

164. Collaborative Workspaces in the DOE Systems Biology Knowledgebase

Christopher S. Henry*¹ (chenry@mcs.anl.gov), Jason Baumohl², Aaron Best¹, Jared Bischof¹, Ben Bowen², Tom Brettin¹, Tom Brown¹, Shane Canon², Stephen Chan², John- Marc Chandonia², Dylan Chivian², Ric Colasanti¹, Neal Conrad¹, Brian Davison³, Matt DeJongh⁶, Paramvir Dehal², Narayan Desai¹, Scott Devoid¹, Terry Disz¹, Meghan Drake³, Janaka Edirisinghe¹, Gang Fang⁷, José Pedro Lopes Faria¹, Mark Gerstein⁷, Elizabeth M. Glass¹, Annette Greiner², Dan Gunter², James Gurtowski⁵, Nomi Harris², Travis Harrison¹, Fei He⁴, Matt Henderson², Adina Howe¹, Marcin Joachimiak², Kevin Keegan¹, Keith Keller², Guruprasad Kora³, Sunita Kumari⁵, Miriam Land³, Folker Meyer¹, Steve Moulton³, Pavel Novichkov², Taeyun Oh⁸, Gary Olsen⁹, Bob Olson¹, Dan Olson¹, Ross Overbeek¹, Tobias Paczian¹, Bruce Parrello¹, Shiran Pasternak⁵, Sarah Poon², Gavin Price², Srivdya Ramakrishnan⁵, Priya Ranjan³, Bill Riehl², Pamela Ronald⁸, Michael Schatz⁵, Lynn Schriml¹⁰, Sam Seaver¹, Michael W. Sneddon², Roman Sutormin², Mustafa Syed³, James Thomason⁵, Nathan Tintle⁶, Will Trimble¹, Daifeng Wang⁷, Doreen Ware⁵, David Weston³, Andreas Wilke¹, Fangfang Xia¹, Shinjae Yoo⁴, Dantong Yu⁴, **Robert Cottingham**³, **Sergei Maslov**⁴, **Rick Stevens**¹, **Adam P. Arkin**²

¹Argonne National Laboratory, Argonne, IL, ²Lawrence Berkeley National Laboratory, Berkeley, CA, ³Oak Ridge National Laboratory, Oak Ridge, TN, ⁴Brookhaven National Laboratory, Upton, NY, ⁵Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, ⁶Hope College, Holland, MI, ⁷Yale University, New Haven, CT, ⁸University of California, Davis, CA, ⁹University of Illinois at Champaign-Urbana, Champaign, IL, ¹⁰University of Maryland, College Park, MD

<http://kbase.us>

Project Goals: The KBase project aims to provide the capabilities needed to address the grand challenge of systems biology: to predict and ultimately design biological function. KBase enables users to collaboratively integrate the array of heterogeneous datasets, analysis tools and workflows needed to achieve a predictive understanding of biological systems. It incorporates functional genomic and metagenomic data for thousands of organisms, and diverse tools for (meta)genomic assembly, annotation, network inference and modeling, allowing researchers to combine diverse lines of evidence to create increasingly accurate models of the physiology and community dynamics of microbes and plants. KBase will soon allow models to be compared to observations and dynamically revised. A new prototype Narrative interface lets users create a reproducible record of the data, computational steps and thought process leading from hypothesis to result in the form of interactive publications.

Collaboration, provenance, data consolidation, data standards, and data validation are all major objectives of the KBase project, and the Workspace is the engine being used to drive towards these goals. The Workspace is an online data store in KBase, where the wide variety of data entities (e.g., metabolic models, genome annotations, phenotype data) being produced and consumed by KBase analysis pipelines are stored.

User-generated data including genomes, contigs, phenotype data, and expression data may be loaded directly into the Workspace for analysis by KBase tools. The Workspace stores the many derived data products (e.g., genome annotations, metabolic models) produced by KBase analyses, with all objects generated from an analysis pipeline being interconnected all the way back to the raw data. For example, a

metabolic model is linked to the genome it was constructed from, the genome is linked to the contigs with the genome sequence, and the contigs are linked to the reads they were assembled from. The Workspace also facilitates collaboration by enabling project data, which could encompass thousands of objects, to be rapidly shared with a few select collaborators or—if desired— with the entire research community (for example, during the publication phase of a research project).

The Workspace supports full versioning of all data objects, enabling a researcher to rapidly access, compare, and restore prior versions of all objects stored in the Workspace. Provenance is also provided for data in the Workspace, providing researchers with detailed information on the analysis pipelines and parameters that produced each data object. We anticipate that these features will greatly facilitate the evaluation, validation, and replication of even complex systems biology studies.

All objects stored in the Workspace are typed, with each data type specified in detail (e.g., required fields, field types, field indexing). For example, a Genome object in the workspace must include a scientific name, a domain, and a list of features with functional annotations and locations on the chromosome. Data objects are validated against these specifications, enabling the Workspace to enforce data standards, as well as metadata requirements. The Workspace also provides a rich API for users to specify their own new data types, supporting the rapid integration of new tools, data, and analysis pipelines by the KBase user community. In this poster, we will highlight the functionality encompassed in the Workspace, we will explore the typed objects already available within the Workspace, and we will demonstrate how these typed objects connect to the many analysis tools and pipelines already available in KBase.

KBase is funded by the U.S. Department of Energy, Office of Science, Office of Biological and Environmental Research.