

KBase: Significant Improvements to the DOE Systems Biology Knowledgebase in 2021

Adam P. Arkin¹, Robert Cottingham³, Chris Henry^{2*}, Benjamin Allen³, Jason Baumohl¹, Kathleen Beilsmith², David Dakota Blair⁴, Jay Bolton¹, Shane Canon¹, Stephen Chan¹, John-Marc Chandonia¹, Dylan Chivian¹, Zachary Crockett³, Paramvir Dehal¹, Ellen Dow¹, Meghan Drake³, Janaka N. Edirisinghe², José P. Faria², Jason Fillman¹, Tianhao Gu², AJ Ireland¹, Marcin P. Joachimiak¹, Sean Jungbluth¹, Roy Kamimura¹, Keith Keller¹, Vivek Kumar⁵, Sunita Kumari⁵, Miriam Land³, Sebastian Le Bras¹, Zhenyuan Lu⁵, Filipe Lui², Dan Murphy-Olson², Erik Pearson¹, Gavin Price¹, Priya Ranjan³, William Riehl¹, Boris Sadkhin², Samuel Seaver², Alan Seleman², Gwyneth Terry¹, Charles Trenholm¹, Sumin Wang¹, Doreen Ware⁵, Pamela Weisenhorn², Elisha Wood-Charlson¹, Ziming Yang⁴, 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; ⁵Cold 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.

Over the last year, substantial new improvements were integrated into the KBase platform, with the goal of empowering upload, integration, and analysis of large-scale multi-omics datasets progressing toward a mechanistic understanding of environmental microbiomes. Given their complexity and diversity, the study of environmental microbiomes may require the analysis of hundreds of samples, often with multi-omics data (e.g., metagenome, amplicon, metabolomics, etc.) collected for each sample. This creates multiple significant challenges in supporting the analysis of such data on a platform like KBase.

First, hundreds of disparate data files must be loaded and integrated into the KBase data-model, a task that was previously quite tedious as most upload tools in KBase operated on one file at a time. A new bulk-upload pipeline has been established, enabling users to load hundreds of objects at once. These tools are now being applied by the Genome Resolved Open Watersheds (GROW) CSP and Ecosystems and Networks Integrated with Genes and Molecular Assemblies (ENIGMA) SFA to load hundreds of metagenomes and thousands of MAGs.

A second challenge involves the association of distinct data types collected from common experimental samples to: (1) allow for the organization of all data related to a complex experiment and (2) organize the data in such a way as to enable sophisticated cross-analysis of various data types. Toward this end, a Samples management system was integrated into KBase, enabling users to load sample metadata into KBase and associate sample IDs with data objects

on system (e.g., linking together the metagenome, metabolome, and MAGs derived from a common sample to facilitate integrated analysis). Tools are now being developed that exploit the relationship amongst data created by these sample associations to permit complex statistical analysis of data. We developed our Samples metadata, templates, and ID management systems in close collaborations with our partners at EMSL, JGI, NMDC, and ESS-DIVE.

A third challenge is handling all the datatypes required for multi-omics microbiome analysis, and to meet this challenge, this year KBase added deeper support and analysis pipelines for metabolomic and amplicon data. Now users can upload these new datatypes and integrate them with their broader microbiome data as well as reference data generally available on KBase (e.g., MGnify). Metabolomics and expression data can be visualized on rich metabolic maps, including new genome-wide maps for plants and fungi. Amplicon data can be visualized, compared, used to predict microbiome functional potential, and analyzed in the context of environmental data related to samples.

A fourth challenge is the difficulty associated with annotating protein functions within microbiomes. To improve support for this activity, and in collaboration with our User Working Groups, new tools have been added to KBase in the area of genome annotation. Users can now upload and compare many alternative annotations using tools developed by the LLNL Biofuels SFA; they can build models from those annotations and compare performance of competing annotations on predicting phenotypes like Biolog growth profiles; and they can apply new annotation algorithms like DRAM (integrated into KBase by the Wrighton lab) (1) and the Exascale Analysis tool (integrated into KBase by the Jacobson lab). Tools for metabolic reconstruction are improving dramatically for plants, fungi, and microbes, including new support for visualizing and painting data on metabolic reconstructions. All of this culminates in KBase offering a sophisticated set of tools to load, store, create, compare, test, and visualize annotations.

While much still remains to be done, this new functionality has significantly lowered the barrier for KBase users to develop a deeper, integrated, mechanistic understanding of complex environmental microbiome data based on rich and diverse experimental datasets.

References

1. Shaffer, M. *et al.* DRAM for distilling microbial metabolism to automate the curation of microbiome function. *Nucleic Acids Research*. 48:16. 8883–8900.
<https://doi.org/10.1093/nar/gkaa621>.

This work is supported as part of the BER Genomic Science Program. The DOE Systems Biology Knowledgebase (KBase) is funded by 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.