

KBase: Microbiome and Phylogenomics Capabilities

Dylan Chivian^{*1} (DCChivian@lbl.gov), **Adam P. Arkin¹**, **Robert Cottingham³**, **Chris Henry²**, Benjamin Allen³, Jason Baumohl¹, Jay Bolton¹, Shane Canon¹, Stephen Chan¹, John-Marc Chandonia¹, Zachary Crockett³, Paramvir Dehal¹, Meghan Drake³, Janaka N. Edirisinghe², José P. Faria², Annette Greiner¹, Tianhao Gu², James Jeffryes², Marcin P. Joachimiak¹, Sean Jungbluth¹, Roy Kamimura¹, Keith Keller¹, Vivek Kumar⁵, Sunita Kumari⁵, Miriam Land³, Sebastian Le Bras¹, Zhenyuan Lu⁵, Akiyo Marukawa¹, Sean McCorkle⁴, Cheyenne Nelson¹, Dan Murphy-Olson², Erik Pearson¹, Gavin Price¹, Priya Ranjan³, William Riehl¹, Boris Sadkhin², Samuel Seaver², Alan Seleman², Gwyenth Terry¹, James Thomason⁵, Doreen Ware⁵, Pamela Weisenhorn², Elisha Wood-Charlson¹, Shinjae Yoo⁴, Qizhi Zhang²

¹Lawrence Berkeley National Laboratory, Berkeley, CA; ²Argonne National Laboratory, Argonne, IL; ³Oak Ridge National Laboratory, Oak Ridge, TN; ⁴Brookhaven National Laboratory, Upton, NY; ⁵ Cold Spring Harbor Laboratory, Cold Spring Harbor, NY.

<http://kbase.us>

Project Goals: The Department of Energy Systems Biology Knowledgebase (KBase) is a knowledge creation and discovery environment designed for both biologists and bioinformaticians. KBase integrates a large variety of data and analysis tools, from DOE and other public services, into an easy-to-use platform that leverages scalable computing infrastructure to perform sophisticated systems biology analyses. KBase is a publicly available and developer extensible platform that enables scientists to analyze their own data within the context of public data and share their findings across the system.

KBase was designed to enable systems biology analysis of communities of microbes and/or plants. KBase is extensible and currently includes powerful tools for metabolic modeling, comparative and phylogenomics of microbial genomes that can be used for developing mechanistic understanding of functional interactions between species in microbial ecosystems. Essential to gaining new insight is obtaining high-quality genomes to annotate, either via cultivation or genome extraction, from metagenome assembly. KBase has incorporated and added to a suite of microbiome analysis apps meant to be used in concert, including sequence QA/QC tools such as Trimmomatic and FastQC, taxonomic structure profiling of shotgun metagenome sequence with Kaiju, custom KBase apps for generating sample-specific *in silico* reads for downstream benchmarking, several metagenome assembly algorithms including MEGAHIT, IDBA-UD, and metaSPAdes, custom KBase apps for comparing metagenome assemblies, grouping assembled genome fragments (contigs) into putative genomes (bins) with MaxBin2 and other binners, and genome completeness and contamination assessment and filtering with CheckM. Tools for fractionation of unassembled reads and unbinned contigs to permit taxonomic and functional assessment of unbinned portions of samples are offered, and can also be used to reassemble individual bins of interest. Additionally, we've recently released tools and services that allow users to search rapidly (seconds to minutes) all reference genome databases, metagenomes and published metagenome-assembled genomes (MAGs) using their

reads, assemblies or MAGs. This is implemented using a MinHash like sketching process that works well for identifying matches above ~90% identity.

We have greatly expanded microbiome analysis in KBase. In addition to support for amplicon-based analyses (please see poster KBase: Leveraging Amplicon Analysis Tools to Generate Testable Hypotheses From Complex Natural Communities), it is now possible to incorporate and use tools that enable users to get from shotgun reads through to MAGs to phylogenomics and metabolic modeling. As an example from our initial set of tools, a user can upload or find data from collaborators or the public and apply one of the metagenome assembly apps and bin the assembled contigs so that individual genomes can be extracted from the bins. Once individual MAGs are extracted, the highest quality fraction can be piped into a wide range of downstream analysis apps in KBase, including genome annotation, phylogenetic placement and genome content comparison with respect to one another, KBase reference genomes, and other public genome and MAG collections. Unbinned metagenome assemblies can also have gene annotation for analysis. For high-quality MAGs, metabolic modeling and RNA-seq alignment can be performed (please see poster KBase: Omics driven discovery of novel functional capabilities in biological systems). Pangenome calculations among related genomes can be combined with phylogenetic and functional analysis to capture evolutionary histories of gene families and allow researchers to investigate functional repertoires and niche roles of microbial lineages.

In addition to efforts by KBase developers to expand the functionality of our Microbiome tool suite, community developers have been adding tools that they use and have developed, including members of the DOE Joint Genome Institute (BBMap, MetaBAT2, RQCFiler, JGI Metagenome Assembly Pipeline), the ENIGMA SFA, the LLNL Soils SFA (vConTACT2, VirSorter), and the LANL Bacterial:Fungal Interactions SFA (GOTTCHA2). All Apps in KBase are openly available for users to apply with their own data.

KBase is funded by the Genomic Science program within the U.S. Department of Energy, Office of Science, Office of Biological and Environmental Research under award numbers DE-AC02-05CH11231, DE-AC02-06CH11357, DE-AC05-00OR22725, and DE-AC02-98CH10886.