

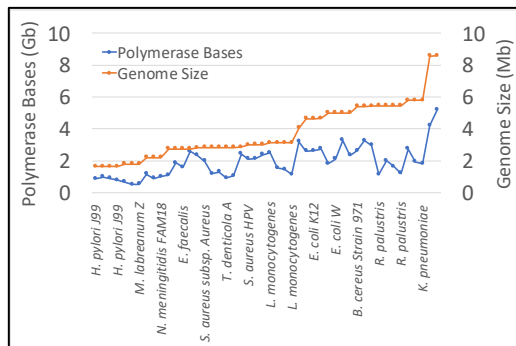
## Background and Motivation

Complete, high-quality microbial genomes are very valuable across a broad array of fields, from environmental studies, to human microbiome health, food pathogen surveillance, etc. Long-read sequencing enables accurate resolution of complex microbial genomes and is becoming the new standard. Here we report our novel Microbial Assembly pipeline to facilitate rapid, large-scale analysis of microbial genomes. We sequenced a 48-plex library with one SMRT Cell 8M on the Sequel II System, demultiplexed, then analyzed the data with Microbial Assembly.

## Strains, Sequencing, and Demultiplexing

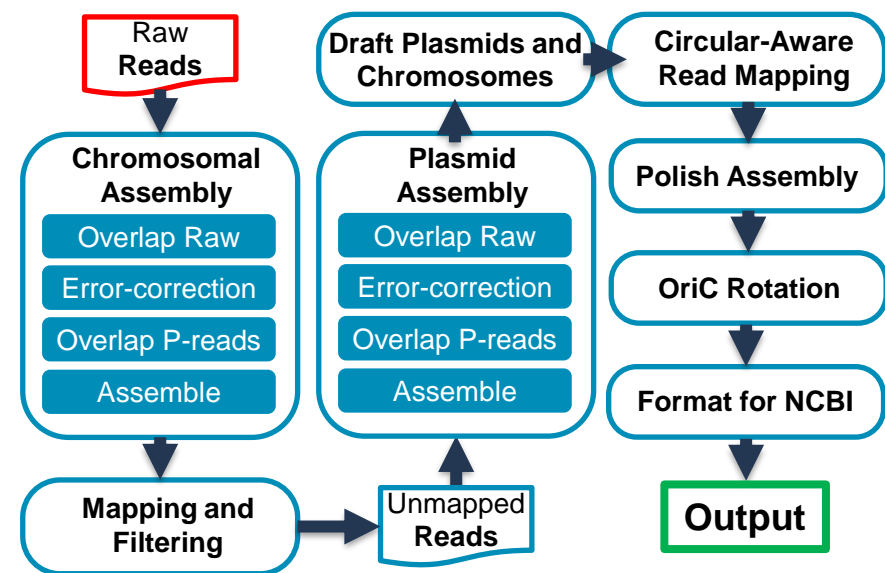
Strain	ATCC ID	Replicates	Genome	Plasmids
<i>Bacillus cereus</i> 971	14579	3	5.4 Mb	15 kb
<i>Bacillus subtilis</i> W23	in-house	1	4.0 Mb	n/a
<i>Burkholderia cepacia</i> *	25416	2	8.4 Mb	209 kb
<i>Enterococcus faecalis</i> OG1RF	47077D-5	4	2.7 Mb	n/a
<i>Escherichia coli</i> K12	in-house	3	4.6 Mb	n/a
<i>Escherichia coli</i> W	9637	4	4.9 Mb	5,102 kb
<i>Helicobacter pylori</i> J99	700824	4	1.7 Mb	n/a
<i>Klebsiella pneumoniae</i>	BAA-2146	3	5.4 Mb	2,85,118,141 kb
<i>Listeria monocytogenes</i>	19117	4	2.9 Mb	n/a
<i>Methanocorpusculum labreanum</i> Z	43576	3	1.8 Mb	n/a
<i>Neisseria meningitidis</i> FAM18	700532	3	2.2 Mb	n/a
<i>Rhodopseudomonas palustris</i>	in-house	4	5.5 Mb	n/a
<i>Staphylococcus aureus</i> HPV	BAA-44	3	2.9 Mb	24 kb
<i>Staphylococcus aureus</i> subsp. aureus	25923	3	2.8 Mb	27 kb
<i>Treponema denticola</i> A	35405	4	2.8 Mb	n/a

Value	Analysis Metric
104,515,141,915	Polymerase Read Bases
2,764,838	Polymerase Reads
37,802	Polymerase Read Length (mean)
87,502	Polymerase Read N50
9,544	Subread Length (mean)
12,192	Subread N50
10,735	Insert Length (mean)
14,812	Insert N50
28,944,687,104	Unique Molecular Yield

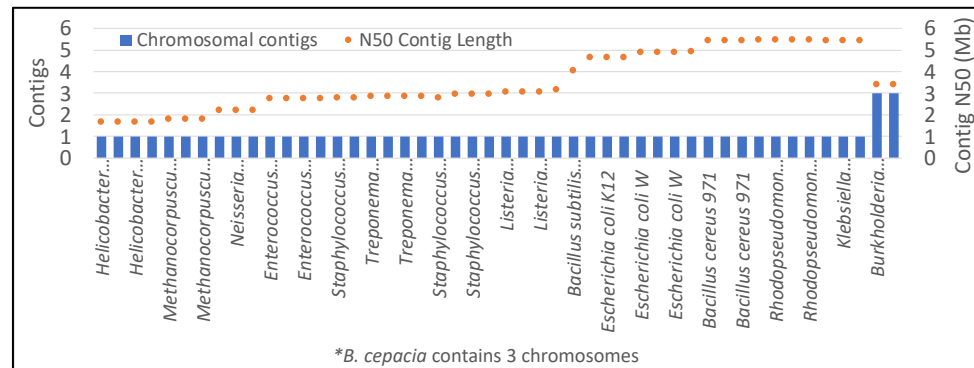


**Results of sequencing and demultiplexing barcodes.** The microbial pooling calculator helps balance coverage across species with different genome sizes, such that each has enough data to be processed with Microbial Assembly. Our lima demultiplex tool is available as a free bioconda package or in SMRT Link.

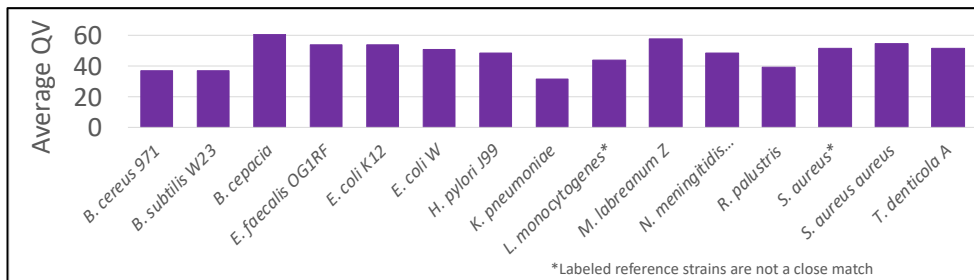
## Brand New Microbial Assembly Tool



## Performance of Chromosome Recovery



**Microbial Assembly results in recovery of chromosomes in complete contigs.** The chromosomes of each bacteria were consistently assembled into complete, circular contigs. This figure does not account for the plasmids assembled by the pipeline and included in the output.



**Analysis of Microbial Assembly results compared with references.** The average nucleotide identity (ANI) was calculated for the Microbial Assembly results and respective reference assemblies from the NCBI database. These data show that our results agree strongly with reference assemblies.

## Implementation of New Features

- Key differences between HGAP4 and Microbial Assembly
  - No repeat masking in Microbial Assembly workflow
  - DALIGNER tool stack is replaced with Raptor tool stack
  - Raptor – New graph-based mapping and alignment tool**
  - New method to detect and filter out chimeric reads**
- String graph and error-correction reused from FALCON

Assembly		HGAP4	Microbial Assembly
Overlap Raw		DALIGNER	Raptor
Error Correction		DAZZ_DB	RaptorDB
Overlap P-reads		fc_consensus	fc_consensus
Assemble		DALIGNER	Raptor
		DAZZ_DB	RaptorDB
		-	Chimera Detection
		Falcon SG	Falcon SG

## Performance of Plasmid Recovery

<i>E. coli</i> W plasmids	Replicate 1	Replicate 2	Replicate 3	Replicate 4
5 kb	X	✓	X	✓
102 kb	✓	✓	✓	✓

<i>K. pneumoniae</i> plasmids	Replicate 1	Replicate 2	Replicate 3
2 kb	X	X	X
85 kb	✓	✓	✓
118 kb	✓	✓	✓
141 kb	✓	✓	✓

**Microbial Assembly results in successful recovery for most plasmids.** The notable exception is a 2 kb plasmid in *K. pneumoniae* and a 5 kb plasmid in *E. coli* W. However, this is likely due to removal of shorter DNA fragments during the library preparation. By default, the pipeline assumes all contigs <300 kb are plasmids.

Bacteria with single plasmid	<i>B. cereus</i> 971 (15 kb)	<i>B. cepacia</i> (209 kb)	<i>S. aureus</i> HPV (24 kb)	<i>S. aureus</i> aureus (27 kb)	<i>L. monocytogenes</i> (no plasmid in NCBI)
Replicate 1	✓	✓ (204 kb)	✓ (32 kb)	✓	New Plasmid (94 kb)
Replicate 2	X	✓ (204 kb)	✓ (32 kb)	✓	New Plasmid (94 kb)
Replicate 3	X	N/A*	✓ (32 kb)	✓	New Plasmid (94 kb)
Replicate 4	N/A*	N/A*	N/A*	N/A*	X

\*N/A indicates lack of replicates for strain (see strain table for number of replicates per strain)

**Overall performance of Microbial Assembly is very good.** With the recovery of complete chromosomes and a majority of plasmids, our tool will help advance the study of microbes in a wide variety of fields. Further developments in the sample preparation protocol and data analysis software will continue to improve results.