

KBase : Omics driven discovery of novel functional capabilities in biological systems

Janaka N. Edirisinghe*² (janakae@anl.gov), Benjamin Allen³, Jason Baumohl¹, Jay Bolton¹, Shane Canon¹, Stephen Chan¹, John-Marc Chandonia¹, Dylan Chivian¹, Zachary Crockett³, Paramvir Dehal¹, Meghan Drake³, José P. Faria², Annette Greiner¹, Tianhao Gu², James Jeffryes², Marcin P. Joachimiak¹, Sean Jungbluth¹, Roy Kamimura¹, Keith Keller¹, Vivek Kumar⁵, Sunita Kumari⁵, Miriam Land³, Sebastian Le Bras¹, Zhenyuan Lu⁵, Akiyo Marukawa¹, Sean McCorkle⁴, Cheyenne Nelson¹, Dan Murphy-Olson², Erik Pearson¹, Gavin Price¹, Priya Ranjan³, William Riehl¹, Boris Sadkhin², Samuel Seaver², Alan Seleman², Gwyenth Terry¹, James Thomason⁵, Doreen Ware⁵, Pamela Weisenhorn², Elisha Wood-Charlson¹, Shinjae Yoo⁴, Qizhi Zhang², **Robert Cottingham³, Chris Henry² and Adam P. Arkin¹**

¹Lawrence Berkeley National Laboratory, Berkeley, CA; ²Argonne National Laboratory, Argonne, IL; ³Oak Ridge National Laboratory, Oak Ridge, TN; ⁴Brookhaven National Laboratory, Upton, NY; ⁵ Cold Spring Harbor Laboratory, Cold Spring Harbor, NY.

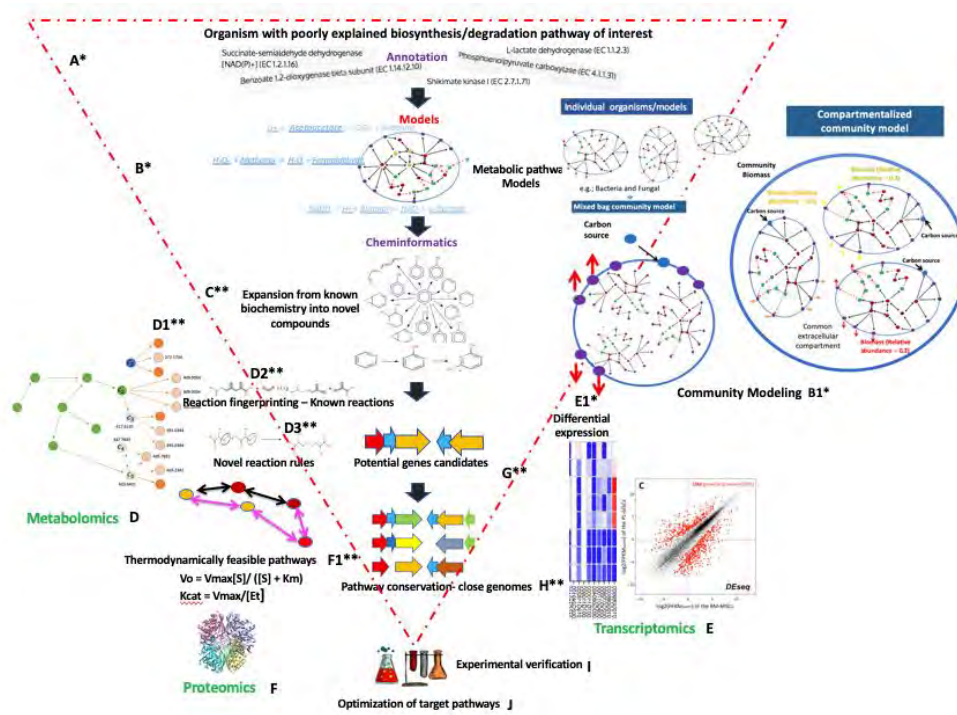
Project Goals: Discovery and characterization of novel biochemistry is key to understanding the behavior of complex biological systems. To explain a microbe's role in biogeochemical cycles, or synthesis of various biofuel products from plant biomass, or to better understand community interactions within complex biological systems, it is essential to identify poorly explained or novel biochemical pathways. Though recent advances in system biology and the exponential growth of reference data suggest that there is an enormous amount of untapped enzyme potential, the traditional approaches for discovery of new functions and pathways are still mostly trial and error experimental processes. Using KBase, we demonstrate a mechanistic-modeling-based approach coupled with multi-omics data to support high throughput discovery of new metabolic pathways.

While there are system-biology approaches that have been developed in isolation for novel pathway identification, it has been a great challenge to translate raw data into an improved understanding of microbe-mediated chemical transformations in degradation and biosynthesis pathways. In KBase, users will be able to combine genomic, metabolic modeling and cheminformatic predictions reconciled with omics-data to predict novel pathway reactions as depicted in Figure 1. Currently, users are able to upload sequencing data then assemble, annotate and build metabolic models from isolate genomes or from metagenomes. Multi-omics data such as chemical abundance data (e.g.; FTICR data, MS2 metabolomics data) (Fig 1D), RNA-seq data (Fig 1E), and proteomics data (Fig 1F) can then be mapped onto models and applied to predict active known pathways, reactions, and genes. This is done by optimizing the extent to which active reactions in the models are associated with genes with high expression, quantitatively measured enzyme levels and positively identified metabolites. However, it is often the case that much of the chemical abundance data does not get mapped onto the existing biochemistry databases, which suggests that many functioning pathways remain unidentified.

To predict these new pathways and compounds, we use a cheminformatics pipeline in order to expand the existing biochemistry based on enzymatic and spontaneous chemical rules, then reconcile against unmapped metabolomics data (Fig. 1C-D). With a continuously interconnected

metabolic network of known and predicted reactions, it is feasible to activate all pathways that are implicated based on omics data. While this strategy can generate multiple probable pathway routes, we would be able to filter out high confidence pathways based on thermodynamic feasibility of the predicted reactions (Fig. 1F) and also by eliminating the pathways that lack metabolomic, transcriptomic, or proteomic evidence (Fig 1. D, F, E). By having a set of high confidence novel reactions, we would be able to map potential gene candidates via chemical and or structural based gene finding approaches (reaction finger-printing) (Fig 1C, D2, D3) coupled with genes that are expressed according to transcriptomic and/or proteomic data. Finally, a subset of these predicted genes can be validated experimentally (Fig 1J).

Figure 1. The figure shows the pathway discovery pipeline that has been implemented in DOE KBase and the integration of omics data at certain levels. The red triangle indicates the process of narrowing down into few target genes from thousands of genes in a genome for experimental validation. Steps labeled with a single asterisk (*) indicate



that the specific component of the pipeline has already been implemented in KBase while the labels with a double asterisk (**) indicate that the functionality is still being developed. **B1: Community Modeling** - Construction of community models from individual models. There are two types of models that can be built in understanding the microbial communities. (i) Mixed-bag models (on left): Able to assess the overall metabolic capability of a community (ii) Compartmentalized community models (on right) able to assess the contribution of each member in community

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.