

The JGI Pipeline for Annotation of Microbial Genomes and Metagenomes

Marcel Huntemann^{1*}, Konstantinos Mavrommatis¹, Natalia Ivanova¹, Natalia Mikhailova¹, Galina Ovchinnikova¹, Andrew Schaumberg¹, Jim Tripp¹, Krishna Palaniappan², Ernest Szeto², I-Min Chen², Nikos Kyrpides¹, Amrita Pati¹,

¹ LBNL - Department of Energy Joint Genome Institute, 2800 Mitchell Drive, Walnut Creek, CA USA

² Lawrence Berkeley National Laboratory Computational Research, Berkeley, CA USA

**To whom correspondence should be addressed:* Email: mhuntemann@lbl.gov

March 21, 2014

ACKNOWLEDGMENTS:

The work conducted by the U.S. Department of Energy Joint Genome Institute is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231.

DISCLAIMER:

LBNL: This document was prepared as an account of work sponsored by the United States Government. While this document is believed to contain correct information, neither the United States Government nor any agency thereof, nor The Regents of the University of California, nor any of their employees, makes any warranty, express or implied, or assumes any legal responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by its trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof, or The Regents of the University of California. The views and opinions of authors expressed herein do not

The JGI Pipeline for Annotation of Microbial Genomes and Metagenomes

necessarily state or reflect those of the United States Government or any agency thereof or The Regents of the University of California.

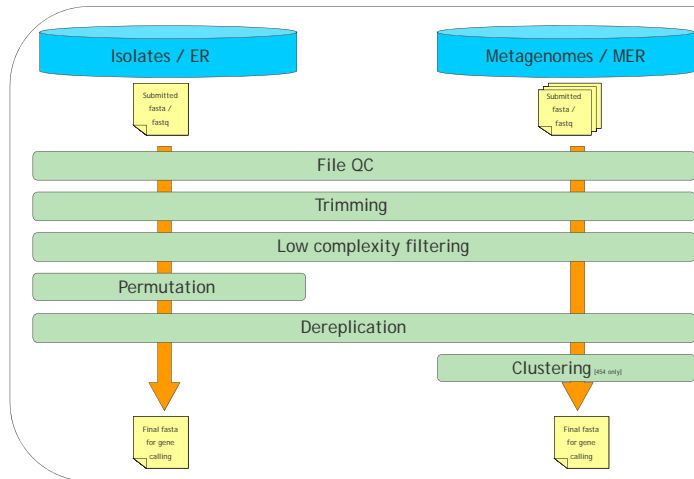
The JGI Pipeline for Annotation of Microbial Genomes and Metagenomes

Marcel Huntemann^{1*}, Konstantinos Mavrommatis¹, Natalia Ivanova¹, Natalia Mikhailova¹, Galina Ovchinnikova¹, Andrew Schaumberg¹, Jim Tripp¹, Krishna Palaniappan², Ernest Szeto², I-Min Chen², Nikos Kyrpides¹, Amrita Pati¹

¹ DOE Joint Genome Institute, Walnut Creek, CA, USA

² Lawrence Berkeley National Laboratory, Berkeley, CA, USA

* Correspondence: Marcel Huntemann (mhuntemann@lbl.gov)



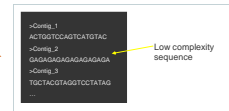
Dataset Pre-Processing

File QC

Validates correctness of input format and sequence alphabet, renames sequences (mapping file provided) and removes sequences that don't meet length thresholds.

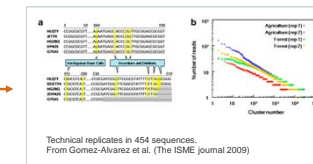
Low complexity filtering

Removes sequences with too many low complexity parts.



Dereplication

Removes very similar sequences that have the same starting sequence.



Clustering

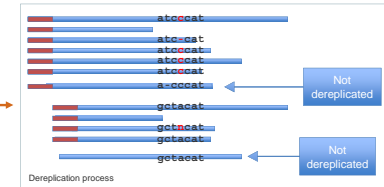
Keeps only the longest one of almost identical sequences (N matches any character and one gap allowed to account for polyN in 454 sequences).

Trimming

Trims off terminal Ns or if fastq file submitted uses appropriate method/cutoff depending on type (Sanger or Solexa).

Permutation

For finished genomes we try to predict the origin of replication and then permute the sequence so that it begins there.



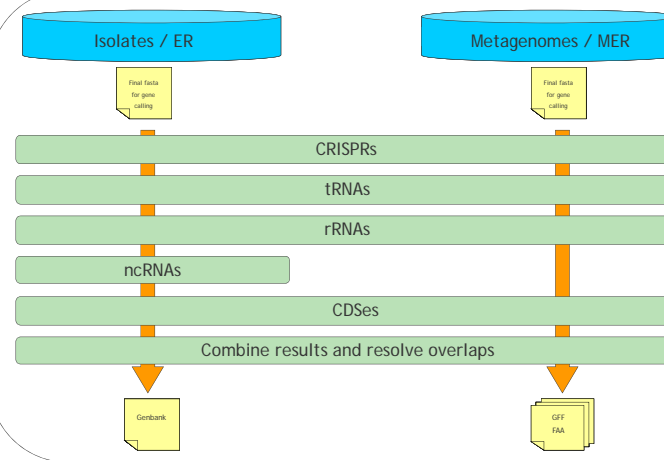
Computational Setup

Structural and functional annotation are embarrassingly parallel steps. JGI-PAMM v1.0, the JGI Pipeline for Annotation of Microbial genomes and Metagenomes, is set up on computational systems at NERSC (www.nersc.gov) and makes efficient use of the provided HPC infrastructure.

The current pipeline has been deployed in an integrated manner since the beginning of 2012 and has processed over **19.5 Billion genes** till date in a fully automated fashion.

The functional annotation of metagenomes is implemented within the Hadoop framework. It uses an on-the-fly approach on the Genepool Phase 2 system and runs out of GPFS instead of HDFS. All compute processes are implemented in the map phase to account for the dynamic set of TaskTracker nodes.

The pipeline is available to users submitting their datasets via IMG's submission system (<http://img.jgi.doe.gov/submit>).



Feature Prediction

CRISPR prediction

Uses a modified **CRT CL1** and **PileRCR2** to check for CRISPR elements.

tRNA prediction

Uses **tRNAscan3** and additionally **Blast4** on the first and last 150 bp of each sequence to account for partial tRNAs.

rRNA prediction

Uses **Hmmer5** and an in-house curated set of rRNA models.

ncRNA prediction

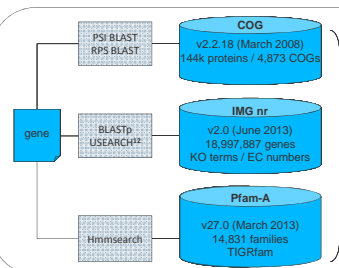
Uses **Blast** to determine which sequence regions have potential hits to **Rfam9** models, extracts those regions and runs **cmsearch7** on them. [genomes side only]

CDS prediction

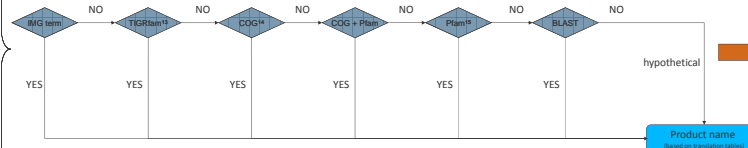
Uses **Prodigal9** for genomes and a combination of **Prodigal**, **GeneMark9**, **Metagene10** and **FragGeneScan11** for metagenomes.

Overlap resolution

Combines prediction results and resolves overlaps based on a curated set of rules for each feature type combination.



Functional Annotation



References

1. Berris, C. et al. CRISPR recognition tool (CRT): a tool for auto-detection of clustered regularly interspaced short palindromic repeats. *BMC Bioinformatics* 8: 231 (2007).
 2. Edgar, R.C. PileUP: fast and accurate identification of CRISPR repeats. *BMC Bioinformatics* 8: 18 (2007).
 3. Lark, T.H. & Eddy, S.R. tRNAscan-SE: a program for improved detection of transfer RNA genes in genomic sequences. *Nucleic Acids Res* 25, 955-64 (1997).
 4. Karch, S.P. et al. Blast: local alignment search tool. *JAMA* 293 (2): 403-410.
 5. Eddy, S.R. A new generation of homology search tools based on probabilistic inference. *Genome Inform* 2009 (2009):201-11.
 6. Collier, P.P. et al. Rfam: updates to the RNA families database. *Nucleic Acids Res* 37, D136-140 (2009).
 7. Naumov, E.P. et al. Internal 1.0: inference of RNA alignments. *Bioinformatics* 25, 1325-1327 (2009).
 8. P. Ye, D. et al. Prodigal: prokaryotic gene recognition and translation initiation site identification. *BMC Bioinformatics* 11, 119 (2010).
 9. Tatusin, A.V. & Bonch-Bruyevich, M. GeneMark-ES: new solutions for gene finding. *Nucleic Acids Res* 36, 1107-15 (2008).
 10. Kapur, H., Park, J. & Tang, J. GeneMark-ES: predicting genes in short and error-prone reads. *Nucleic Acids Res* 34, e151 (2006).
 11. Park, H., Park, H. & Yeo, Y. FragGeneScan: predicting genes in short and error-prone reads. *Nucleic Acids Res* 34, e151 (2006).
 12. Edgar, R.C. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics* 24, 2460-2461 (2009).
 13. Finn, R.D., Clough, J.D. & Wilke, C. The TrEMBL database of protein families. *Nucleic Acids Res* 39, 371-373 (2011).
 14. Pearson, R.L. et al. The COG database: an updated version includes eukaryotes. *BMC Bioinformatics* 4: 41 (2003).
 15. Ye, R.D. et al. The Pfam protein families database. *Nucleic Acids Res* 36, D202-209 (2008).