

166. Design of KBase Infrastructure

Thomas Brettin*¹ (brettin@cels.anl.gov), Daniel Olson¹, Jason Baumohl², Aaron Best³, Jared Bischof¹, Ben Bowen², Tom Brown¹, Shane Canon¹, Stephen Chan², John-Marc Chandonia², Dylan Chivian², Ric Colasanti¹, Neal Conrad¹, Brian Davison⁴, Matt DeJongh³, Paramvir Dehal², Narayan Desai¹, Scott Devoid¹, Terry Disz¹, Meghan Drake⁴, Janaka Edirisinghe¹, Gang Fang⁷, José Pedro Lopes Faria¹, Mark Gerstein⁷, Elizabeth M. Glass¹, Annette Greiner², Dan Gunter², James Gurtowski⁶, Nomi Harris², Travis Harrison¹, Fei He⁵, Matt Henderson², Chris Henry¹, Adina Howe¹, Marcin Joachimiak², Kevin Keegan¹, Keith Keller², Guruprasad Kora⁴, Sunita Kumari⁶, Miriam Land⁴, Folker Meyer¹, Steve Moulton⁴, Pavel Novichkov², Taeyun Oh⁸, Gary Olsen⁹, Bob Olson¹, Dan Olson¹, Ross Overbeek¹, Tobias Paczian¹, Bruce Parrello¹, Shiran Pasternak⁶, Sarah Poon², Gavin Price², Srividya Ramakrishnan⁶, Priya Ranjan⁴, Bill Riehl², Pamela Ronald⁸, Michael Schatz⁶, Lynn Schriml¹⁰, Sam Seaver¹, Michael W. Sneddon², Roman Sutormin², Mustafa Syed⁴, James Thomason⁶, Nathan Tintle³, Will Trimble¹, Daifeng Wang⁷, Doreen Ware^{5,6}, David Weston⁴, Andreas Wilke¹, Fangfang Xia¹, Shinjae Yoo⁵, Dantong Yu⁵, **Robert Cottingham⁴, Sergei Maslov⁵, Rick Stevens¹, Adam P. Arkin²**

¹Argonne National Laboratory, Argonne, IL, ²Lawrence Berkeley National Laboratory, Berkeley, CA, ³Hope College, Holland, MI, ⁴Oak Ridge National Laboratory, Oak Ridge, TN, ⁵Brookhaven National Laboratory, Upton, NY, ⁶Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, 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.

At the core of the KBase architecture is a set of rich data models and stores, scalable computing, and workflow management. Our KBase physical infrastructure is built on the successes of DOE investment in our national scientific cyber-infrastructure and therefore leverages enormous intellectual resources present in the DOE community. Building on ESNet allows us to construct a wide area network between the partner labs that enables a virtual hardware infrastructure. Our use of cloud-computing supports development of new tools and provides compute resources for production services. The acceptance of virtualization technology is growing, and the use of machine images produced by others is already visible in our core services. Additionally, machine images are now provided which contain multiple components of the KBase infrastructure and services. Cluster Computing has long been a critical part of biological data

analysis. In collaboration with computing centers created by the Office of Advanced Computing Research such as NERSC, our underlying cluster services can leverage these resources and scale to meet needs.

KBase aims to power the next wave of biological research in DOE and beyond. Enabling these capabilities requires a software and hardware infrastructure that is integrated, extensible, and scalable. The architecture is designed to meet these needs and support user functionality to visualize data, create models or design experiments based on KBase- generated suggestions.

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