

KBase: The United States Department of Energy Systems Biology Knowledgebase

To the Editor: Over the past two decades, the scale and complexity of genomics technologies and data have advanced from sequencing genomes of a few organisms to generating metagenomes, genome variation, gene expression, metabolites, and phenotype data for thousands of organisms and their communities. A major challenge in this data-rich age of biology is integrating heterogeneous and distributed data into predictive models of biological function, ranging from a single gene to entire organisms and their ecologies. The US Department of Energy (DOE) has invested substantially in efforts to understand the complex interplay between biological and abiotic processes that influence soil, water, and environmental dynamics of our biosphere. The community that has grown around these efforts recognizes the need for scientists of diverse backgrounds to have access to sophisticated computational tools that enable them to analyze complex and heterogeneous data sets and integrate their data and results effectively with the work of others. In this way, new data and conclusions can be rapidly propagated across existing, related analyses and easily discovered by the community for evaluation and comparison with previous results^{1–3}.

Here we present the DOE Systems Biology Knowledgebase (KBase, <http://kbase.us>), an open-source software and data platform that enables data sharing, integration, and analysis of microbes, plants, and their communities. KBase maintains an internal reference database that consolidates information from widely used external data repositories. This includes over 90,000 microbial genomes from RefSeq⁴, over 50 plant genomes from Phytozome⁵, over 300 Biolog media formulations⁶, and >30,000 reactions and compounds from KEGG⁷, BIGG⁸, and MetaCyc⁹. These public data are available for integration with user data where appropriate (e.g., genome comparison or building species trees). KBase links these diverse data types with a range of analytical functions within a web-based user interface. This extensive community resource facilitates large-scale analyses on scalable computing infrastructure and has

the potential to accelerate scientific discovery, improve reproducibility, and foster open collaboration (**Supplementary Note 1**).

Although similar integrative tools exist (**Supplementary Note 2**), no other open platform shares all KBase's features, which include the following: (i) comprehensive support for data provenance and analysis reproducibility; (ii) a flexible system for sharing data and workflows; (iii) an integrated database of genomes and biochemistry; (iv) a point-and-click interface that enables users to build, store, run, and share complex scientific analyses of fully integrated data; (v) built-in support for the use of custom code interleaved with point-and-click apps; and (vi) a software development kit that enables external developers to add applications to KBase (**Supplementary Table 1**). KBase has a suite of scientific applications that enables users to build and share sophisticated workflows. For example, a user can predict species interactions from metagenomic data by assembling raw reads, binning assembled contigs by species, annotating genomes, aligning RNA-seq reads, and reconstructing and analyzing individual and community metabolic models. KBase supports numerous branch points, alternative pipelines, alternative entry points, and internal curation loops that facilitate a wide range of scientific analyses, some of which are not available elsewhere (e.g., merging individual metabolic models into community models and using these to predict interspecies interactions). Although KBase was developed to support analysis of microbes, plants, and their communities, it is potentially applicable to any area of science. There is, however, a policy on use restriction for projects that require HIPAA compliance.

KBase's primary user interface, the Narrative Interface, provides a user experience distinct from other analysis platforms available today, although it shares some common features with a few other systems (**Supplementary Note 2**). From this interface, which is built on the Jupyter^{10,11} platform, users can upload their private data, search and retrieve extensive public reference data, access data shared

by others, share their data with others, select and run applications on their data, view and analyze the results from those applications, and record their thoughts and interpretations along with the analysis steps. These activities take place within a point-and-click 'notebook' environment (**Fig. 1**). When a user begins a new computational experiment in KBase, they create a new 'notebook' (referred to as a Narrative in KBase) to hold this experiment. Every action performed by a user appears as a 'cell' in the Narrative. App cells show the chosen input parameters for the application and the results of the analysis. Markdown cells allow users to add formatted text and figures to a Narrative to describe the thought process behind the scientific workflow being crafted.

A finished Narrative is a precise record of everything the authors did to complete their analysis. Although Narratives are private by default, users may choose to make their Narratives public, or share them with other individual users. This recording of a user's KBase activities within a sharable Narrative is a central pillar of KBase's support for reproducible, transparent research (**Supplementary Note 1**). Once a Narrative has been shared or made public, other users can copy the Narrative and rerun it on their own data, or modify it to suit their scientific needs. Thus, public Narratives serve as resources for the user community by capturing valuable data sets, associated computational analyses, and scientific context describing the rationale behind a scientific study in a form that is immediately reproducible and reusable. A growing number of public Narratives are available in KBase, some of which are showcased in the Narrative Library (<https://kbase.us/narrative-library/>).

The data model in KBase is fundamental to supporting reproducibility and collaboration. KBase is built upon an object-oriented data model where each object instance is automatically versioned and linked to provenance information describing how it was generated. Each data object is also associated with the specific Narrative in which it was uploaded or generated. When a Narrative is shared or cop-

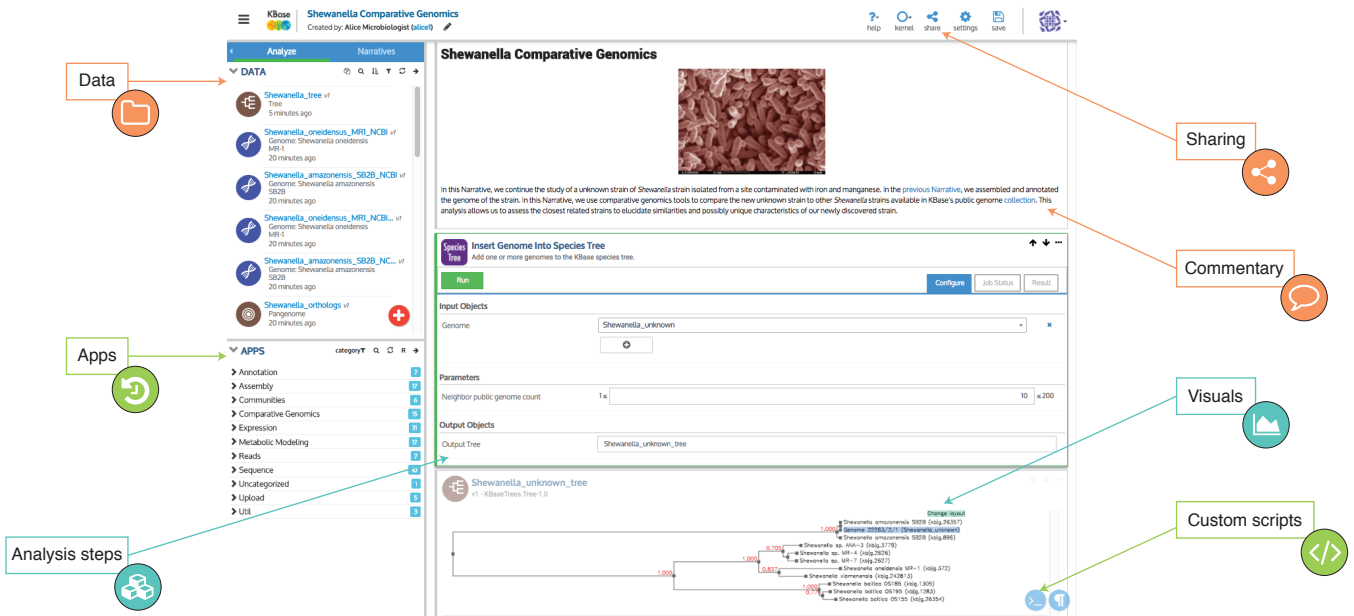


Figure 1 KBase Narratives. A Narrative is an interactive, dynamic, and persistent document created by users that promotes open, reproducible, and collaborative science.

ied, all its input and output data are shared or copied with it. Currently supported data types include reads, contigs, genomes, metabolic models, growth media, RNA-seq, expression, growth phenotype data, and flux balance analysis solutions. This set of types can be extended to support new apps and functionality.

Many existing systems (**Supplementary Note 2**) provide similar support for object-level sharing and provenance, but these systems operate on raw files only, without integration into a common data model. In KBase, objects are not simple files—they are explicitly defined and validated data structures, within which associated objects are linked to one another. For example, a metabolic model is linked to its associated genome, which is linked to its associated taxonomy. This data model enhances interoperability by requiring apps to operate on a common data representation. Furthermore, it enhances awareness of interdependence so users can be notified when an object on which an analysis is based has been updated, and it will ultimately enable data discovery and meta-analysis across the KBase platform.

Presently, KBase has over 160 apps (<https://narrative.kbase.us/#appcatalog>) offering diverse scientific functionality for (meta)genome assembly, contig binning, genome annotation, sequence homology analysis, tree building, comparative genomics, metabolic modeling, community modeling, gap-filling, RNA-seq processing, and expression analysis (see **Supplementary Note 2** for references). Apps interoperate seamlessly to enable a range of scientific workflows (**Fig. 2**). For reproducibility, all apps in KBase

are containerized in versioned Docker modules, enabling a user to run any version at any time.

In addition to running apps, users can create and run blocks of code within a Narrative using “code cells.” KBase has an application-programming interface (API) that allows users to call any KBase app programmatically from within these code cells. This enables users to, for example, run large-scale studies in KBase (e.g., building thousands of models at once) by using loops within a code cell (**Supplementary Note 1; Fig. 1**). Users can also leverage the flexibility of code cells to add custom analysis steps that are not yet available as KBase apps.

Although there are other systems that allow users to create workflows consisting of a series of analysis tool runs and code blocks, the app functionality in KBase differs from these systems in several ways (**Supplementary Note 2**). Currently, KBase’s capabilities for community model reconstruction, plant model reconstruction, community model gap-filling, and expression data model integration are unique to the KBase platform (**Supplementary Note 1; Fig. 2**).

KBase was designed to be an extensible community resource. This extensibility is supported by the KBase Software Development Kit (SDK), which is a set of command-line tools and a web interface that enable any developer to build, register, and deploy new or existing software as KBase apps, thereby extending the platform’s scientific capabilities. All software contributed to the central KBase software repository must adhere to a standard open-source license (<https://opensource.org/licenses>). Information about the app developer is maintained in the documenta-

tion for that app so credit can be given to the contributor. Data provenance, job management, usage logging, and app versioning are handled automatically by the platform, allowing developers with minimal training to package scientific tools in a form that makes them accessible to users within KBase. Other existing platforms offer similar support for third-party development (**Supplementary Note 2**), but KBase’s data model provides the additional benefit of improving interoperability of third-party applications by imposing a single data format and specification on all data types consumed or produced by each app. More information about the KBase SDK is available at https://github.com/kbase/kb_sdk/blob/master/README.md.

Many users have already discovered and applied KBase to meet their scientific needs. As of September 2017, over 3,000 users have KBase accounts, and users have created over 5,000 Narratives. These Narratives contain a total of over 250,000 data objects, or an average of 96 data objects and five apps per Narrative. Science done within Kbase, which has been published in over 30 peer-reviewed publications (**Supplementary Note 1; http://kbase.us/publications**), includes reconstruction of >8,000 models of core metabolism across the microbial tree of life¹², reconstruction of semi-curated metabolic models for 773 gut microbes¹³, prediction of trophic interactions within a microbial community¹⁴, and reconstruction of regulons from expression data¹⁵.

Much of the research performed within KBase has been publicly shared as Narratives that any user can view, copy, and rerun.

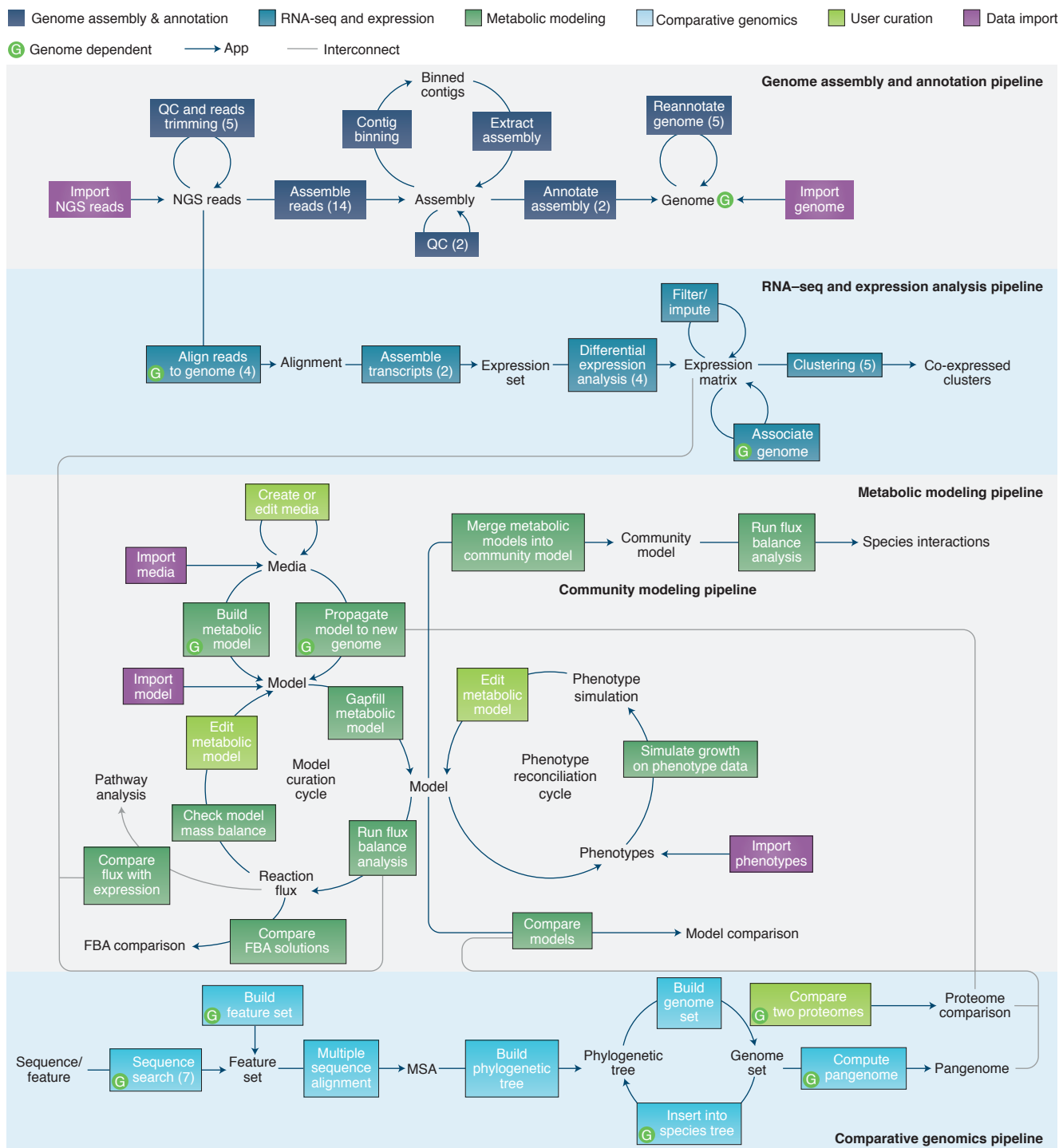


Figure 2 Major workflows and data types in KBase. The unboxed labels represent data types, while each colored box represents a single app. The box colors signify the category of functionality, and the numbers in parentheses indicate the number of alternative apps that implement each function. Apps that require a genome data type as input are marked with a green 'G' icon. For more information see <http://kbase.us/apps/>.

Through these public Narratives, scientists can rapidly follow the examples set by their peers to apply similar approaches to new data and scenarios. Thus KBase goes beyond supporting reproducible science to enabling rapid repurposing, reapplication, and extension of scientific techniques. As more users apply the

system to address their scientific questions, and share their resulting Narratives, KBase will have a continually growing body of experiments, results, and scientific use cases that can be adapted and extended by other researchers.

KBase's integration of data and tools and the ease of creating and running large-scale analy-

sis workflows have the potential to empower scientists in a broad range of application areas for systems biology, including environmental analysis, biosystems design, and human health. KBase's sharing capabilities amplify this potential by enabling scientists with differing expertise to easily work together and leverage

each other's work (**Supplementary Note 1; Supplementary Fig. 3**).

Future development of KBase will build upon the concept of KBase as a knowledgebase. The social aspects of the platform will be enhanced, enabling scientists to discover colleagues with complementary talents. New data-discovery features will allow the platform to suggest data sets and Narratives that may be of interest to a particular user, based on interconnections found in the data in KBase. These features will ultimately evolve into knowledge-discovery features, enabling KBase to propose new hypotheses by making connections across the system. Information on data and code availability can be found in **Supplementary Note 3**.

Note: Any Supplementary Information and Source Data files are available in the online version of the paper.

Editor's note: This article has been peer-reviewed.

ACKNOWLEDGMENTS

This work is supported by the Office of Biological and Environmental Research's Genomic Science program within the US Department of Energy Office of Science, under award numbers DE-AC02-05CH11231, DE-AC02-06CH11357, DE-AC05-00OR22725, and DE-AC02-98CH10886.

AUTHOR CONTRIBUTIONS

A.P.A., R.W.C., C.S.H., R.L.S., S.M., P.D., D.W., and F.P. developed the concept and vision. A.P.A., C.S.H., R.L.S., S.C., M.W.S., M.L.H., W.J.R., D.M.O., S.Y.C., T.S.B., D.C., D.G., J.B., A.A.B., B.P.B., S.E.B., C.C.B., J.M.C., J.C., R.C., N.C., J.J.D., M.D.J., S.D., A.G., F.H., M.P.J., K.P.K., F.M., P.S.N., R.O., E.P., S.P., G.A.P., S.R., P.R., S.M.D.S., M.S., R.A.S., M.H.S., J.T., F.X., H.Y., S.J.Y., and D.Y. designed and developed the system. R.W.C., N.L.H., R.T.K., S.K., M.M.D., E.M.G., D.C., D.J.W., B.H.A., B.H.D., E.D., I.D., J.N.E., G.F., J.P.F., P.M.F., W.G., M.G., J.G., R.J., S.N.K., V.K., M.L.L., M.M., T.Y.O., G.J.O., B.P., S.S.P., P.C.R., M.C.S., N.L.T., and D.F.W. developed, documented and conducted testing and validation. A.P.A., C.S.H., and N.L.H. drafted the manuscript. N.L.H., H.L.H., B.H.A., M.M.D., M.P.J., A.A.B., J.M.C., D.C., R.O., B.H.D., N.L.T., S.M., P.C.R., M.D.J., and V.K. revised the manuscript and provided important intellectual content. J.B., M.P.J., J.M.C., V.K., J.N.E., J.P.F., S.M.D.S. provided content for the supplementary information. A.P.A., R.W.C., C.S.H., and N.L.H. reviewed and approved the final version to be published.

COMPETING INTERESTS

F.P. declares competing interest related to his work for Plot.ly and research funding from Microsoft, Google, and Anaconda Inc. S.E.B. receives funding and has a research collaboration with Tata Consultancy Service that is unrelated to the KBase project.

Adam P Arkin^{1,2}, Robert W Cottingham³, Christopher S Henry⁴, Nomi L Harris², Rick L Stevens^{5,6}, Sergei Maslov^{7,24}, Paramvir Dehal², Doreen Ware⁸, Fernando Perez^{9-11,24}, Shane Canon¹², Michael W Sneddon², Matthew L Henderson², William J Riehl², Dan Murphy-Olson⁴, Stephen Y Chan², Roy T Kamimura², Sunita Kumari⁸, Meghan M Drake³, Thomas S Brettin⁶, Elizabeth M Glass⁴,

Dylan Chivian², Dan Gunter⁹, David J Weston³, Benjamin H Allen³, Jason Baumohl², Aaron A Best¹³, Ben Bowen², Steven E Brenner¹⁴, Christopher C Bun⁴, John-Marc Chandonia², Jer-Ming Chia⁸, Ric Colasanti⁴, Neal Conrad⁶, James J Davis⁶, Brian H Davison³, Matthew DeJongh¹⁵, Scott Devoid⁴, Emily Dietrich⁶, Inna Dubchak², Janaka N Edirisinghe^{4,16}, Gang Fang^{17,24}, José P Faria⁴, Paul M Frybarger⁴, Wolfgang Gerlach⁴, Mark Gerstein¹⁷, Annette Greiner¹², James Gurtowski⁸, Holly L Haun³, Fei He^{7,24}, Rashmi Jain^{18,19}, Marcin P Joachimiak², Kevin P Keegan⁴, Shinnosuke Kondo¹⁵, Vivek Kumar⁸, Miriam L Land³, Folker Meyer⁴, Marissa Mills³, Pavel S Novichkov², Taeyun Oh^{18,19,24}, Gary J Olsen²⁰, Robert Olson⁴, Bruce Parrello⁴, Shiran Pasternak⁸, Erik Pearson², Sarah S Poon⁹, Gavin A Price², Srividya Ramakrishnan^{8,24}, Priya Ranjan^{3,21}, Pamela C Ronald^{18,19}, Michael C Schatz^{8,24}, Samuel M D Seaver⁴, Maulik Shukla⁶, Roman A Sutormin², Mustafa H Syed^{3,24}, James Thomason⁸, Nathan L Tintle^{22,24}, Daifeng Wang^{17,24}, Fangfang Xia⁶, Hyunseung Yoo⁶, Shinjae Yoo²³ & Dantong Yu^{23,24}

¹Department of Bioengineering, University of California, Berkeley, California, USA.

²Environmental Genomics and Systems Biology Division, E.O. Lawrence Berkeley National Laboratory, Berkeley, California, USA.

³Biosciences Division, Oak Ridge National Laboratory, Oak Ridge, Tennessee, USA.

⁴Mathematics and Computer Science Division, Argonne National Laboratory, Argonne, Illinois, USA.

⁵Computer Science Department and Computation Institute, University of Chicago, Chicago, Illinois, USA.

⁶Computing, Environment, and Life Sciences Directorate, Argonne National Laboratory, Argonne, Illinois, USA.

⁷Biology Department, Brookhaven National Laboratory, Upton, New York, USA.

⁸Cold Spring Harbor Laboratory, Cold Spring Harbor, New York, USA.

⁹Computational Research Division, E.O. Lawrence Berkeley National Laboratory, Berkeley, California, USA.

¹⁰Berkeley Institute for Data Science, University of California, Berkeley, California, USA.

¹¹Department of Statistics, University of California, Berkeley, California, USA.

¹²National Energy Research Scientific Computing Center, E.O. Lawrence Berkeley National Laboratory, Berkeley, California, USA.

¹³Department of Biology, Hope College, Holland, Michigan, USA.

¹⁴Department of Plant and Microbial Biology, University of California, Berkeley, California, USA.

¹⁵Department of Computer Science, Hope College, Holland, Michigan, USA.

¹⁶Computation Institute, University of Chicago, Chicago, Illinois, USA.

¹⁷Program in Computational Biology and Bioinformatics, Yale University, New Haven, Connecticut, USA.

¹⁸Department of Plant Pathology and Genome Center, University of California, Davis, Davis, California, USA.

¹⁹Joint BioEnergy Institute, Lawrence Berkeley National Laboratory, Berkeley, California, USA.

²⁰Department of Microbiology, University of Illinois at Urbana-Champaign, Urbana, Illinois, USA.

²¹Department of Plant Sciences, University of Tennessee, Knoxville, Tennessee, USA. ²²Department of Mathematics, Hope College, Holland, Michigan, USA. ²³Computer Science and Math, Computer Science Initiative, Brookhaven National Laboratory, Upton, New York, USA. ²⁴Present addresses: Department of Bioengineering and Carl R. Woese Institute for Genomic Biology, University of Illinois at Urbana-Champaign, Urbana, Illinois, USA (S.M.); Department of Statistics, University of California, Berkeley, California, USA (F.P.); New York University Shanghai Campus, Pudong, Shanghai, China (G.F.); Department of Plant Pathology, Kansas State University, Manhattan, Kansas, USA (F.H.); Insilicogen, Inc., Giheung-gu, Yongin-si, Gyeonggi-do, Korea (T.O.); Department of Computer Science, Johns Hopkins University, Baltimore, Maryland, USA (S.R., M.C.S.); Memorial Sloan Kettering Cancer Center, New York, New York, USA (M.H.S.); Dordt College, Sioux Center, Iowa, USA (N.L.T.); Department of Biomedical Informatics, Stony Brook University, Stony Brook, New York, USA (D.W.); Martin Tuchman School of Management, New Jersey Institute of Technology, Newark, New Jersey, USA (D.Y.). Correspondence should be addressed to A.P.A. (aparkin@lbl.gov).



This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution

and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>

1. Stodden, V. *et al. Science* **354**, 1240–1241 (2016).
2. Prlić, A. & Procter, J.B. *PLoS Comput. Biol.* **8**, e1002802 (2012).
3. Millman, K.J. & Pérez, F. in *Developing Open-Source Scientific Practice*. Implementing reproducible research. (eds. F.L.V. Stodden & R.D. Peng) 149–183 (CRC Press, Boca Raton, FL; 2014).
4. Pruitt, K.D., Tatusova, T. & Maglott, D.R. *Nucleic Acids Res.* **35**, D61–D65 (2007).
5. Goodstein, D.M. *et al. Nucleic Acids Res.* **40**, D1178–D1186 (2012).
6. Bochner, B.R. *FEMS Microbiol. Rev.* **33**, 191–205 (2009).
7. Kanehisa, M. & Goto, S. *Nucleic Acids Res.* **28**, 27–30 (2000).
8. Schellenberger, J., Park, J.O., Conrad, T.M. & Palsson, B.O. *BMC Bioinformatics* **11**, 213 (2010).
9. Caspi, R. *et al. Nucleic Acids Res.* **34**, D511–D516 (2006).
10. Perez, F. & Granger, B.E. *Comput. Sci. Eng.* **9**, 21–29 (2007).
11. Kluyver, T. *et al.* in *Positioning and Power in Academic Publishing: Players, Agents and Agendas* (F. Loizides & B. Schmidt) 87–90 (IOS Press, Amsterdam, the Netherlands; 2016).
12. Edirisinghe, J.N. *et al. BMC Genomics* **17**, 568 (2016).
13. Magnúsdóttir, S. *et al. Nat. Biotechnol.* **35**, 81–89 (2017).
14. Henry, C.S. *et al. J. Cell. Physiol.* **231**, 2339–2345 (2016).
15. Faria, J.P. *et al. Front. Microbiol.* **7**, 1819 (2016).