

## 166. Design of KBase Infrastructure

Thomas Brettin\*<sup>1</sup> ([brettin@cels.anl.gov](mailto:brettin@cels.anl.gov)), Daniel Olson<sup>1</sup>, Jason Baumohl<sup>2</sup>, Aaron Best<sup>3</sup>, Jared Bischof<sup>1</sup>, Ben Bowen<sup>2</sup>, Tom Brown<sup>1</sup>, Shane Canon<sup>1</sup>, Stephen Chan<sup>2</sup>, John-Marc Chandonia<sup>2</sup>, Dylan Chivian<sup>2</sup>, Ric Colasanti<sup>1</sup>, Neal Conrad<sup>1</sup>, Brian Davison<sup>4</sup>, Matt DeJongh<sup>3</sup>, Paramvir Dehal<sup>2</sup>, Narayan Desai<sup>1</sup>, Scott Devoid<sup>1</sup>, Terry Disz<sup>1</sup>, Meghan Drake<sup>4</sup>, Janaka Edirisinghe<sup>1</sup>, Gang Fang<sup>7</sup>, José Pedro Lopes Faria<sup>1</sup>, Mark Gerstein<sup>7</sup>, Elizabeth M. Glass<sup>1</sup>, Annette Greiner<sup>2</sup>, Dan Gunter<sup>2</sup>, James Gurtowski<sup>6</sup>, Nomi Harris<sup>2</sup>, Travis Harrison<sup>1</sup>, Fei He<sup>5</sup>, Matt Henderson<sup>2</sup>, Chris Henry<sup>1</sup>, Adina Howe<sup>1</sup>, Marcin Joachimiak<sup>2</sup>, Kevin Keegan<sup>1</sup>, Keith Keller<sup>2</sup>, Guruprasad Kora<sup>4</sup>, Sunita Kumari<sup>6</sup>, Miriam Land<sup>4</sup>, Folker Meyer<sup>1</sup>, Steve Moulton<sup>4</sup>, Pavel Novichkov<sup>2</sup>, Taeyun Oh<sup>8</sup>, Gary Olsen<sup>9</sup>, Bob Olson<sup>1</sup>, Dan Olson<sup>1</sup>, Ross Overbeek<sup>1</sup>, Tobias Paczian<sup>1</sup>, Bruce Parrello<sup>1</sup>, Shiran Pasternak<sup>6</sup>, Sarah Poon<sup>2</sup>, Gavin Price<sup>2</sup>, Srividya Ramakrishnan<sup>6</sup>, Priya Ranjan<sup>4</sup>, Bill Riehl<sup>2</sup>, Pamela Ronald<sup>8</sup>, Michael Schatz<sup>6</sup>, Lynn Schriml<sup>10</sup>, Sam Seaver<sup>1</sup>, Michael W. Sneddon<sup>2</sup>, Roman Sutormin<sup>2</sup>, Mustafa Syed<sup>4</sup>, James Thomason<sup>6</sup>, Nathan Tintle<sup>3</sup>, Will Trimble<sup>1</sup>, Daifeng Wang<sup>7</sup>, Doreen Ware<sup>5,6</sup>, David Weston<sup>4</sup>, Andreas Wilke<sup>1</sup>, Fangfang Xia<sup>1</sup>, Shinjae Yoo<sup>5</sup>, Dantong Yu<sup>5</sup>, **Robert Cottingham<sup>4</sup>, Sergei Maslov<sup>5</sup>, Rick Stevens<sup>1</sup>, Adam P. Arkin<sup>2</sup>**

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