

UC Berkeley

UC Berkeley Previously Published Works

Title

KBase: The United States Department of Energy Systems Biology Knowledgebase

Permalink

<https://escholarship.org/uc/item/7p39r37f>

Journal

Nature Biotechnology, 36(7)

ISSN

1087-0156

Authors

Arkin, Adam P
Cottingham, Robert W
Henry, Christopher S
[et al.](#)

Publication Date

2018-08-01

DOI

10.1038/nbt.4163

Peer reviewed

1 **KBase: The United States Department of Energy Systems Biology**
2 **Knowledgebase**

3
4

5 To the Editor:

6

7 Over the past two decades, the scale and complexity of genomics technologies and data have
8 advanced from sequencing genomes of a few organisms to generating metagenomes, genome
9 variation, gene expression, metabolites, and phenotype data for thousands of organisms and
10 their communities. A major challenge in this data-rich age of biology is integrating
11 heterogeneous and distributed data into predictive models of biological function, ranging from a
12 single gene to entire organisms and their ecologies. The US Department of Energy (DOE,
13 Washington, DC) has invested substantially in efforts to understand the complex interplay
14 between biological and abiotic processes that influence soil, water, and environmental dynamics
15 of our biosphere. The community that has grown around these efforts recognizes the need for
16 scientists of diverse backgrounds to have access to sophisticated computational tools that
17 enable them to analyze complex and heterogeneous data sets and integrate their data and
18 results effectively with the work of others. In this way, new data and conclusions could be
19 rapidly propagated across existing, related analyses and easily discovered by the community for
20 evaluation and comparison with previous results¹⁻³.

21

22 Here we present the DOE Systems Biology Knowledgebase (KBase, <http://kbase.us>), an
23 open-source software and data platform that enables data sharing, integration and analysis of
24 microbes, plants and their communities. KBase maintains an internal reference database that
25 consolidates information from widely used external data repositories. This includes over 90,000
26 microbial genomes from RefSeq⁴, over 50 plant genomes from Phytozome⁵, over 300 Biolog
27 media formulations⁶, and >30,000 reactions and compounds from KEGG⁷, BIGG⁸, and

28 MetaCyc⁹. These public data are available for integration with user data where appropriate (e.g.,
29 genome comparison or building species trees). KBase links these diverse data types with a
30 range of analytical functions within a web-based user interface. This extensive community
31 resource facilitates large-scale analyses on scalable computing infrastructure and has the
32 potential to accelerate scientific discovery, improve reproducibility and foster open collaboration.

33

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

51

52 KBase's primary user interface, the Narrative Interface, provides a user experience
53 distinct from other analysis platforms available today, although it shares some common features

54 with a few other systems (**Supplementary Note 2**). From this interface, which is built on the
55 Jupyter^{10, 11} platform, users can upload their private data, search and retrieve extensive public
56 reference data, access data shared by others, share their data with others, select and run
57 applications on their data, view and analyze the results from those applications and record their
58 thoughts and interpretations along with the analysis steps. These activities take place within a
59 point-and-click 'notebook' environment (Fig. 1). When a user begins a new computational
60 experiment in KBase, they create a new "notebook" (referred to as a 'Narrative' in KBase) to
61 hold this experiment. Every action performed by a user appears as a 'cell' in the Narrative. App
62 cells show the chosen input parameters for the application and the results of the analysis.
63 Markdown cells allow users to add formatted text and figures to a Narrative to describe the
64 thought process behind the scientific workflow being crafted.

65

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

77

78 The data model in KBase is fundamental to supporting reproducibility and collaboration.
79 KBase is built upon an object-oriented data model where each object instance is automatically

80 versioned and linked to provenance information describing how it was generated. Each data
81 object is also associated with the specific Narrative in which it was uploaded or generated.
82 When a Narrative is shared or copied, all its input and output data is shared or copied with it.
83 Currently supported data types include reads, contigs, genomes, metabolic models, growth
84 media, RNA-seq, expression, growth phenotype data, and flux balance analysis solutions. This
85 set of types can be extended to support new apps and functionality.

86
87 Many existing systems (**Supplementary Note 2**) provide similar support for object-level
88 sharing and provenance, but these systems operate on raw files only, without integration into a
89 common data model. In KBase, objects are not simple files—they are explicitly defined and
90 validated data structures, within which associated objects are linked to one another. For
91 example, a metabolic model is linked to its associated genome, which is linked to its associated
92 taxonomy. This data model enhances interoperability by requiring apps to operate on a common
93 data representation. Furthermore, it enhances awareness of interdependence so users could be
94 notified when an object on which an analysis is based has been updated and it will ultimately
95 enable data discovery and meta-analysis across the KBase platform.

96
97 Presently, KBase has over 160 apps (narrative.kbase.us/#appcatalog) offering diverse
98 scientific functionality for (meta)genome assembly, contig binning, genome annotation,
99 sequence homology analysis, tree building, comparative genomics, metabolic modeling,
100 community modeling, gap-filling, RNA-seq processing, and expression analysis (see
101 Supplementary Note 2 for references). Apps interoperate seamlessly to enable a range of
102 scientific workflows (**Fig. 2**). For reproducibility, all apps in KBase are containerized in versioned
103 Docker modules, enabling a user to run any version at any time.

104

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

111

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

117

118 KBase was designed to be an extensible community resource. This extensibility is
119 supported by the KBase Software Development Kit (SDK), which is a set of command-line tools
120 and a web interface that enable any developer to build, test, register, and deploy new or existing
121 software as KBase apps, thereby extending the platform's scientific capabilities. All software
122 contributed to the central KBase software repository must adhere to a standard open-source
123 license (opensource.org/licenses). Information about the app developer is maintained in the
124 documentation for that app so credit can be given to the contributor. Data provenance, job
125 management, usage logging, and app versioning are handled automatically by the platform,
126 allowing developers to wrap new scientific tools quickly with minimal KBase-specific training.
127 Other existing platforms offer similar support for third-party development (**Supplementary Note**
128 **2**), but KBase’s data model provides the additional benefit of improving interoperability of third-
129 party applications by imposing a single data format and specification on all data types

130 consumed or produced by each app. More information about the KBase SDK is available in the
131 supplement and at https://github.com/kbase/kb_sdk/blob/master/README.md.

132
133 Many users have already discovered and applied KBase to meet their scientific needs.
134 As of September 2017, over 3000 users have KBase accounts, and users have created over
135 5000 Narratives. These Narratives contain a total of over 250,000 data objects, or an average of
136 96 data objects and five apps per Narrative. Science done within KBase has been published in
137 over 30 peer-reviewed publications (**Supplementary Note 1**; <http://kbase.us/publications>),
138 including reconstruction of >8000 models of core metabolism across the microbial tree of life¹²;
139 reconstruction of semi-curated metabolic models for 773 gut microbes¹³; predicting trophic
140 interactions within a microbial community¹⁴; and reconstruction of regulons from expression
141 data¹⁵.

142
143 Much of the research performed within KBase has been publicly shared as Narratives that any
144 user can view, copy, and re-run. Through these public Narratives, scientists can rapidly follow
145 the examples set by their peers to apply similar approaches to new data and scenarios. Thus
146 KBase goes beyond supporting reproducible science to enable rapid re-purposing, re-
147 application, and extension of scientific techniques. As more users apply the system to address
148 their scientific questions, and share their resulting Narratives, KBase will have a continually
149 growing body of experiments, results and scientific use cases that can be adapted and
150 extended by other researchers.

151
152 KBase's integration of data and tools and the ease of creating and running large-scale
153 analysis workflows have the potential to empower scientists in a broad range of application
154 areas for systems biology, including environmental analysis, biosystems design, and human

155 health. KBase's sharing capabilities amplify this potential by enabling scientists with differing
 156 expertise to easily work together and leverage each other's work.

157
 158 Future development of KBase will build upon the concept of KBase as a knowledgebase.
 159 The social aspects of the platform will be enhanced, enabling scientists to discover colleagues
 160 with complementary talents. New 'data-discovery' features will allow the platform to suggest
 161 datasets and Narratives that may be of interest to a particular user based on interconnections
 162 found in the data in KBase. These features will ultimately evolve into 'knowledge-discovery'
 163 features, enabling KBase to propose new hypotheses by making connections across the
 164 system.

165
 166 Adam P Arkin^{1,2}, Robert W Cottingham³, Christopher S Henry⁴, Nomi L Harris², Rick L
 167 Stevens^{5,6}, Sergei Maslov^{7,24}, Paramvir Dehal², Doreen Ware⁸, Fernando Perez^{9-11,25}, Shane
 168 Canon¹², Michael W Sneddon², Matthew L Henderson², William J Riehl², Dan Murphy-Olson⁴,
 169 Stephen Y Chan², Roy T Kamimura², Sunita Kumari⁸, Meghan M Drake³, Thomas S Brettin⁶,
 170 Elizabeth M Glass⁴, Dylan Chivian², Dan Gunter⁹, David J Weston³, Benjamin H Allen³, Jason
 171 Baumohl², Aaron A Best¹³, Ben Bowen², Steven E Brenner¹⁴, Christopher C Bun⁴, John-Marc
 172 Chandonia², Jer-Ming Chia⁸, Ric Colasanti⁴, Neal Conrad⁶, James J Davis⁶, Brian H Davison³,
 173 Matthew DeJongh¹⁵, Scott Devoid⁴, Emily Dietrich⁶, Inna Dubchak², Janaka N Edirisinghe^{4,16},
 174 Gang Fang^{17,26}, José P Faria⁴, Paul M Frybarger⁴, Wolfgang Gerlach⁴, Mark Gerstein¹⁷, Annette
 175 Greiner¹², James Gurtowski⁸, Holly L Haun³, Fei He^{7,27}, Rashmi Jain^{18,19}, Marcin P Joachimiak²,
 176 Kevin P Keegan⁴, Shinnosuke Kondo¹⁵, Vivek Kumar⁸, Miriam L Land³, Folker Meyer⁴, Marissa
 177 Mills³, Pavel S Novichkov², Taeyun Oh^{18,19,28}, Gary J Olsen²⁰, Robert Olson⁴, Bruce Parrello⁴,
 178 Shiran Pasternak⁸, Erik Pearson², Sarah S Poon⁹, Gavin A Price², Srividya Ramakrishnan^{8,29},
 179 Priya Ranjan^{3,21}, Pamela C Ronald^{18,19}, Michael C Schatz^{8,29}, Samuel M D Seaver⁴, Maulik
 180 Shukla⁶, Roman A Sutormin², Mustafa H Syed^{3,30}, James Thomason⁸, Nathan L Tintle^{22,31},
 181 Daifeng Wang^{17,32}, Fangfang Xia⁶, Hyunseung Yoo⁶, Shinjae Yoo²³, Dantong Yu^{23,33}

182 183 **Affiliations**

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

185
 186 ²Environmental Genomics and Systems Biology Division, E. O. Lawrence Berkeley National
 187 Laboratory, Berkeley, California, USA.

188
 189 ³Biosciences Division, Oak Ridge National Laboratory, Oak Ridge, Tennessee, USA.

190

- 191 ⁴Mathematics and Computer Science Division, Argonne National Laboratory, Argonne, Illinois,
192 USA.
193
- 194 ⁵Computer Science Department and Computation Institute, University of Chicago, Chicago,
195 Illinois, USA.
196
- 197 ⁶Computing, Environment, and Life Sciences Directorate, Argonne National Laboratory,
198 Argonne, Illinois, USA.
199
- 200 ⁷Biology Department, Brookhaven National Laboratory, Upton, New York, USA.
201
- 202 ⁸Cold Spring Harbor Laboratory, Cold Spring Harbor, New York, USA.
203
- 204 ⁹Computational Research Division, E. O. Lawrence Berkeley National Laboratory, Berkeley,
205 California, USA.
206
- 207 ¹⁰Berkeley Institute for Data Science, University of California, Berkeley, California, USA.
208
- 209 ¹¹Department of Statistics, University of California, Berkeley, California, USA.
210
- 211 ¹²National Energy Research Scientific Computing Center, E. O. Lawrence Berkeley National
212 Laboratory, Berkeley, California, USA.
213
- 214 ¹³Department of Biology, Hope College, Holland, Michigan, USA.
215
- 216 ¹⁴Department of Plant and Microbial Biology, University of California, Berkeley, California, USA.
217
- 218 ¹⁵Department of Computer Science, Hope College, Holland, Michigan, USA.
219
- 220 ¹⁶Computation Institute, University of Chicago, Chicago, Illinois, USA.
221
- 222 ¹⁷Program in Computational Biology and Bioinformatics, Yale University, New Haven,
223 Connecticut, USA.
224
- 225 ¹⁸Department of Plant Pathology and Genome Center, University of California, Davis, Davis
226 California, USA.
227
- 228 ¹⁹Joint BioEnergy Institute, Lawrence Berkeley National Laboratory, Berkeley, California, USA.
229
- 230 ²⁰Department of Microbiology, University of Illinois at Urbana-Champaign, Urbana, Illinois, USA.
231
- 232 ²¹Department of Plant Sciences, University of Tennessee, Knoxville, Tennessee, USA.
233
- 234 ²²Department of Mathematics, Hope College, Holland, Michigan, USA.

235

236 ²³Computer Science and Math, Computer Science Initiative, Brookhaven National Laboratory,
237 Upton, New York, USA.

238

239 **Present Addresses**

240 ²⁴Department of Bioengineering and Carl R. Woese Institute for Genomic Biology, University of
241 Illinois at Urbana-Champaign, Urbana, Illinois, USA.

242

243 ²⁵Department of Statistics, University of California, Berkeley, California, USA.

244

245 ²⁶New York University Shanghai Campus, Pudong, Shanghai, China.

246

247 ²⁷Department of Plant Pathology, Kansas State University, Manhattan, Kansas, USA.

248

249 ²⁸Insilicogen. Inc., Giheung-gu, Yongin-si, Gyeonggi-do, Korea.

250

251 ²⁹Department of Computer Science, Johns Hopkins University, Baltimore, MD, USA.

252

253 ³⁰Memorial Sloan Kettering Cancer Center, New York, New York, USA.

254

255 ³¹Dordt College, Sioux Center, Iowa, USA.

256

257 ³²Department of Biomedical Informatics, Stony Brook University, Stony Brook, New York, USA.

258

259 ³³Martin Tuchman School of Management, New Jersey Institute of Technology, Newark, New
260 Jersey, USA.

261

262 **Corresponding author**

263 Correspondence and requests for materials should be addressed to: AParkin@lbl.gov.

264

265

266 **Code Availability**

267 The KBase code, available at github.com/kbase, is open source and freely distributed under the

268 MIT License. The web-accessible KBase system (narrative.kbase.us) is run on DOE computing

269 infrastructure and is freely available for anyone to use. KBase adheres to the FAIR (Findable,

270 Accessible, Interoperable, Re-usable) data principles endorsed by many funding agencies and
271 scientific organizations¹⁶.

272

273 **Data Availability**

274 All data generated or analyzed during this study are included in this published article and
275 Supplementary Note 1 as links to the original work, or in the associated KBase Narratives linked
276 here. An earlier version of this paper was published as a preprint¹⁷.

277 **Acknowledgements**

278 *This work is supported by the Office of Biological and Environmental Research's Genomic*
279 *Science program within the U.S. Department of Energy Office of Science, under award numbers*
280 *DE-AC02-05CH11231, DE-AC02-06CH11357, DE-AC05-00OR22725, and DE-AC02-*
281 *98CH10886.*

282

283 Editors Note: This manuscript has been peer reviewed.

284 **Competing financial interests**

285 FP declares competing financial interest related to his work for Plot.ly and research funding from
286 Microsoft, Google, and Anaconda Inc.

287 SEB receives funding and has a research collaboration with Tata Consultancy Service that is
288 unrelated to the KBase project.

289 All other authors declare no competing financial interests.

290

291 **Author Contributions**

292 APA, RWC, CSH, RLS, SM, PD, DW and FP developed the concept and vision.

293

294 APA, CSH, RLS, SC, MWS, MLH, WJR, DMO, SYC, TSB, DC, DG, JB, AAB, BPB, SEB, CCB,
 295 JMC, JC, RC, NC, JJD, MDJ, SD, AG, FH, MPJ, KPK, FM, PSN, RO, EP, SP, GAP, SR, PR,
 296 SMDS, MS, RAS, MHS, JT, FX, HY, SJY and DY designed and developed the system.

297
 298 RWC, NLH, RTK, SK, MMD, EMG, DC, DJW, BHA, BHD, ED, ID, JNE, GF, JPF, PMF, WG,
 299 MG, JG, RJ, SNK, VK, MLL, MM, TYO, GJO, BP, SSP, PCR, MCS, NLT and DFW developed,
 300 documented and conducted testing and validation.

301
 302 APA, CSH and NLH drafted the manuscript.

303
 304 NLH, HLH, BHA, MMD, MPJ, AAB, JMC, DC, RO, BHD, NLT, SM, PCR, MDJ and VK revised
 305 the manuscript and provided important intellectual content.

306
 307 JB, MPJ, JMC, VK, JNE, JPF, SMDS provided content for the supplemental material.

308
 309 APA, RWC CSH and NLH reviewed and approved the final version to be published.

310

311

312

313 References

314

- 315 1. Stodden, V. et al. Enhancing reproducibility for computational methods. *Science* **354**,
 316 1240-1241 (2016).
- 317 2. Prlic, A. & Procter, J.B. Ten simple rules for the open development of scientific software.
 318 *PLoS Comput Biol* **8**, e1002802 (2012).
- 319 3. Millman, K.J. & Pérez, F. in *Developing Open-Source Scientific Practice. Implementing*
 320 *reproducible research*. (eds. F.L.V. Stodden & R.D. Peng) 149-183 (CRC Press, Boca
 321 Raton, FL; 2014).
- 322 4. Pruitt, K.D., Tatusova, T. & Maglott, D.R. NCBI reference sequences (RefSeq): a
 323 curated non-redundant sequence database of genomes, transcripts and proteins.
 324 *Nucleic Acids Res* **35**, D61-65 (2007).
- 325 5. Goodstein, D.M. et al. Phytozome: a comparative platform for green plant genomics.
 326 *Nucleic Acids Res* **40**, D1178-1186 (2012).
- 327 6. Bochner, B.R. Global phenotypic characterization of bacteria. *Fems Microbiology*
 328 *Reviews* **33**, 191-205 (2009).
- 329 7. Kanehisa, M. & Goto, S. KEGG: Kyoto encyclopedia of genes and genomes. *Nucleic*
 330 *Acids Research* **28**, 27-30 (2000).
- 331 8. Schellenberger, J., Park, J.O., Conrad, T.M. & Palsson, B.O. BiGG: a Biochemical
 332 Genetic and Genomic knowledgebase of large scale metabolic reconstructions. *BMC*
 333 *Bioinformatics* **11**, 213.
- 334 9. Caspi, R. et al. MetaCyc: a multiorganism database of metabolic pathways and
 335 enzymes. *Nucleic Acids Research* **34**, D511-D516 (2006).

- 336 10. Perez, F. & Granger, B.E. IPython: A system for interactive scientific computing.
 337 *Computing in Science & Engineering* **9**, 21-29 (2007).
- 338 11. Kluyver, T. et al. Jupyter Notebooks—a publishing format for reproducible computational
 339 workflows. *Positioning and Power in Academic Publishing: Players, Agents and*
 340 *Agendas*, 87-90 (2016).
- 341 12. Edirisinghe, J.N. et al. Modeling central metabolism and energy biosynthesis across
 342 microbial life. *BMC Genomics* **17**, 568 (2016).
- 343 13. Magnusdottir, S. et al. Generation of genome-scale metabolic reconstructions for 773
 344 members of the human gut microbiota. *Nat Biotechnol* **35**, 81-89 (2017).
- 345 14. Henry, C.S. et al. Microbial community metabolic modeling: A community data-driven
 346 network reconstruction. *J Cell Physiol* (2016).
- 347 15. Faria, J.P. et al. Computing and Applying Atomic Regulons to Understand Gene
 348 Expression and Regulation. *Front Microbiol* **7**, 1819 (2016).
- 349 16. Wilkinson, M.D. et al. The FAIR Guiding Principles for scientific data management and
 350 stewardship. *Sci Data* **3**, 160018 (2016).
- 351 17. Arkin, A.P. et al. The DOE Systems Biology Knowledgebase (KBase). *bioRxiv preprint*
 352 **first posted online Dec. 22, 2016** (2016).
- 353
- 354
- 355

356

357 Figure Legends

358

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

361

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

367