

# Lawrence Berkeley National Laboratory

## Recent Work

### Title

Microbial Finishing at JGI

### Permalink

<https://escholarship.org/uc/item/4048v7mv>

### Authors

Lapidus, Alla  
Chain, Patrick  
Han, Cliff  
[et al.](#)

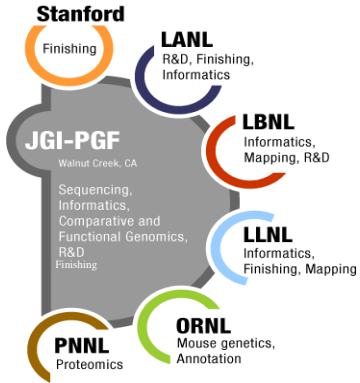
### Publication Date

2005-09-09

## Microbial Finishing at JGI

Alla Lapidus, Patrick Chain, Cliff Han, Eugene Goltsman, Michele Martinez, Stephanie Malfatty, Olga Chertkov, Stephan Trong, Tom Brettin, Roxanne Tapia, Alex Copeland, Paul Richardson.

\*Correspondence alapidus@lbl.gov



Project ID	Project Name	Lead	Status	Start Date	End Date	Completion %
000001	...	...	...	...	...	...
000002	...	...	...	...	...	...
000003	...	...	...	...	...	...
000004	...	...	...	...	...	...
000005	...	...	...	...	...	...
000006	...	...	...	...	...	...
000007	...	...	...	...	...	...
000008	...	...	...	...	...	...
000009	...	...	...	...	...	...
000010	...	...	...	...	...	...
000011	...	...	...	...	...	...
000012	...	...	...	...	...	...
000013	...	...	...	...	...	...
000014	...	...	...	...	...	...
000015	...	...	...	...	...	...
000016	...	...	...	...	...	...
000017	...	...	...	...	...	...
000018	...	...	...	...	...	...
000019	...	...	...	...	...	...
000020	...	...	...	...	...	...
000021	...	...	...	...	...	...
000022	...	...	...	...	...	...
000023	...	...	...	...	...	...
000024	...	...	...	...	...	...
000025	...	...	...	...	...	...
000026	...	...	...	...	...	...
000027	...	...	...	...	...	...
000028	...	...	...	...	...	...
000029	...	...	...	...	...	...
000030	...	...	...	...	...	...
000031	...	...	...	...	...	...
000032	...	...	...	...	...	...
000033	...	...	...	...	...	...
000034	...	...	...	...	...	...
000035	...	...	...	...	...	...
000036	...	...	...	...	...	...
000037	...	...	...	...	...	...
000038	...	...	...	...	...	...
000039	...	...	...	...	...	...
000040	...	...	...	...	...	...
000041	...	...	...	...	...	...
000042	...	...	...	...	...	...
000043	...	...	...	...	...	...
000044	...	...	...	...	...	...
000045	...	...	...	...	...	...
000046	...	...	...	...	...	...
000047	...	...	...	...	...	...
000048	...	...	...	...	...	...
000049	...	...	...	...	...	...
000050	...	...	...	...	...	...

### Finishing Statistics

Category	Total	Description
Complete Finished	44	Project completed to finish standards
Complete Draft	4	Project completed to draft standards
Finishing	55	Draft phase complete, Finish phase underway
Production	27	Libraries completed and draft sequencing underway
Library	19	Libraries under construction for drafting
Awaiting DNA	40	Awaiting DNA for initiation of library construction
Pending	16	Statement of Work incomplete for project
<b>Grand Total</b>	<b>205</b>	

Taxa	Total
archaea	16
bacterium	150
environmental	18
eukaryon	7
symbiont	2
virus	11
<b>Grand Total</b>	<b>205</b>

### Microbial Finishing steps

- QD -> draft assembly
- Repeat resolution (automated and manual)
- Primer walk
- More repeat resolution
- Assembled genome -> closed
- Polishing (quality improvement) -> polished
- Assembly QA (Stanford) -> FINISHED

### Finishing Check List/Standards

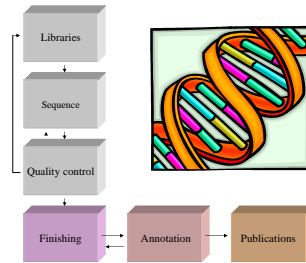
- All low quality areas (<Q30) reviewed and re sequenced.
- Final error rate < 0.2 per 10 Kb.
- No single clone coverage (minimum of 2X depth everywhere).
- Manually inspected and quantified single stranded regions.
- Checked all high quality discrepancies.
- Final sequence has no strings of xxxx anywhere
- All repeats verified (paired ends and PCR if necessary).
- Check ends of final contigs (chromosomes, plasmids)
- Final Assembly QC

The U.S. Department of Energy (DOE) established the Microbial Genome and GTL programs to determine the complete genome sequence of a number of microbes that may be useful to DOE in carrying out its missions (which include research of new energy sources, sequester excess atmospheric carbon affecting global climate, and to clean up contaminated environments). Another program was established to study hard-to-culture individual microbes and microbial communities. They are very difficult to study but play critical roles in the Earth's ecology. JGI is a leader in performing sequences to support these programs. Another important specialization of the JGI is the recently started Community Sequencing Program (<http://www.jgi.doe.gov/CSP/index.html>).

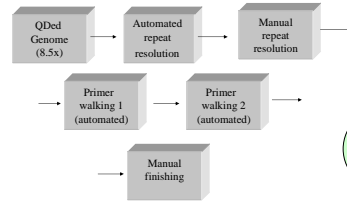
All of the above projects include sequence and detailed analysis of the genomes of the different representatives of the microbial world. A completely sequenced, high quality genome is a perfect starting point for the genome annotation (<http://img.jgi.doe.gov/v1.1/main.cgi>), microarrays, knockout experiments etc. Despite the fact that drafted genome contains a sufficiently large amount of information to be usable, a completed and polished one is overall a better product especially if it will be used to analyze previously unknown and difficult-to-cultivate microbes; for the comparative analysis of the clinical isolates or for the creation of microbial strains overproducing different proteins and amino acids. Knowledge of the completely finished genome will allow scientists to modify specific regions of the genome and therefore to affect the expression of the gene being studied. The study of the potential usage of microorganisms as new energy sources requires a complete knowledge of the genomes sequence as well.

Thus in order to be able to realize these and many other studies, it is necessary to close most (if not all) of the genomes being sequenced at JGI. To this date JGI has accumulated a significant amount of experience closing the microbial genomes. During the last two years the combined efforts of three groups (LANL, LLNL and PGF) allowed to finish 45 microbial projects. Our goal is to fulfill the needs of all of the projects undertaken by JGI and complete no less than 50 genomes per year.

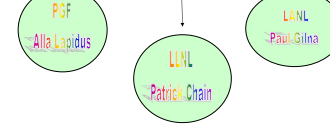
### Genomics Project



### Microbial Finishing

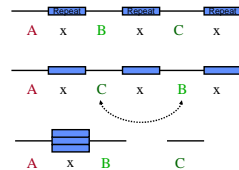


### Finishing Groups



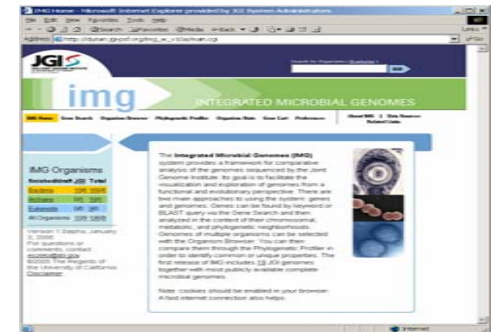
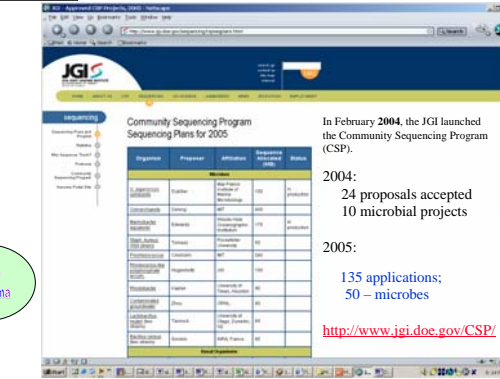
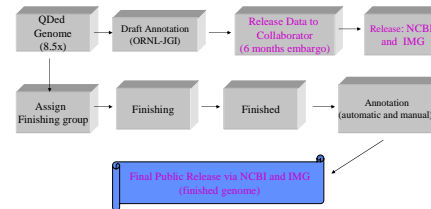
### The Problem With Repeats

- Repeats can make assembly ambiguous, with several possible layouts being equally likely
- Repeats can be up to 5 kb and more in length - plasmid clone inserts are too short to span them (~3 kb). Libraries with larger inserts may provide an answer



An average sequencing read is ~650-700 bp. When in a repeat it is impossible to tell which copy it belongs to

### Data flow - Data Release



(Poster: "Integrated Microbial Genomes (IMG): New Data and Functionality in Version 1.2"; Victor M. Markowitz, Frank Korzeniewski et al)

**Finished Genomes: 44 Draft Genomes: 128**