GUNC provides a fast and robust tool to substantially improve prokaryotic genome quality.

Keywords: Genome quality, Genome contamination, Metagenomics, Metagenome-assembled genomes, Bioinformatics

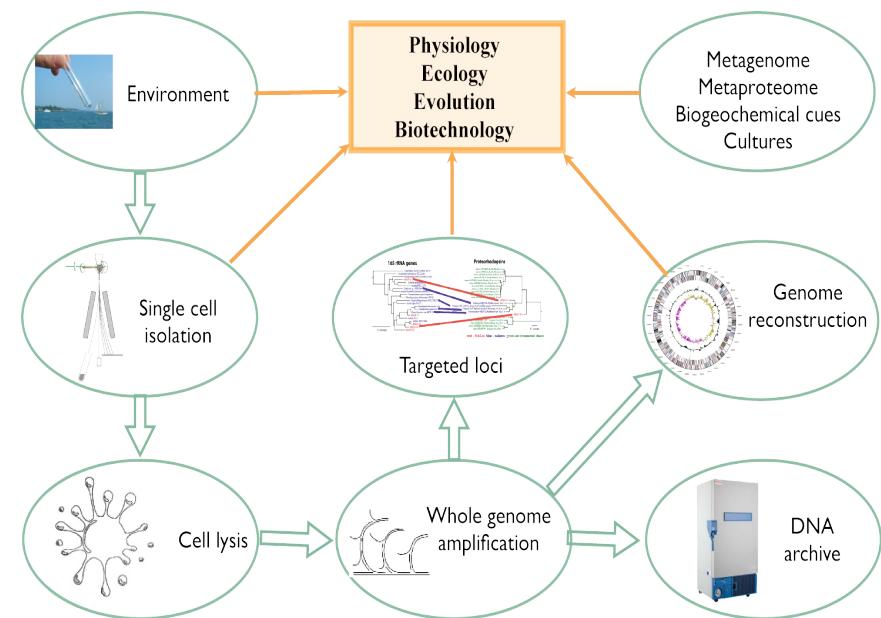
Contamination over-estimation by checkM



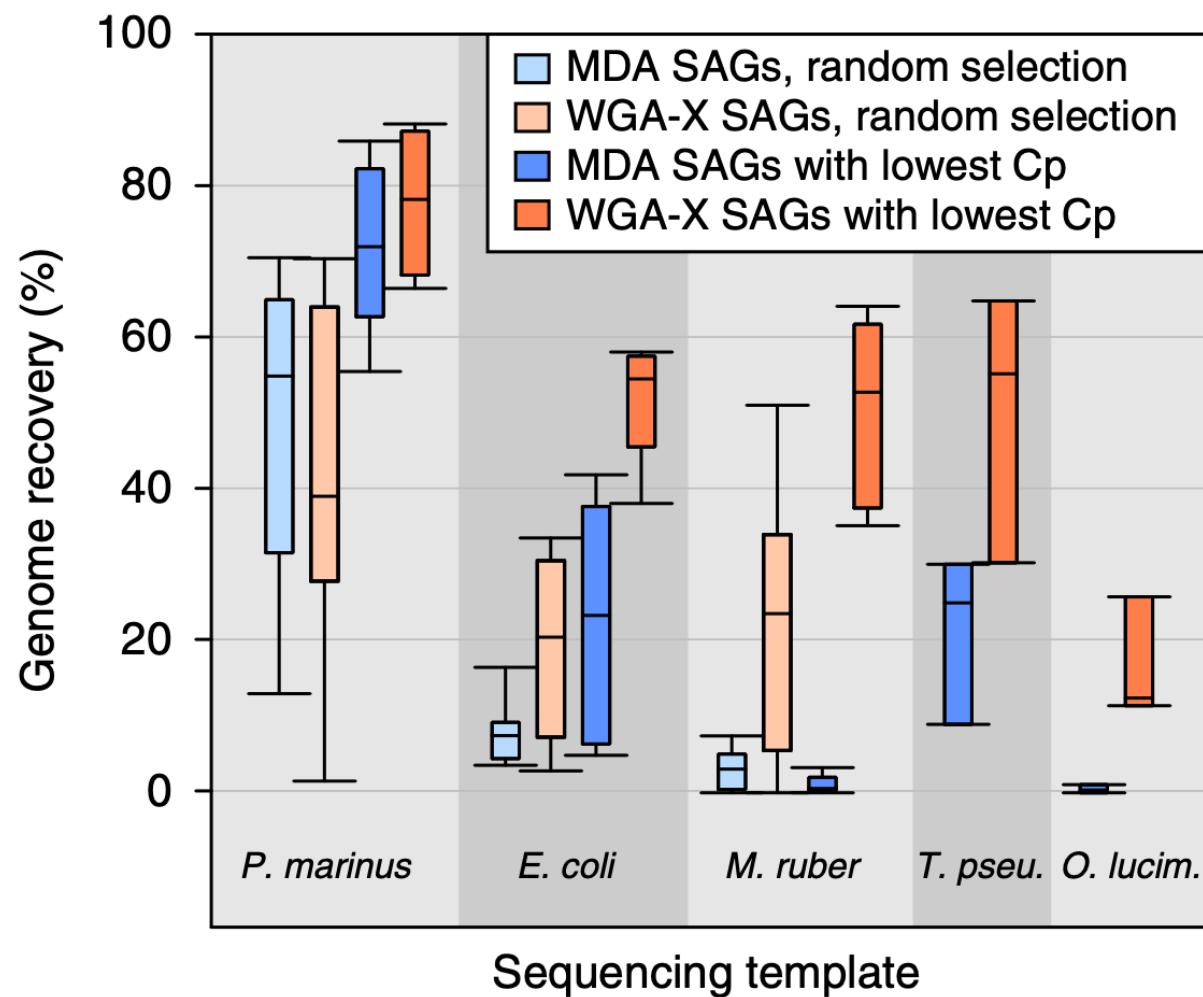
QA and QC at SCGC

SCGC strategy:

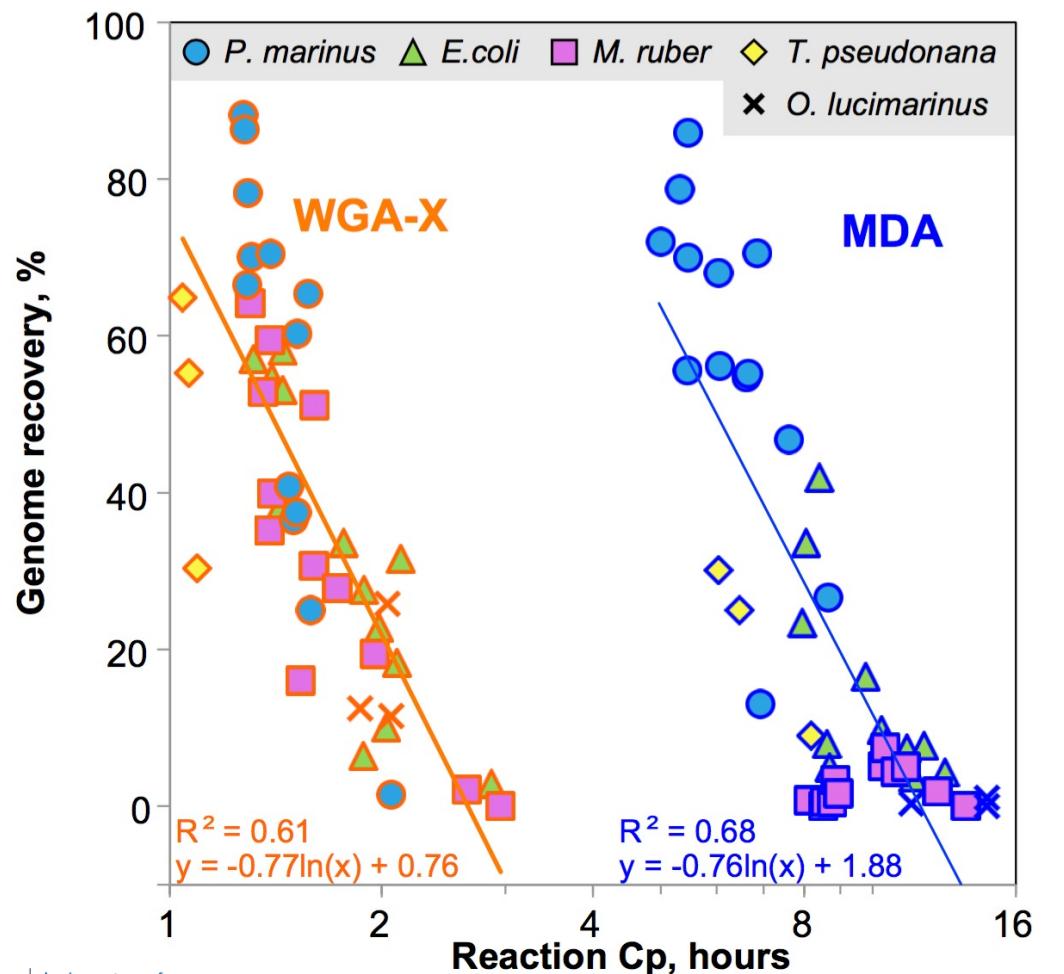
1. Validate the entire process
2. Minimize physical DNA contamination and damage
3. Maximize accuracy of the computational workflow
4. Generate data outputs for manual QC



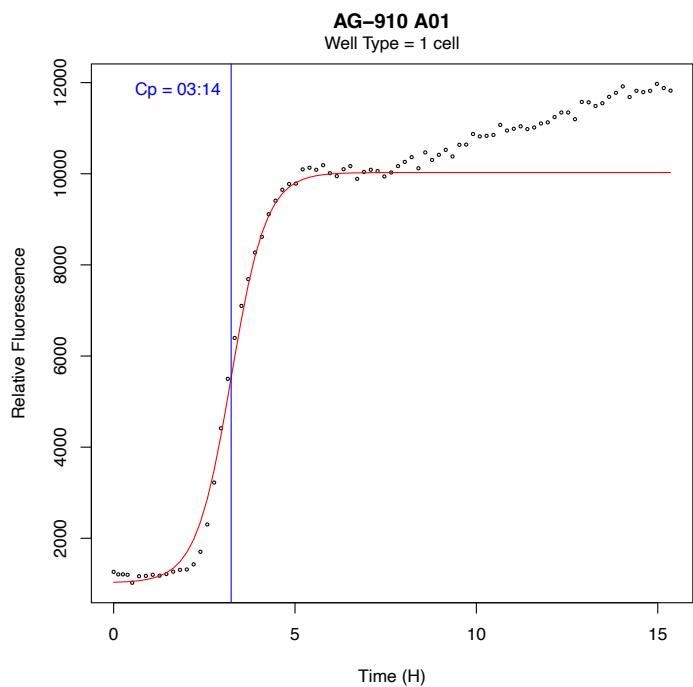
Validation of entire SCGC workflow



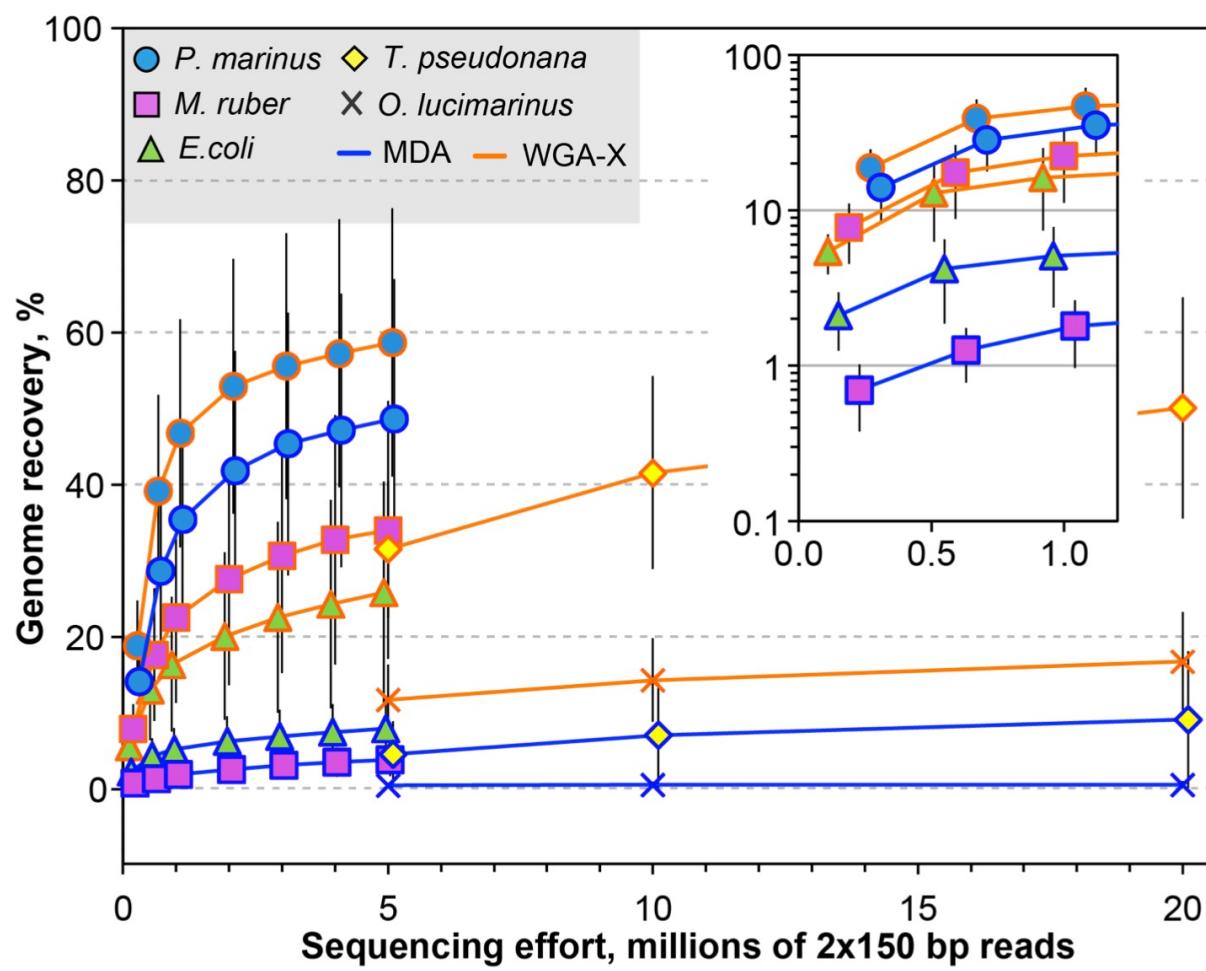
Validation of entire SCGC workflow



From: Stepanauskas et al. 2017

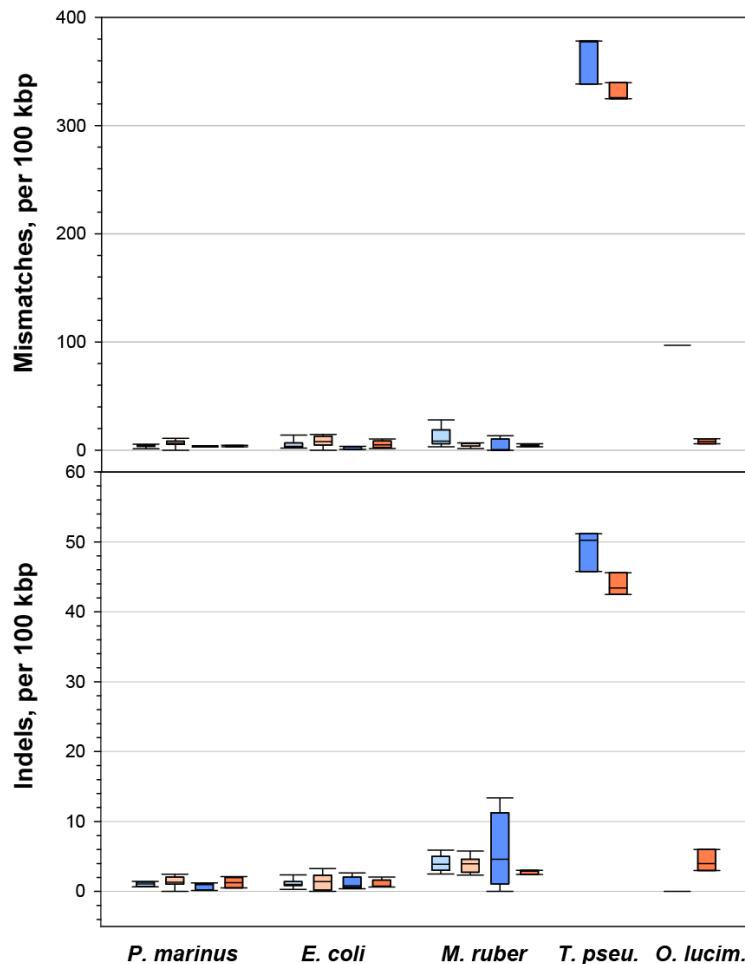
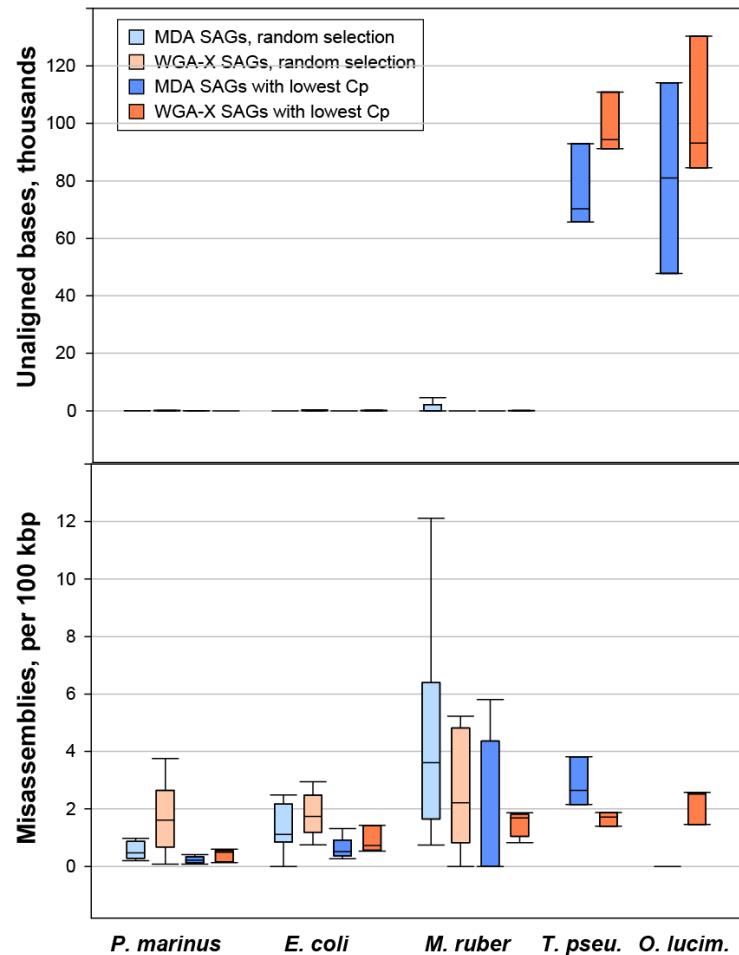


Relationship between sequencing depth and genome recovery



From: Stepanauskas et al. 2017

Validation of entire SCGC workflow

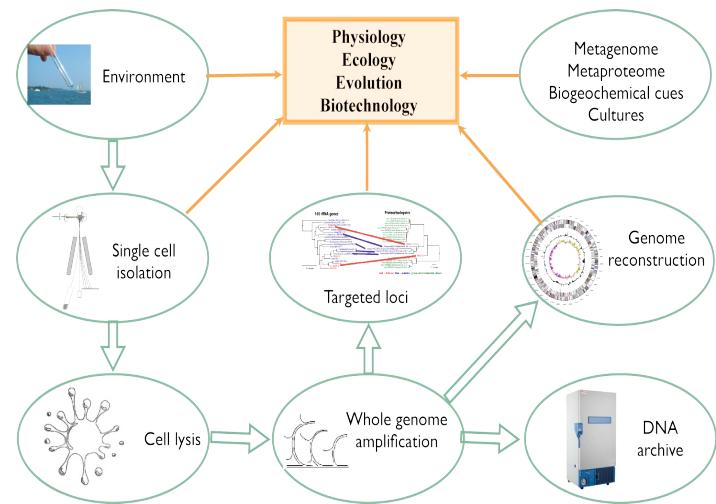


From: Stepanauskas et al. 2017

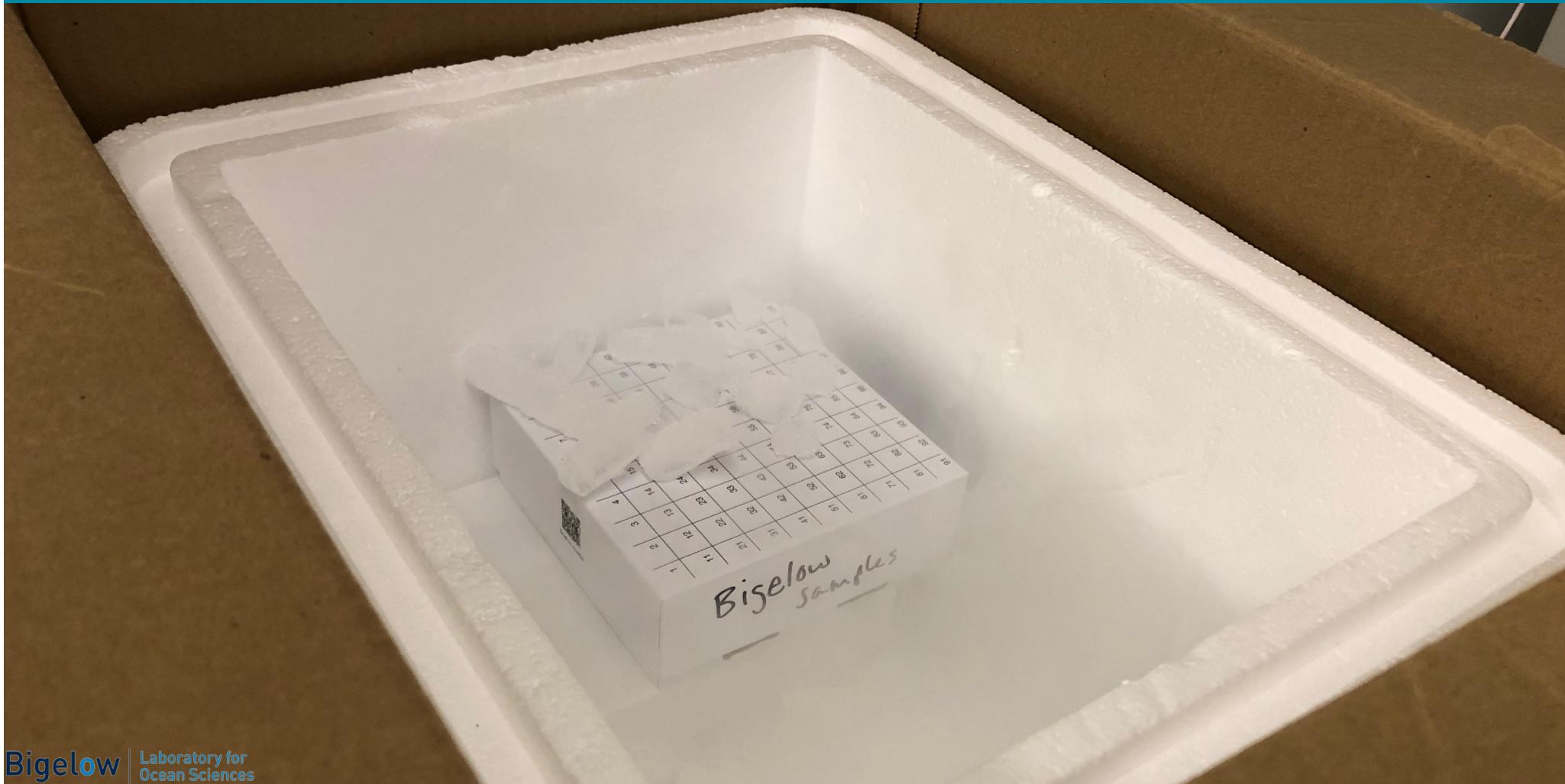
QA and QC at SCGC

SCGC strategy:

1. Validate the entire process
2. Minimize physical DNA contamination and damage
3. Maximize accuracy of the computational workflow
4. Generate data outputs for manual QC



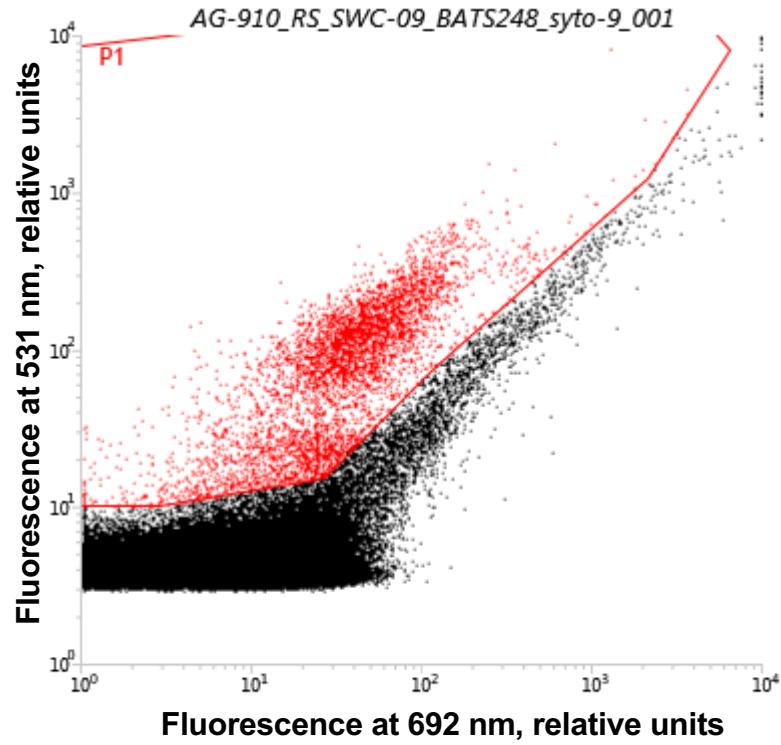
Pre-WGS QA: Field sample handling



Pre-WGS QA: Cleanroom environment



Pre-WGS QA: FACS gate



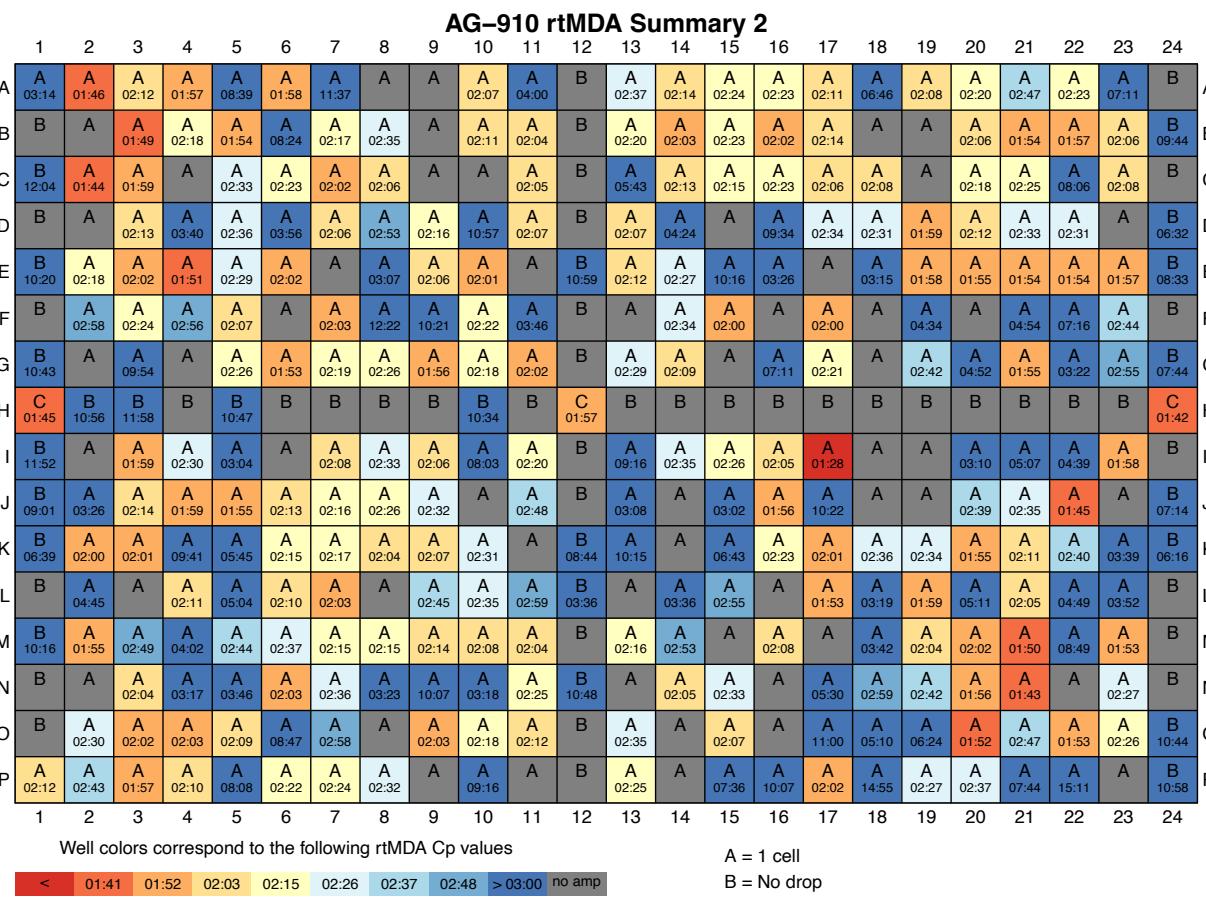
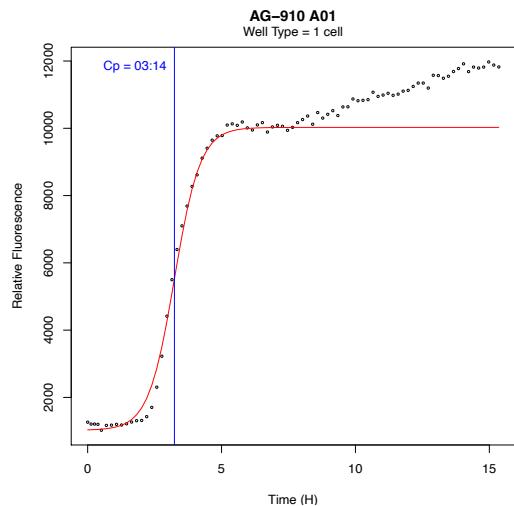
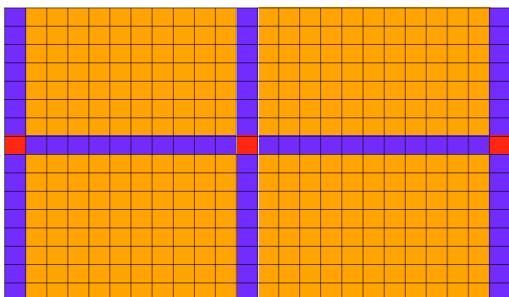
Pre-WGS QC: Whole genome amplification kinetics

384-well plate layout

Negative controls: empty wells

Positive controls: 10 cells per well

Wells containing single cells



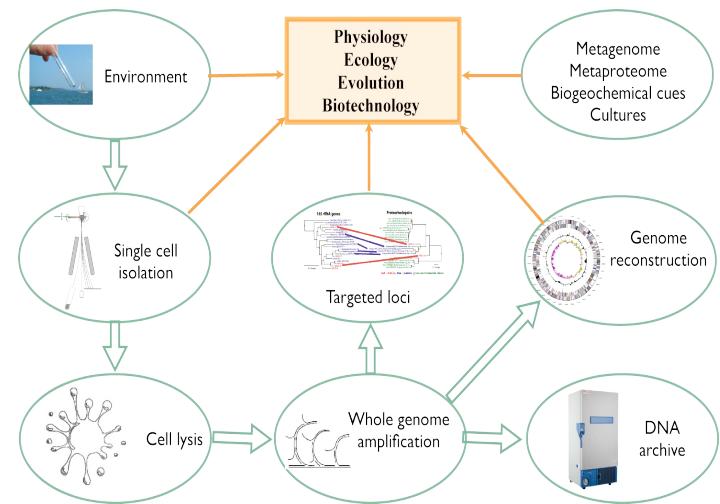
Contamination prevention and removal in SCGC's workflow

Measure	Type
Decontamination of instruments and reagents	Lab
Cleanroom for cell sorting and DNA amplification	Lab
Construction of a reagent contaminant database	Lab + Comp
Removal of reagent contaminants from reads and assemblies	Computation
Removal of contig ends, <2 kbp contigs, <20 kbp assemblies	Computation
Provision of tools for manual curation	Computation

QA and QC at SCGC

SCGC strategy:

1. Validate the entire process
2. Minimize physical DNA contamination and damage
3. Maximize accuracy of the computational workflow
4. Generate data outputs for manual QC



Tools for SCGC SAG assembly manual curation

Clue	Tool
Taxonomic assignment	Silva, (GTDBtk)
Extra copies of conserved single-copy genes	checkM
K-mer heterogeneity	Tetramer PCA
Heterogeneity of closest homologs	BLAST, (GUNC)
Annotation anomalies	Prokka, (DRAM)

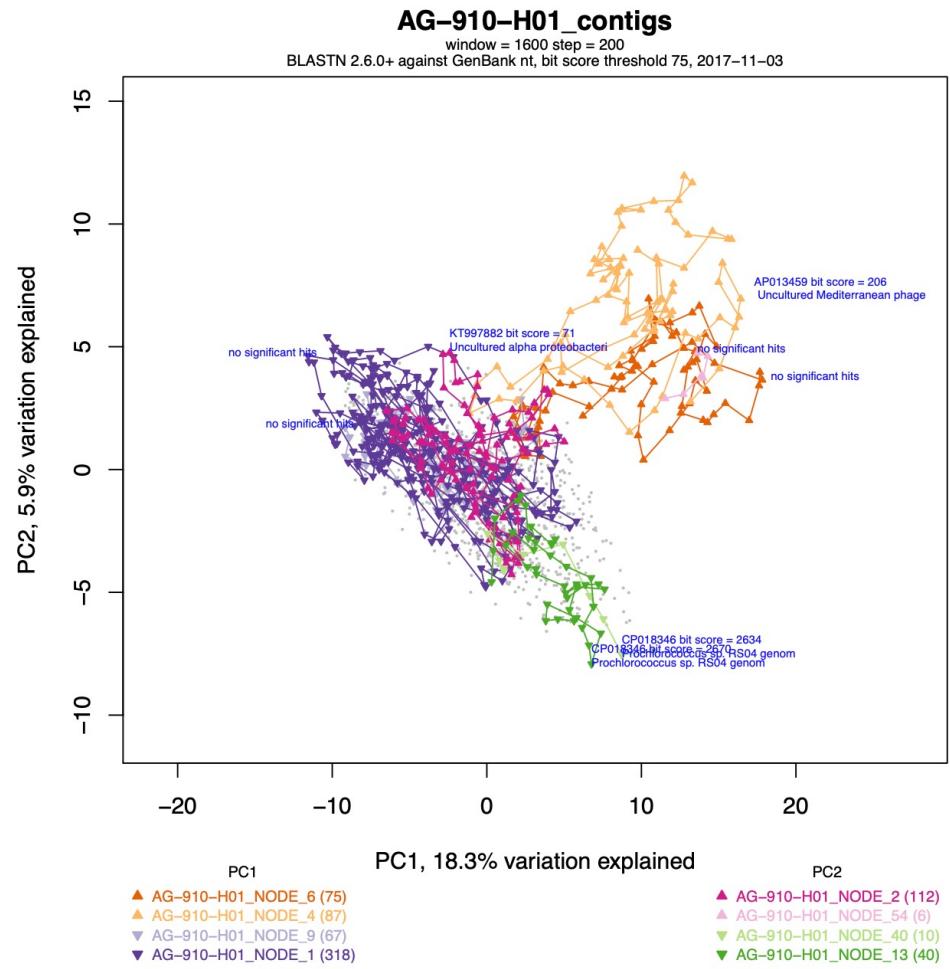
Tools for SAG curation: Taxonomic assignments

	Sample_ID	well_type	wga_cp	SSU_classification_1
Name	AG-910-A16	1 cell	2.38	k_Bacteria;p_Actinobacteria;c_Acidimicrobia;o_Acidimicrobiales;f_OM1_clade;g_Candidatus_Actinomarina;s_?
	AG-910-L02	1 cell	4.75	k_Bacteria;p_Actinobacteria;c_Acidimicrobia;o_Acidimicrobiales;f_OM1_clade;g_Candidatus_Actinomarina;s_?
	AG-910-M10	1 cell	2.15	k_Bacteria;p_Actinobacteria;c_Acidimicrobia;o_Acidimicrobiales;f_OM1_clade;g_Candidatus_Actinomarina;s_?
	AG-910-D03	1 cell	2.23	k_Bacteria;p_Actinobacteria;c_Acidimicrobia;o_Acidimicrobiales;f_OM1_clade;g_Candidatus_Actinomarina;s_?
	AG-910-I07	1 cell	2.14	k_Bacteria;p_Bacteroidetes;c_Flavobacteriia;o_Flavobacteriales;f_Flavobacteriaceae;g_NS4_marine_group;s_?
	AG-910-M02	1 cell	1.92	k_Bacteria;p_Bacteroidetes;c_Flavobacteriia;o_Flavobacteriales;f_Flavobacteriaceae;g_NS4_marine_group;s_?
	AG-910-P04	1 cell	2.18	k_Bacteria;p_Bacteroidetes;c_Flavobacteriia;o_Flavobacteriales;f_Flavobacteriaceae;g_NS5_marine_group;s_?
	AG-910-E10	1 cell	2.02	k_Bacteria;p_Cyanobacteria;c_Oxyphotobacteria;o_Synechococcales;f_Synechococcaceae;g_Prochlorococcus;s_?
	AG-910-J04	1 cell	1.99	k_Bacteria;p_Cyanobacteria;c_Oxyphotobacteria;o_Synechococcales;f_Synechococcaceae;g_Prochlorococcus;s_?
	AG-910-B21	1 cell	1.9	k_Bacteria;p_Proteobacteria;c_Alphaproteobacteria;o_Rhodobacterales;f_Rhodobacteraceae;g_?;s_?
	AG-910-B07	1 cell	2.3	k_Bacteria;p_Proteobacteria;c_Alphaproteobacteria;o_Rhodospirillales;f_Rhodospirillaceae;g_AEGEAN-169_marine_group;s_?
	AG-910-G06	1 cell	1.89	k_Bacteria;p_Proteobacteria;c_Alphaproteobacteria;o_Rhodospirillales;f_Rhodospirillaceae;g_AEGEAN-169_marine_group;s_?
	AG-910-G09	1 cell	1.94	k_Bacteria;p_Proteobacteria;c_Alphaproteobacteria;o_Rhodospirillales;f_Rhodospirillaceae;g_AEGEAN-169_marine_group;s_?
	AG-910-J16	1 cell	1.94	k_Bacteria;p_Proteobacteria;c_Alphaproteobacteria;o_Rhodospirillales;f_Rhodospirillaceae;g_AEGEAN-169_marine_group;s_?
	AG-910-K03	1 cell	2.03	k_Bacteria;p_Proteobacteria;c_Alphaproteobacteria;o_Rhodospirillales;f_Rhodospirillaceae;g_AEGEAN-169_marine_group;s_?
	AG-910-K06	1 cell	2.26	k_Bacteria;p_Proteobacteria;c_Alphaproteobacteria;o_Rhodospirillales;f_Rhodospirillaceae;g_AEGEAN-169_marine_group;s_?
	AG-910-N14	1 cell	2.09	k_Bacteria;p_Proteobacteria;c_Alphaproteobacteria;o_Rhodospirillales;f_Rhodospirillaceae;g_AEGEAN-169_marine_group;s_?
	AG-910-F04	1 cell	2.94	k_Bacteria;p_Proteobacteria;c_Alphaproteobacteria;o_Rhodospirillales;f_Rhodospirillaceae;g_AEGEAN-169_marine_group;s_?
	AG-910-E22	1 cell	1.91	k_Bacteria;p_Proteobacteria;c_Alphaproteobacteria;o_Rickettsiales;f_SAR116_clade;g_?;s_?
	AG-910-D22	1 cell	2.52	k_Bacteria;p_Proteobacteria;c_Alphaproteobacteria;o_Rickettsiales;f_SAR116_clade;g_?;s_?
	AG-910-M21	1 cell	1.85	k_Bacteria;p_Proteobacteria;c_Alphaproteobacteria;o_SAR11_clade;f_?;g_?;s_?
	AG-910-C06	1 cell	2.38	k_Bacteria;p_Proteobacteria;c_Alphaproteobacteria;o_SAR11_clade;f_?;g_?;s_?
	AG-910-C16	1 cell	2.38	k_Bacteria;p_Proteobacteria;c_Alphaproteobacteria;o_SAR11_clade;f_?;g_?;s_?
	AG-910-D17	1 cell	2.57	k_Bacteria;p_Proteobacteria;c_Alphaproteobacteria;o_SAR11_clade;f_?;g_?;s_?
	AG-910-E13	1 cell	2.21	k_Bacteria;p_Proteobacteria;c_Alphaproteobacteria;o_SAR11_clade;f_?;g_?;s_?
	AG-910-M20	1 cell	2.05	k_Bacteria;p_Proteobacteria;c_Alphaproteobacteria;o_SAR11_clade;f_?;g_?;s_?
	AG-910-N08	1 cell	3.39	k_Bacteria;p_Proteobacteria;c_Alphaproteobacteria;o_SAR11_clade;f_?;g_?;s_?
	AG-910-N20	1 cell	1.94	k_Bacteria;p_Proteobacteria;c_Alphaproteobacteria;o_SAR11_clade;f_?;g_?;s_?
	AG-910-E06	1 cell	2.04	k_Bacteria;p_Proteobacteria;c_Alphaproteobacteria;o_SAR11_clade;f_Surface_1;g_?;s_?

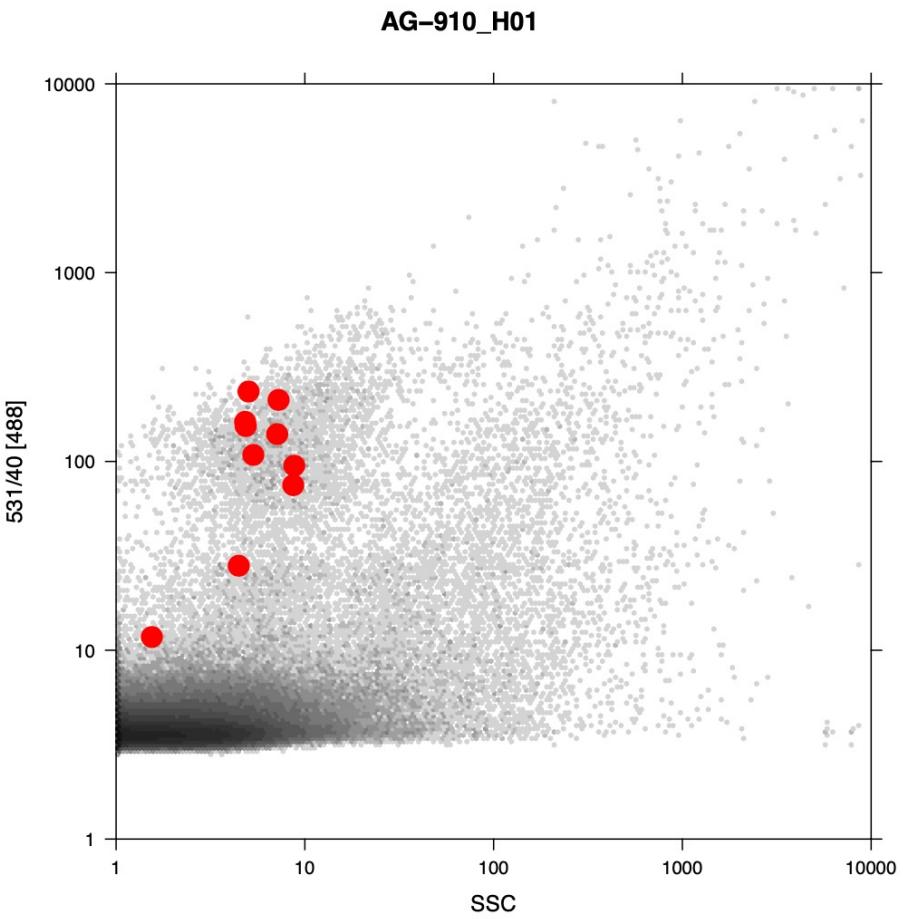
Tools for SAG curation: checkM

Sample_ID	well_type	wga_cp	final_assembly_length	max_contig_length	gc_content	checkM_estimated_completeness	number_multi_copy
AG-910-H24	10 cells	1.72	1,703,043	115,921	32	93	28
AG-910-N04	1 cell	3.29	214,674	27,694	29	22	2
AG-910-H12	10 cells	1.96	1,247,299	82,663	33	66	2
AG-910-K03	1 cell	2.03	1,067,614	103,436	30	56	1
AG-910-C18	1 cell	2.15	979,531	268,683	30	80	1
AG-910-I05	1 cell	3.08	83,974	27,163	32	4	1
AG-910-A16	1 cell	2.38	649,299	87,176	33	45	0
AG-910-L02	1 cell	4.75	134,319	41,293	32	14	0
AG-910-M10	1 cell	2.15	757,654	179,180	33	54	0
AG-910-D03	1 cell	2.23	127,585	56,880	29	11	0
AG-910-I07	1 cell	2.14	956,113	137,377	31	55	0
AG-910-M02	1 cell	1.92	1,128,449	191,928	31	73	0
AG-910-P04	1 cell	2.18	1,088,902	84,609	31	49	0
AG-910-E10	1 cell	2.02	1,247,571	148,843	31	84	0
AG-910-J04	1 cell	1.99	1,196,796	116,970	31	68	0
AG-910-B21	1 cell	1.9	1,930,694	375,705	31	75	0
AG-910-B07	1 cell	2.3	751,372	109,892	31	61	0
AG-910-G06	1 cell	1.89	1,155,577	330,630	31	70	0

Tetramer PCA: Multiple cells



FACS: Particle light scatter and green fluorescence



Assembly QC: BLASTN on contigs

Query Seq-id	Sub	Percentage of identical matches	Alignment length	All Subject Title(s)	Screenshot
AG-910-H01_NODE_4	gi 5	92	1182 82	Uncultured virus clone ctg_DTF_polA_1102 putative DNA polymerase A gene, partial cds	
AG-910-H01_NODE_4	gi 5	74	574 25	Uncultured Mediterranean phage uvMED DNA, complete genome, group G17, isolate: uvMED-CGR-U-MedDCM-OCT-S39-C21	
AG-910-H01_NODE_4	gi 1	75	337 97	Uncultured virus clone vSAG-41-A4-2 genomic sequence	
AG-910-H01_NODE_4	gi 1	75	334 97	Uncultured virus clone vSAG-41-A4-2 genomic sequence	
AG-910-H01_NODE_4	gi 5	81	101 64	Uncultured Mediterranean phage uvMED DNA, complete genome, group G19, isolate: uvMED-CGR-U-MedDCM-OCT-S24-C60	
AG-910-H01_NODE_4	gi 1	91	43 05	Pichia membranifaciens NRRL Y-2026 hypothetical protein partial mRNA	
AG-910-H01_NODE_4	gi 2	94	691 91	Uncultured marine organism clone JORCPF001_22_POF_SPF_gb1 genomic sequence	

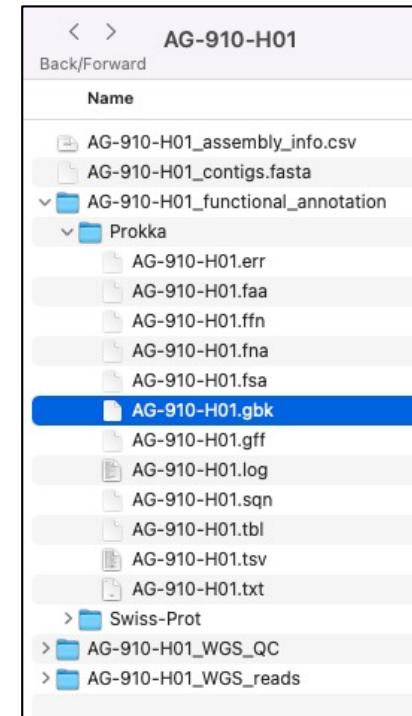
< > AG-910-H01
Back/Forward

Name

- AG-910-H01_assembly_info.csv
- AG-910-H01_contigs.fasta
- AG-910-H01_functional_annotation
- AG-910-H01_WGS_QC
- AG-910-H01_all_contigs_blastn
 - AG-910-H01_AG-665_masked.fasta_all_contigs ASN1
 - AG-910-H01_AG-665_masked.fasta_all_contigs.html
 - AG-910-H01_nt_all_contigs ASN1
 - AG-910-H01_nt_all_contigs.html
 - AG-910-H01_nt_all_contigs.tsv
 - AG-910-H01_Simons_negatives_170807.fasta_all_contigs ASN1
 - AG-910-H01_Simons_negatives_170807.fasta_all_contigs.html
 - AG-910-H01_Simons_negatives_170807.fasta_all_contigs.tsv
 - AG-910-H01_all_contigs.fasta
 - AG-910-H01_all_contigs.fasta.count
 - AG-910-H01_contaminated_contigs.fasta
- AG-910-H01_fastqc
- AG-910-H01_final_contigs_checkm
- AG-910-H01_final_contigs_tetramer_pca
 - AG-910-H01_length_failing_contigs.fasta
 - AG-910-H01_length_passing_contigs.fasta
 - AG-910-H01_length_passing_contigs.fasta.count
 - AG-910-H01_raw_trimmed_contigs.fasta
- SSU_rRNA_recovery
- AG-910-H01_WGS_reads

Assembly QC: Prokka annotation

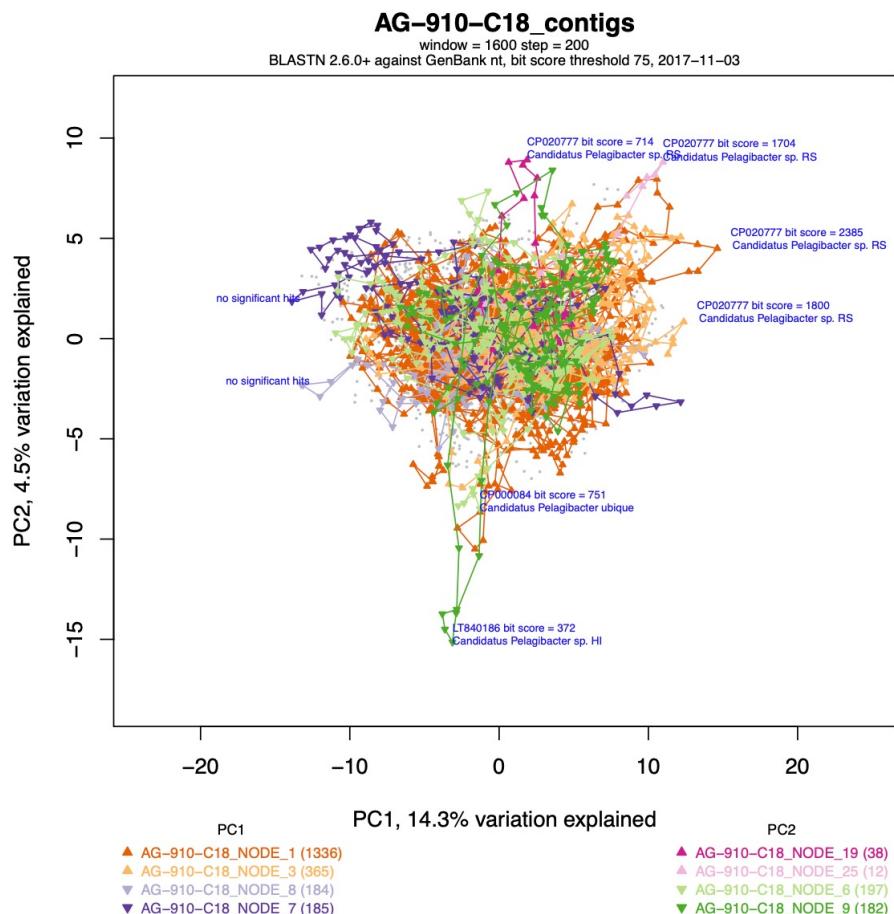
```
//  
LOCUS    AG-910-H01_NODE_4      18923 bp    DNA     linear      03-NOV-2017  
DEFINITION Genus species strain strain.  
ACCESSION  
VERSION  
KEYWORDS .  
SOURCE   Genus species  
ORGANISM Genus species  
Unclassified.  
COMMENT  Annotated using prokka 1.12 from  
https://github.com/tseemann/prokka.  
FEATURES  Location/Qualifiers  
source    1..18923  
          /organism="Genus species"  
          /mol_type="genomic DNA"  
          /strain="strain"  
          CDS       21..131  
          /locus_tag="AG-910-H01_00107"  
          /inference="ab initio prediction:Prodigal:2.6"  
          /codon_start=1  
          /transl_table=11  
          /product="hypothetical protein"  
          /translation="MSNKCKYCNGNCPNDSYMCMDGYAGDIDNLYENNEE"  
          CDS       134..571  
          /locus_tag="AG-910-H01_00108"  
          /inference="ab initio prediction:Prodigal:2.6"  
          /codon_start=1  
          /transl_table=11  
          /product="hypothetical protein"  
          /translation="MNMDRLLESVKKHEGYRNKVLDTLGKRTVGVGHLCVEDFWEDN  
KEYEEKFLMTILEHDLQTAIKGAKELMEDHGACADIDEQAEELIEMVFQLGKNGVSFK  
KNMWKALAEKNYIGASYEMLSRWAKQTPNRAKSMAKTMKEIT"  
          CDS       568..753  
          /locus_tag="AG-910-H01_00109"  
          /inference="ab initio prediction:Prodigal:2.6"  
          /codon_start=1  
          /transl_table=11  
          /product="hypothetical protein"  
          /translation="MIYKGQKLSDNLTKSQRDMFKYMIDRGDSIEEIVKIYTLDDKVN  
QALREPDVEEKIEALK"
```



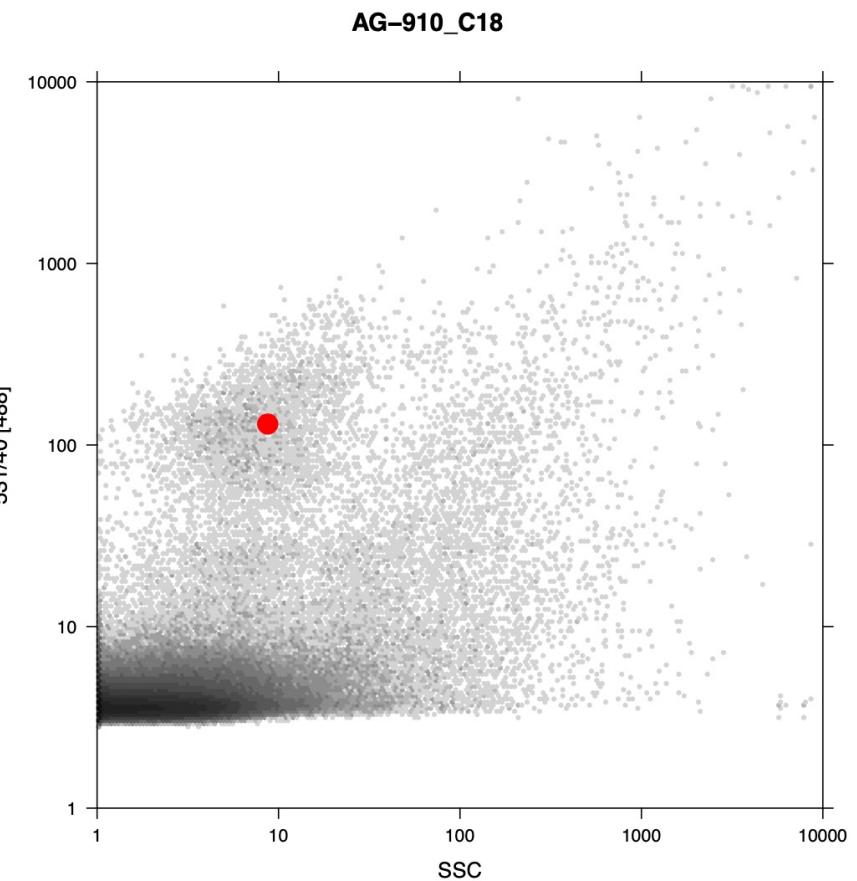
Tools for SAG curation: checkM

Sample_ID	well_type	wga_cp	final_assembly_length	max_contig_length	gc_content	checkM_estimated_completeness	number_multi_copy
AG-910-H24	10 cells	1.72	1,703,043	115,921	32	93	28
AG-910-N04	1 cell	3.29	214,674	27,694	29	22	2
AG-910-H12	10 cells	1.96	1,247,299	82,663	33	66	2
AG-910-K03	1 cell	2.03	1,067,614	103,436	30	56	1
AG-910-C18	1 cell	2.15	979,531	268,683	30	80	1
AG-910-I05	1 cell	3.08	83,974	27,163	32	4	1
AG-910-A16	1 cell	2.38	649,299	87,176	33	45	0
AG-910-L02	1 cell	4.75	134,319	41,293	32	14	0
AG-910-M10	1 cell	2.15	757,654	179,180	33	54	0
AG-910-D03	1 cell	2.23	127,585	56,880	29	11	0
AG-910-I07	1 cell	2.14	956,113	137,377	31	55	0
AG-910-M02	1 cell	1.92	1,128,449	191,928	31	73	0
AG-910-P04	1 cell	2.18	1,088,902	84,609	31	49	0
AG-910-E10	1 cell	2.02	1,247,571	148,843	31	84	0
AG-910-J04	1 cell	1.99	1,196,796	116,970	31	68	0
AG-910-B21	1 cell	1.9	1,930,694	375,705	31	75	0
AG-910-B07	1 cell	2.3	751,372	109,892	31	61	0
AG-910-G06	1 cell	1.89	1,155,577	330,630	31	70	0

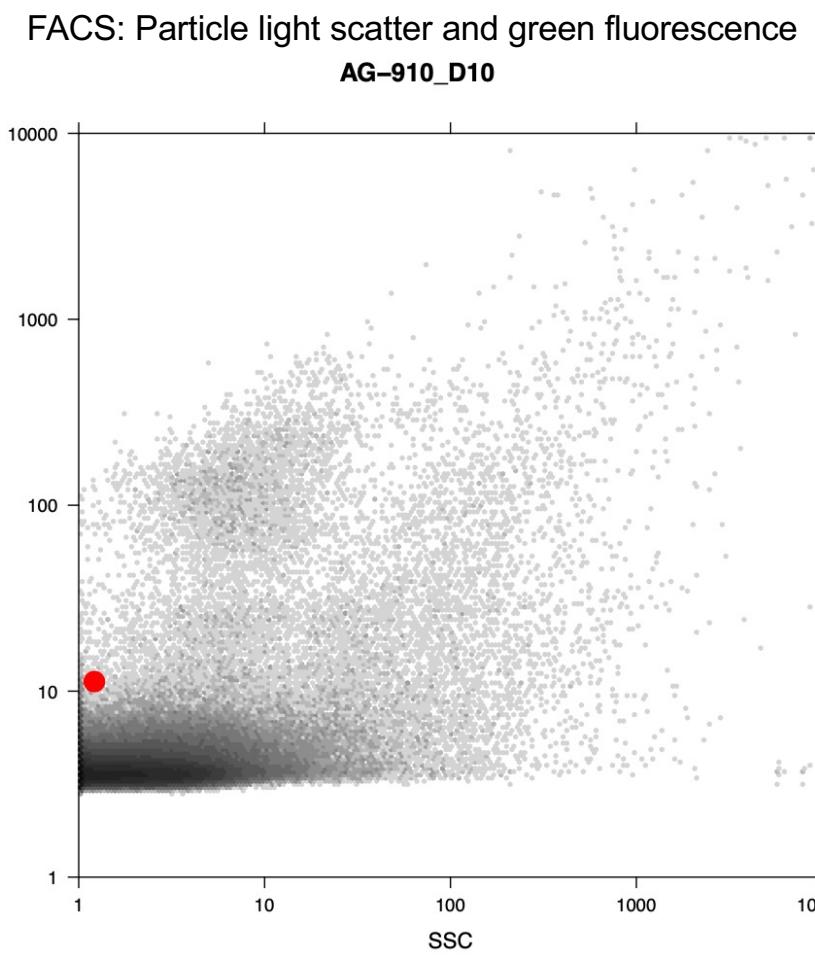
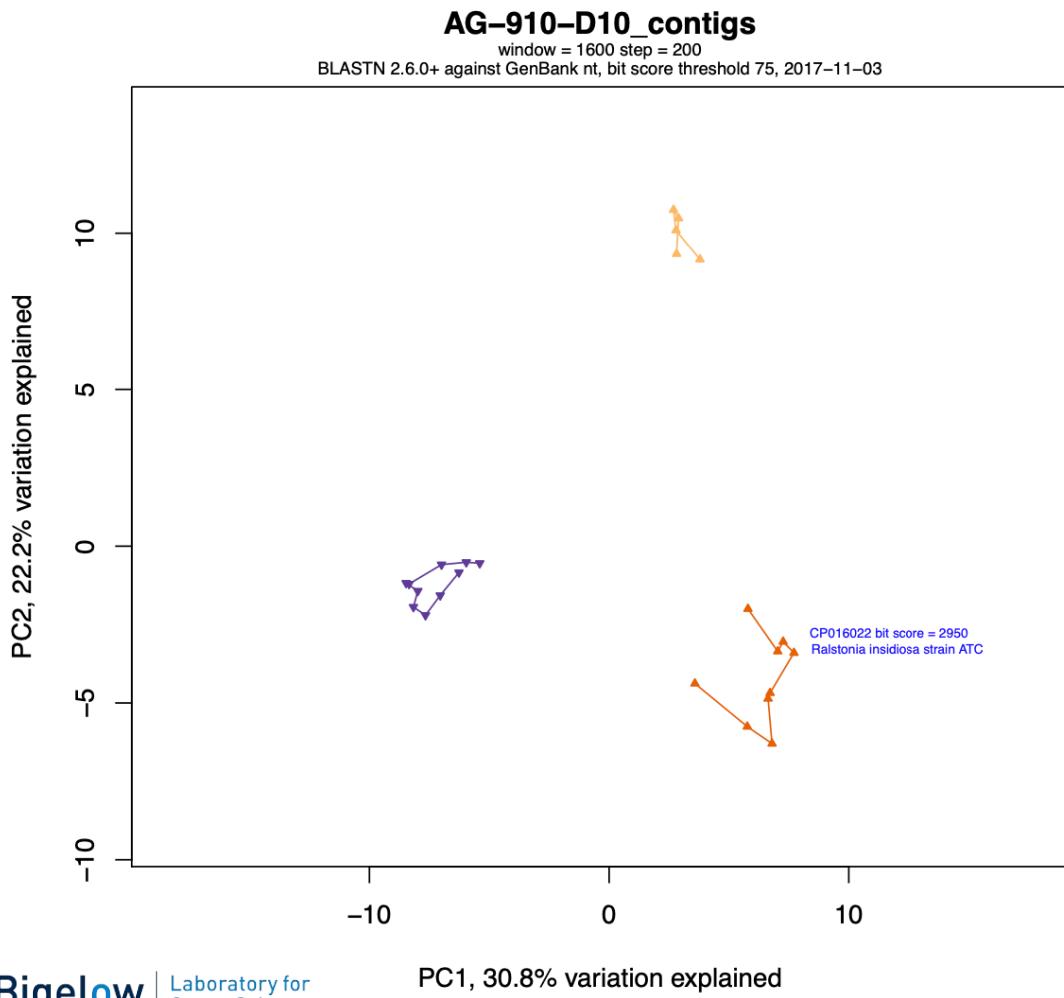
Tetramer PCA: Clean SAG



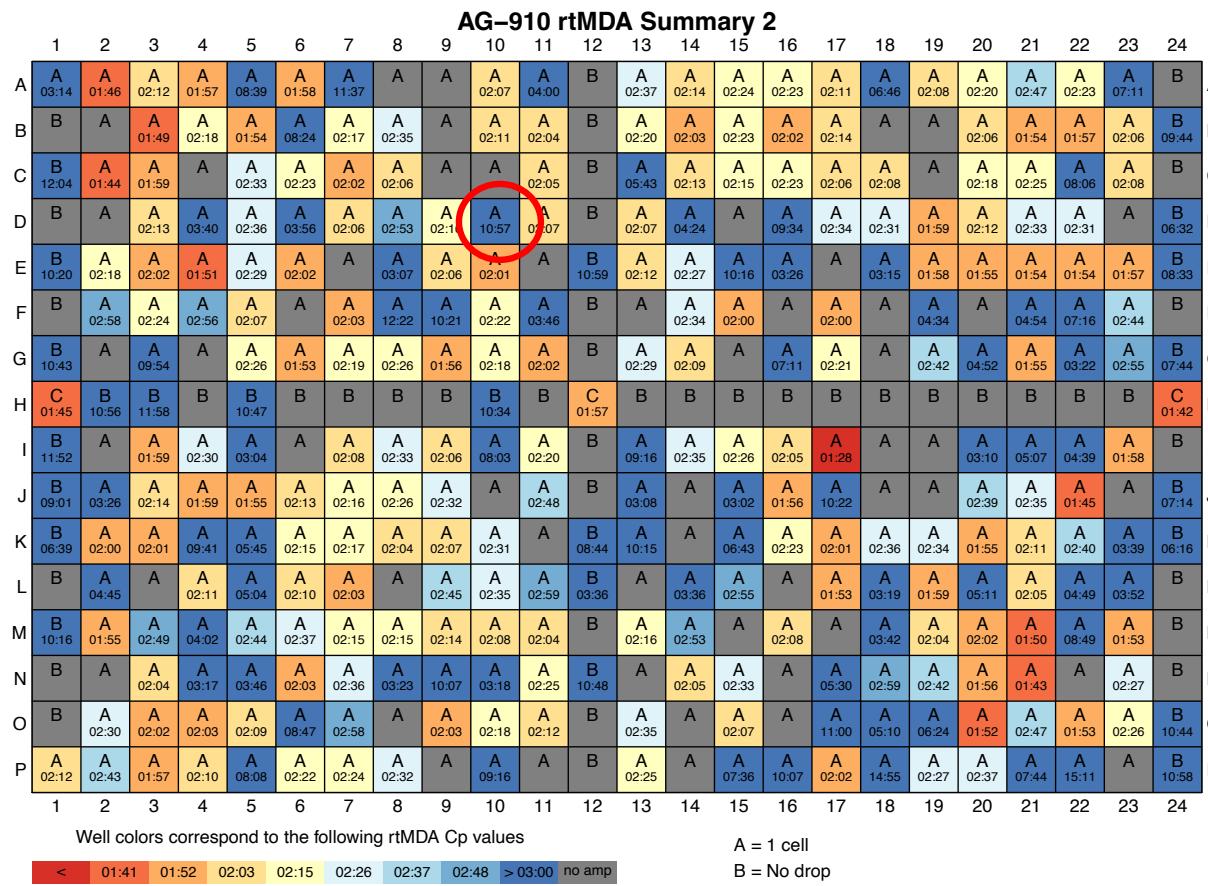
FACS: Particle light scatter and green fluorescence



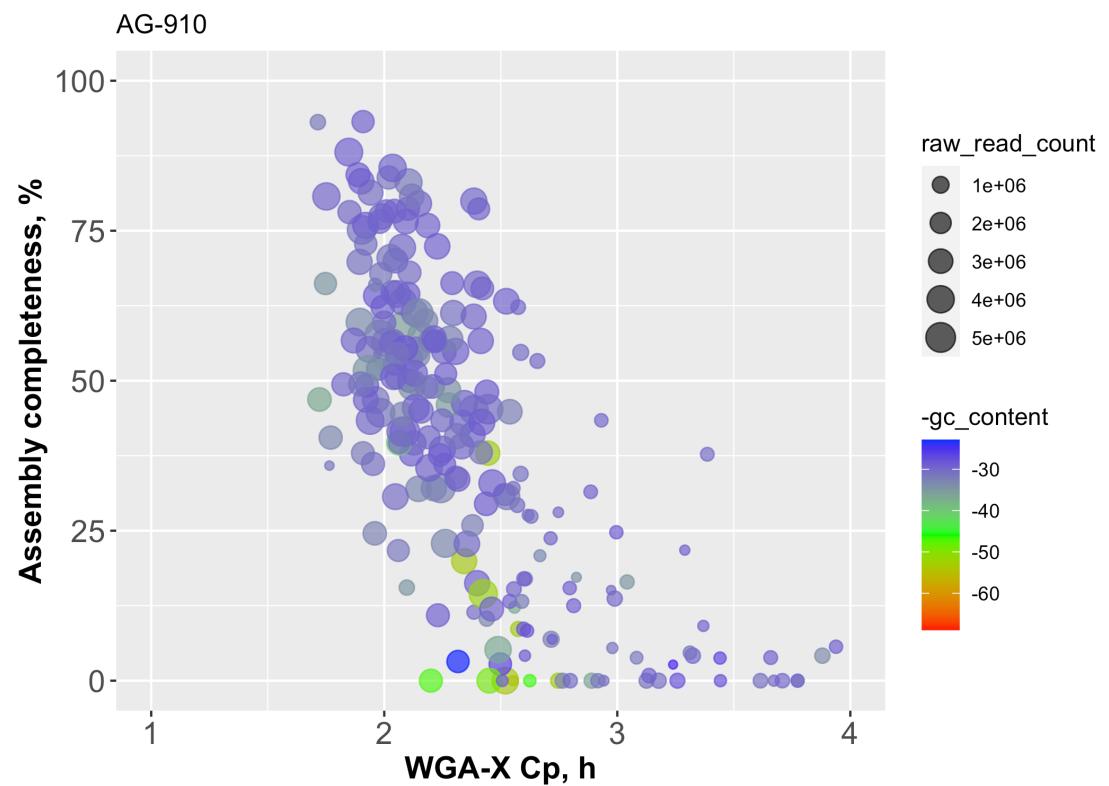
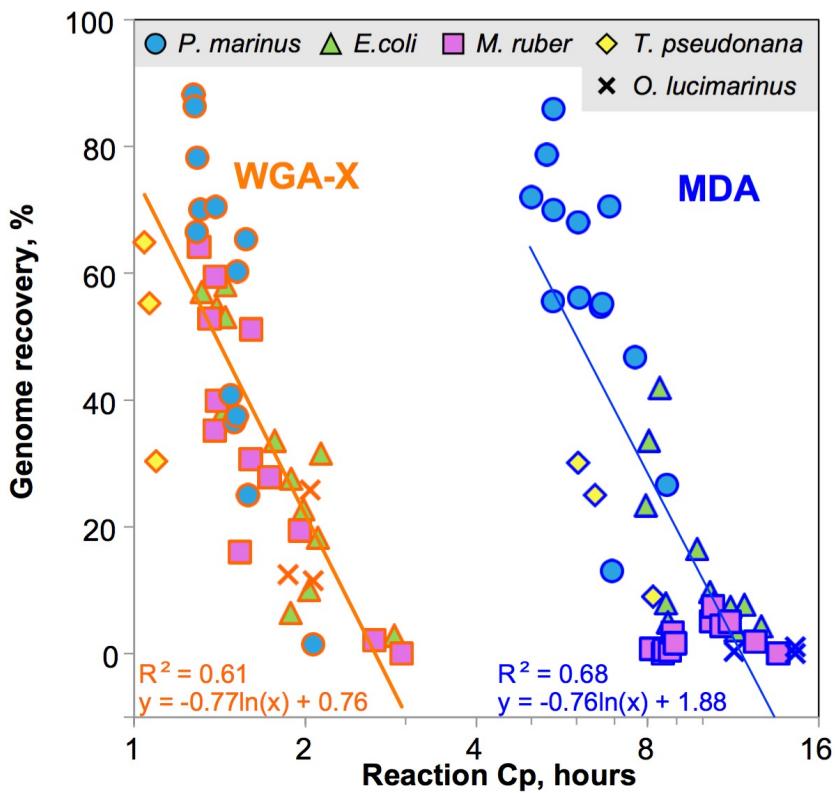
Tetramer PCA: Reagent contaminants



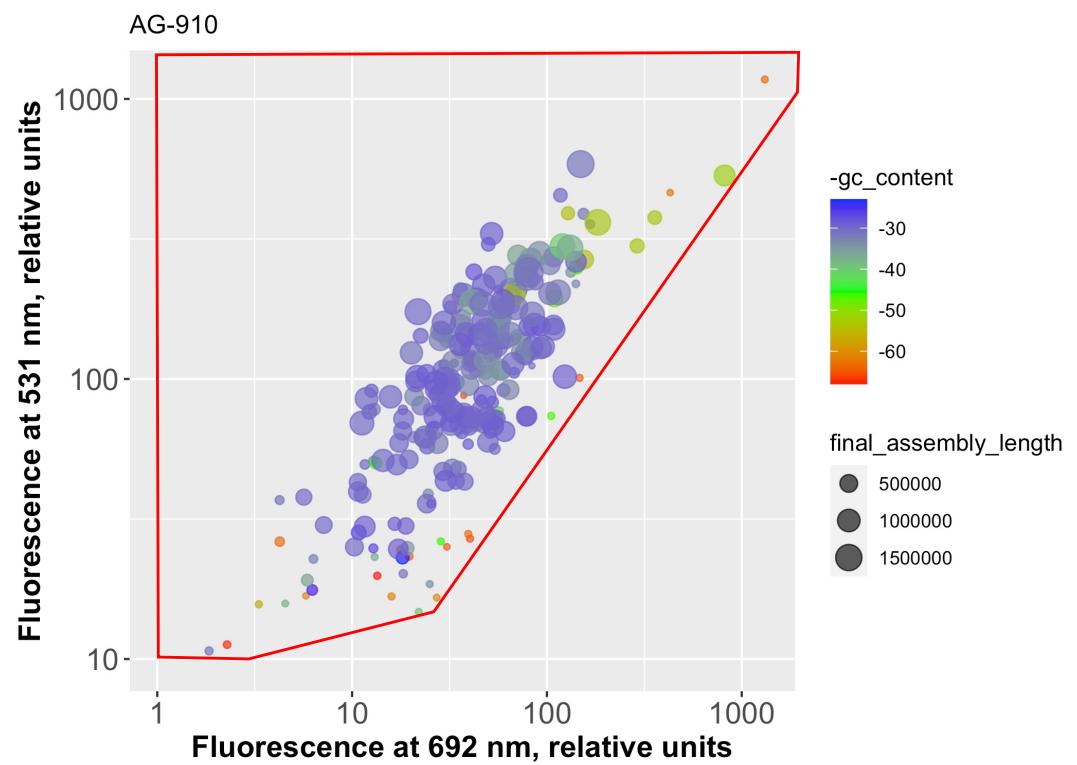
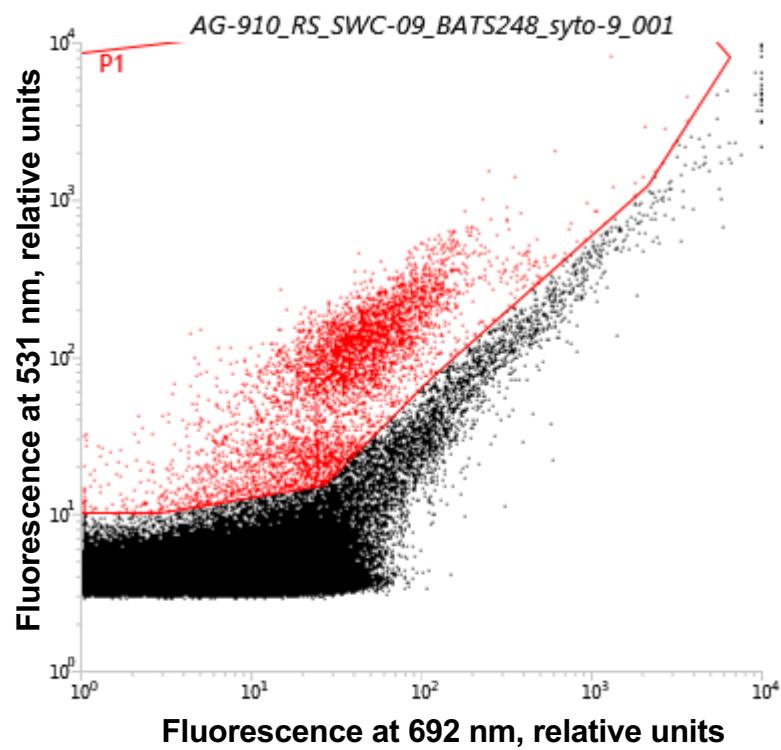
Pre-WGS QC: Whole genome amplification kinetics



Multi-parameter visualization of AG-910 SAGs

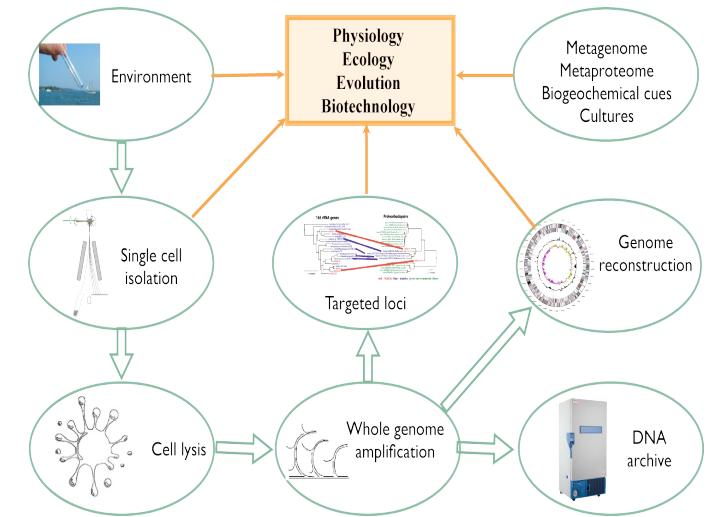


Multi-parameter visualization of AG-910 SAGs

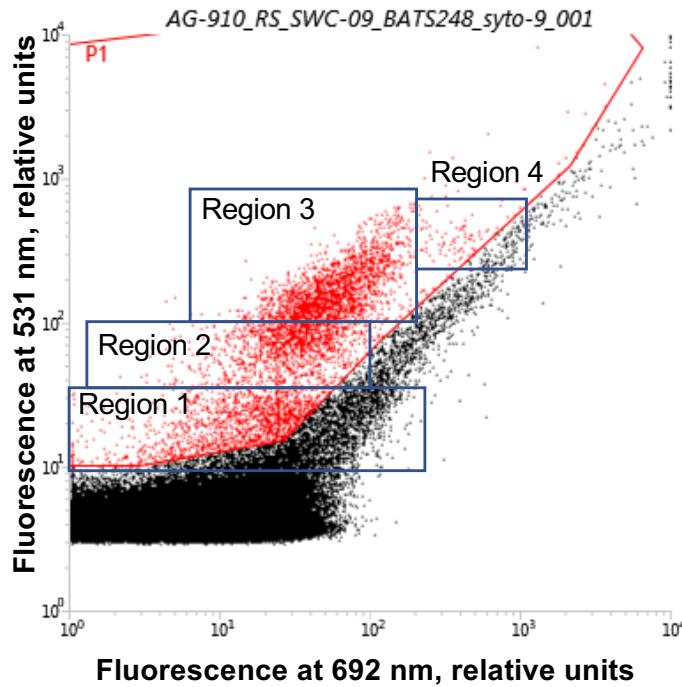


Take-home message

1. Use genome QC metrics and thresholds that match your research goals.
2. Understand how the data were produced, field -> lab -> computation.
3. Be critical: no tool is perfect; run reality checks and use all available information.



SCGC data exploration exercise



What particles were in FACS regions 1-4?