

## **KBase: The Systems Biology Knowledgebase for Predictive Biological and Environmental Research in an Integrated Data Platform**

**Adam P. Arkin<sup>1</sup>, Robert Cottingham<sup>3</sup>, Chris Henry<sup>2</sup>**, Benjamin Allen\*<sup>3</sup> (allenbh@ornl.gov), Jason Baumohl<sup>1</sup>, Jay Bolton<sup>1</sup>, Shane Canon<sup>1</sup>, Stephen Chan<sup>1</sup>, John-Marc Chandonia<sup>1</sup>, Dylan Chivian<sup>1</sup>, Zachary Crockett<sup>3</sup>, Paramvir Dehal<sup>1</sup>, Meghan Drake<sup>3</sup>, Janaka N. Edirisinghe<sup>2</sup>, José P. Faria<sup>2</sup>, Annette Greiner<sup>1</sup>, Tianhao Gu<sup>2</sup>, James Jeffryes<sup>2</sup>, Marcin P. Joachimiak<sup>1</sup>, Sean Jungbluth<sup>1</sup>, Roy Kamimura<sup>1</sup>, Keith Keller<sup>1</sup>, Vivek Kumar<sup>5</sup>, Sunita Kumari<sup>5</sup>, Miriam Land<sup>3</sup>, Sebastian Le Bras<sup>1</sup>, Zhenyuan Lu<sup>5</sup>, Akiyo Marukawa<sup>1</sup>, Sean McCorkle<sup>4</sup>, Cheyenne Nelson<sup>1</sup>, Dan Murphy-Olson<sup>2</sup>, Erik Pearson<sup>1</sup>, Gavin Price<sup>1</sup>, Priya Ranjan<sup>3</sup>, William Riehl<sup>1</sup>, Boris Sadkhin<sup>2</sup>, Samuel Seaver<sup>2</sup>, Alan Seleman<sup>2</sup>, Gwyenth Terry<sup>1</sup>, James Thomason<sup>5</sup>, Doreen Ware<sup>5</sup>, Pamela Weisenhorn<sup>2</sup>, Elisha Wood-Charlson<sup>1</sup>, Shinjae Yoo<sup>4</sup>, Qizhi Zhang<sup>2</sup>

<sup>1</sup>Lawrence Berkeley National Laboratory, Berkeley, CA; <sup>2</sup>Argonne National Laboratory, Argonne, IL; <sup>3</sup>Oak Ridge National Laboratory, Oak Ridge, TN; <sup>4</sup>Brookhaven National Laboratory, Upton, NY; <sup>5</sup> 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.**

The U.S. Department of Energy (DOE) supports biological and environmental research to investigate the complex interactions within biological systems and the processes that shape soil, water, and ecological dynamics of our biosphere, and to harness these processes for sustainable production of energy and materials. KBase is an open-source data-science platform funded by DOE to enable sharing, integration, and analysis of many types of data associated with microbes, plants, and their communities using scalable computing infrastructure. This extensive public resource is designed to facilitate large-scale analyses of biological and environmental systems while accelerating scientific discovery, improve reproducibility, and foster open collaboration.

KBase offers a suite of scientific applications to enable users to build sophisticated analytical workflows, share their findings, and organize their projects. Over 200 apps in KBase offer diverse scientific functionality across the realms of comparative genomics, microbial communities, metabolic modeling, and transcriptomics. Several tools and services in KBase have been co-developed with the DOE Joint Genome Institute. Users can build and share sophisticated workflows through a combination of chaining together multiple analysis tools, writing scripts for automation, and using batch processing, all within notebook-style *Narratives* that contain the employed data and tools. Projects, laboratories, and even whole institutions can organize their users and associated *Narratives* into a shared *Organization* with multiple permission levels and

management features. Developers can build, test, register, and deploy new or existing software as KBase apps using the Software Development Kit, thereby extending the platform's scientific capabilities.

Newly added services include several tools collaboratively developed by DOE SFA programs and KBase staff, including NWchem, GOTTCHA2, VIRSorter, and vConTACT2. KBase is unique in offering these diverse and integrated capabilities to a growing community of several thousand users. A central premise of KBase is to maximize the impact of data shared and developed between all users on system. This is the premise of KBase's emerging Knowledge Engine technology, which draws inferences between user data and public repositories, so scientists can better understand their work in the context of public knowledge. As data propagates across the system, it can be continually updated with new information like predictions of gene function, allowing the user to see how these changes scale from genes to ecologies and better predict outcomes.

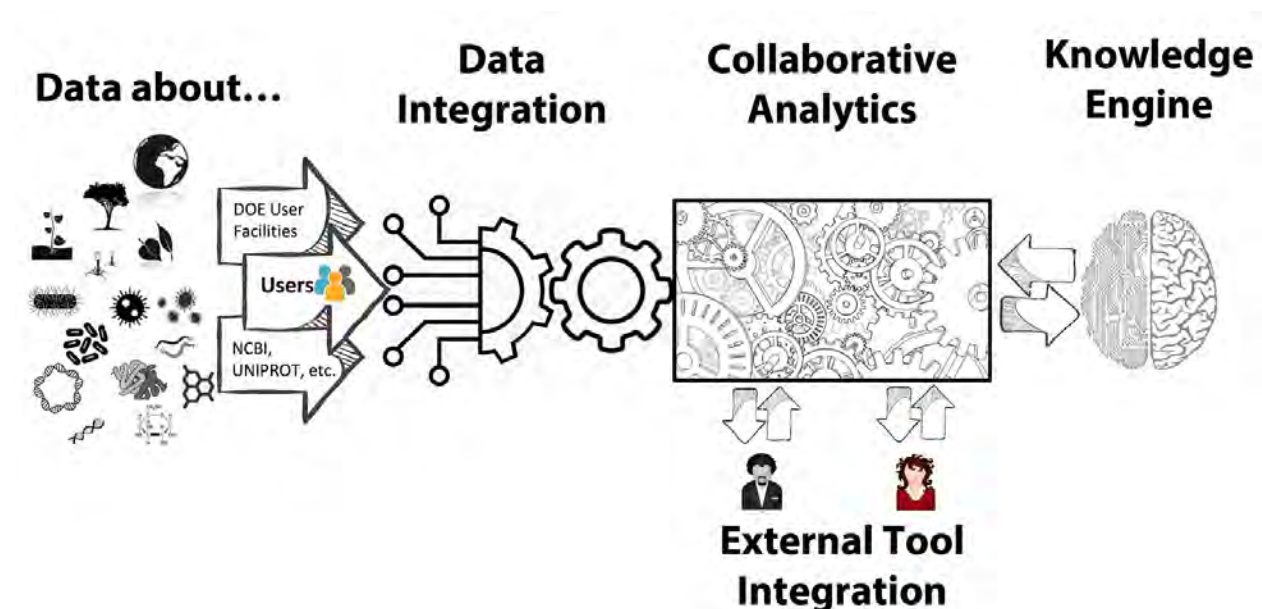


Figure 1. Integration of data and tools into KBase.

*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.*