

The DOE Systems Biology Knowledgebase (KBase): Progress towards a system for collaborative and reproducible inference and modeling of biological function

Adam P. Arkin^{*1} (aparkin@lbl.gov), Jason Baumohl¹, Aaron Best², Jared Bischof², Ben Bowen¹, Tom Brettin², Tom Brown², Shane Canon¹, Stephen Chan¹, John-Marc Chandonia¹, Dylan Chivian¹, Ric Colasanti², Neal Conrad², Brian Davison³, Matt DeJongh⁶, Paramvir Dehal¹, Narayan Desai², Scott Devoid², Terry Disz², Meghan Drake³, Janaka Edirisinghe², Gang Fang⁷, José Pedro Lopes Faria², Mark Gerstein⁷, Elizabeth M. Glass², Annette Greiner¹, Dan Gunter¹, James Gurtowski⁵, Nomi Harris¹, Travis Harrison², Fei He⁴, Matt Henderson¹, Chris Henry², Adina Howe², Marcin Joachimiak¹, Kevin Keegan², Keith Keller¹, Guruprasad Kora³, Sunita Kumari⁵, Miriam Land³, Folker Meyer², Steve Moulton³, Pavel Novichkov¹, Taeyun Oh⁸, Gary Olsen⁹, Bob Olson², Dan Olson², Ross Overbeek², Tobias Paczian², Bruce Parrello², Shiran Pasternak⁵, Sarah Poon¹, Gavin Price¹, Srividya Ramakrishnan⁵, Priya Ranjan³, Bill Riehl¹, Pamela Ronald⁸, Michael Schatz⁵, Lynn Schriml¹⁰, Sam Seaver², Michael W. Sneddon¹, Roman Sutormin¹, Mustafa Syed³, James Thomason⁵, Nathan Tintle⁶, Will Trimble², Daifeng Wang⁷, Doreen Ware⁵, David Weston³, Andreas Wilke², Fangfang Xia², Shinjae Yoo⁴, Dantong Yu⁴, Bob Cottingham³, Sergei Maslov⁴, Rick Stevens²

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, 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 computational 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 including (meta)genomic assembly, annotation, network inference and modeling, thereby 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.

Systems biology is driven by the ever-increasing wealth of data resulting from new generations of genomics-based technologies. With the success of genome sequencing, biology began to generate and accumulate data at an exponential rate. In addition to the massive stream of sequencing data, each type of technology that researchers use to analyze a sequenced organism adds another layer of complexity to the challenge of understanding how different biological components work together to form a functional living system. Achieving this systemslevel understanding of biology will enable researchers to predict and ultimately design how biological systems will function under certain conditions. A collaborative computational environment is needed to bring researchers together so they can share and integrate large, heterogeneous datasets and readily use this information to develop predictive models that drive scientific discovery.

The advancement of systems biology relies not only on sharing the results of projects through traditional