# Supplementary Notes

## Supplementary Note 1. Examples of Use of KBase Infrastructure and Tools

The KBase platform has been utilized by a variety of scientific groups, with citations in over 30 peer-reviewed publications. These publications cover a range of topics and demonstrate the novel ways scientists are applying tools within KBase to their research. Several of these studies have publicly shared KBase Narratives associated with them. These serve as a mechanism to show how the KBase platform was applied to perform the described work as well as a means of releasing the data produced in the published analysis. The KBase website maintains a curated list of published Narratives at http://www.kbase.us/narrative-library/. Here we will review a selection of these research Narratives that demonstrate many (but not all) of the workflows and capabilities available in the KBase platform (see Figure 2 in the main manuscript).

## 1A. Reconstruction of 8000 Models of Core Metabolism Across the Microbial Tree of Life

Link: https://narrative.kbase.us/narrative/ws.20186.obj.18

In a recent publication in *BMC Genomics*[1] researchers developed a framework in KBase to reconstruct and analyze core metabolic models (CMMs). Core metabolic models representing 48 major phylogenetic microbial groups were constructed based on a core model template consisting of a highly curated set of biochemical reactions derived from a diverse set of model organisms, and about 200 unique reactions were selected from this set, comprising 12 key energy biosynthesis pathways linked to central metabolism and variations of bacterial electron transport chains. This framework was applied in KBase to build core models for over 8000 prokaryotic genomes that span the prokaryotic tree of life. The authors used CMMs to determine: (i) accurate ATP yields based on different growth/environmental conditions; (ii) ETC variations and respiration types; (iii) ability to produce fermentation products; (iv) presence and absence of classical biochemical pathways in central metabolism; and (v) ability to produce key metabolic pathway intermediates in central metabolism which are precursors of essential biomass components of the cell.

The Narrative associated with this work utilizes the core model reconstruction, gapfilling, and flux balance analysis (FBA) workflows in KBase. It also demonstrates the capacity of KBase to perform large-scale analyses. The authors used a *code cell* to programmatically run the model reconstruction, gapfilling, and FBA workflow on over 8000 genomes (Figure 1). Code cells make it possible to run any KBase App programmatically from within a Narrative. One can write a code cell in the Python programming language within a Narrative to run a single app or an entire pipeline of apps. While code cells require knowledge of Python programming to write, Narratives with useful code cells can be shared with any user, who can re-run the code cells with ease. This work also inspired the addition of a new batch app for model reconstruction called *Build*

42     *Multiple Metabolic Models.* With this new app, it is no longer necessary to use a code cell to
43     build thousands of models. Such batch apps demonstrate a further mechanism beyond code
44     cells through which the KBase platform can be used to conduct large studies. To download the
45     8000+ models from KBase in SBML form, the authors used the *Bulk Download Modeling*
46     *Objects* app, which permits bulk download of any supported type modeling object.

## 1B. Identifying Violacein Synthesis Genes in a New Isolate Genome

48     Link: https://narrative.kbase.us/narrative/ws.21546.obj.1
49
50     With its many apps for microbial genome assembly and annotation, KBase is an ideal platform
51     for the analysis of new isolate genomes. This is demonstrated in a recent *Genome*
52     *Announcement* publication by Romy Chakraborty and colleagues[2]. In this study, groundwater
53     samples from the ORNL FRC site were collected from multiple wells, where researchers
54     identified a *Janthinobacterium* isolate that produces violacein. Violacein is a naturally-occurring
55     bis-indole pigment with antibiotic (antibacterial, anti-fungal and anti-tumor) properties. The
56     Chakraborty team loaded their isolate reads into KBase, performing QC and assembly. They
57     next applied both the Prokka[3] and RAST[4] annotation tools to the assembled genome, with
58     Prokka calling 5531 genes, while RAST called 5998. They also ran OrthoMCL[5] to compare the
59     two alternative annotations side by side. From this analysis, the researchers determined that
60     both RAST and Prokka properly annotated 4 out of 5 genes in the violacein biosynthesis
61     pathway. However, these algorithms differed in the specific violacein gene they missed. RAST
62     missed VioC, while Prokka missed VioE. By combining both sets of annotations, the
63     researchers were able to identify and confirm the presence of all violacein genes in their isolate
64     of interest. The Chakraborty team's analysis is captured in the linked Narrative, including their
65     thought process at each step. This study demonstrates the value of having multiple apps for
66     some of the key steps of the isolate analysis pipeline. During genome assembly and genome
67     annotation, the researchers were able to apply multiple algorithms, evaluating and comparing
68     the results, and selecting the best-performing algorithm for each step.

## 1C. Predicting Trophic Interactions Within a Microbial Community

70     Narrative 1: https://narrative.kbase.us/narrative/ws.13807.obj.1
71     Narrative 2: https://narrative.kbase.us/narrative/ws.13806.obj.1
72     Narrative 3: https://narrative.kbase.us/narrative/ws.13838.obj.1
73
74     KBase enables researchers to study interactions between multiple species in a microbial
75     community. In three interconnected Narratives associated with a recent cover article in the
76     *Journal of Cellular Physiology*[6], the authors applied KBase to construct new metabolic models
77     for a photoautotroph, *Thermosynechococcus elongatus* BP1, and a heterotroph, *Meiothermus*
78     *ruber* strain A. They used these models to predict trophic interactions between species,
79     exploring a wide range of potential methods for predicting these interactions. The authors
80     utilized metatranscriptomic data to validate predictions and identify which prediction method
81     performed the best (see Figure 2).

82

83  These Narratives collectively demonstrate a large fraction of the capabilities presently available
84  in KBase, including: (i) proteome comparison; (ii) building species trees; (iii) annotating a
85  genome; (iv) building and gapfilling a modeling; (v) propagating a model; (vi) loading expression
86  data; (v) building a community metabolic model; (vi) running flux balance analysis; (vii)
87  comparing flux with expression data (Figure 2); and (viii) using code cells to sift data from
88  objects in the KBase workspace. In this case, the code cell is used to print out the trophic
89  interactions among species in a community model based on the flux profile generated by the
90  FBA app.

91

## 1D. Modeling Metabolic Interactions Between Cyanobacteria-Sphagnum

94  Link: https://narrative.kbase.us/narrative/ws.9667.obj.2

95

96  KBase has tools and datasets that are useful for analyzing plant metabolism and plant/microbe
97  interactions. In a study published in *Plant, Cell & Environment*[7], researchers explored a mutually
98  beneficial metabolic partnership between a moss and a bacterium, using KBase data and tools
99  to build a merged community metabolic model in which a nitrogen-fixing diazotrophic microbe
100 (Anabena) fixes nitrogen to allow a plant (Sphagnum) to grow. Sphagnum (peat moss), a genus
101 of plants (Bryophyta) that associate with nitrogen-fixing diazotrophs is a quintessential
102 ecosystem engineer. All the analysis steps were performed in KBase, culminating with merging
103 the two models into a community model which exhibits nitrogen fixation and exchange, showing
104 that the plant portion of the model consumes the nitrogen fixed by the microbial portion of the
105 model, and predicting that the Sphagnum will grow more when utilizing nitrogen fixed by the
106 microbe than when fixing nitrates on its own.

107

108 This Narrative demonstrates the plant reference genomes in KBase, as well as the plant model
109 reconstruction pipeline. It also shows how models of any type (for example, microbial and plant
110 models) can be merged together to form a community model.

## 1E. Sharing Workflows Easily with Collaborators and the General Public

113 Link 1: https://narrative.kbase.us/narrative/ws.18152.obj.1
114 Link 2: https://narrative.kbase.us/narrative/ws.18153.obj.1
115 Link 3: https://narrative.kbase.us/narrative/ws.18155.obj.1
116 Link 4: https://narrative.kbase.us/narrative/ws.18156.obj.1
117 Link 5: https://narrative.kbase.us/narrative/ws.18157.obj.1

118

119 These Narratives were specifically crafted to demonstrate how the KBase platform can facilitate
120 interaction, data sharing, and collaboration between two scientists working on a common
121 problem. The ability to share data, results, and workflows is extremely important and more or
122 less a requirement in collaborative studies[8]. KBase facilitates copying and sharing of

123     collaborative data, workflows and related commentary. A user can share any Narrative that they
124     own with other KBase users or make it publically available. Importantly, when a user shares a
125     Narrative, they also are sharing all the data objects loaded, used, or generated within the
126     Narrative, complete with versioning and provenance. Copying and sharing are performed in a
127     controlled manner where users with read privileges for a Narrative can create their own copy,
128     which they own and can edit. Users are able to quickly replicate and expand on any KBase
129     Narrative shared with them. This approach facilitates reproducible interdisciplinary science by
130     allowing researchers with different expertise to quickly and easily exchange data, results,
131     methodologies, and workflows to address complex biological problems.
132
133     We demonstrate collaborative sharing of data, commentary (notes) and workflows in KBase
134     using a series of example Narratives (Figure 3) featuring two hypothetical scientists: Alice, an
135     experimental biologist with expertise in assembly, annotation, and comparative genomics, and
136     Bob, a computational biologist with expertise in metabolic modeling.
137
138     Alice the experimentalist uploads raw reads that she has sequenced from a strain of
139     *Shewanella* on which she wants to perform comparative genomic analysis in order to
140     understand the similarities and unique features compared to phylogenetically closely related
141     genomes.  She uploads sequencing reads to KBase and then assembles and annotates the
142     reads generating an annotated genome (Alice assembly and annotation Narrative:
143     https://narrative.kbase.us/narrative/ws.18152.obj.1). In a separate Narrative (Alice Comparative
144     Genomics Narrative: https://narrative.kbase.us/narrative/ws.18153.obj.1), Alice identifies
145     genomes that are phylogenetically close to her *Shewanella* strain. She also finds growth
146     phenotype data for *Shewanella oneidensis* MR-1, which is phylogenetically close to her strain.
147     This inspires Alice to perform a growth phenotype array on her own strain. Alice then compares
148     phenotype arrays of her own strain vs *S.oneidensis* MR-1 and notices many differences that she
149     cannot explain.
150
151     In order to understand more about Biolog phenotype data, Alice contacts Bob the modeler who
152     is able to analyze the metabolic differences between the two strains, which in turn helps Alice to
153     interpret the Biolog phenotype data. Alice shares her Narratives with Bob, which allows Bob to
154     copy her genomes into a new Narrative. In this Narrative (Bob build metabolic models Narrative:
155     https://narrative.kbase.us/narrative/ws.18155.obj.1) Bob loads a published model of *S.
156 oneidensis* MR-1, which he then propagates to Alice's genome, producing a high quality
157     metabolic model using the published model as a template. Bob compares the two models,
158     identifying some interesting metabolic differences. Then Bob creates a separate Narrative (Bob
159     and Alice shared Narrative of Phenotype Data Analysis:
160     https://narrative.kbase.us/narrative/ws.18156.obj.1) in which he imports Alice's Biolog data and
161     simulates the data with his *Shewanella* models. He optimizes his models to fit the Biolog data
162     and shares the results with Alice. Finally, the two scientists build another shared Narrative (Bob
163     and Alice shared Narrative of Phenotype Data Reconciliation:
164     https://narrative.kbase.us/narrative/ws.18157.obj.1) in which Alice improves the quality of Bob's
165     models by curating the models further by replacing some of the gapfilled reactions with more

166  biologically meaningful selections, gaining a complete understanding of the differences between
167  her strain and MR-1.
168
169  This example demonstrates how KBase can facilitate a collaborative study between
170  scientists/groups with different but complementary expertise who are able to accomplish more
171  together than they could individually. By saving their data, workflows along with commentary as
172  Narratives, scientists who use KBase also enable other researchers to quickly reproduce their
173  work with a minimum effort and also aid in further extending similar scientific research.

174  ## Supplementary Note 2. Comparison of KBase with Other Online
175  ## Systems Biology Resources

176  KBase has over 160 apps offering diverse scientific functionality for (meta)genome assembly,
177  contig binning[9], genome annotation[10], sequence homology analysis[5], tree building[11],
178  comparative genomics, metabolic modeling[12], community modeling[13], gap-filling[14, 15], RNA-seq
179  processing[16], and expression analysis[17]. Apps in KBase interoperate seamlessly to enable a
180  range of scientific workflows. Numerous tools exist today that are similar to KBase in that they
181  offer web-based access to a variety of systems biology workflows. Five of the most similar and
182  widely used tools are Galaxy[18], CyVerse[19], Pathway Tools[20], BaseSpace
183  (http://basespace.illumina.com) and GenePattern[21]. In order to highlight the ways in which
184  KBase is different from each of these frameworks, we conducted a detailed comparison
185  between KBase and these tools, focusing on the areas of functionality where KBase is most
186  distinctive. Here we report on the results of our comparison, organizing our analysis into six
187  fundamental areas: (a) user experience; (b) data model and provenance; (c) built in reference
188  data; (d) sharing of data and workflows; (e) third-party development and custom code support;
189  and (f) available scientific functionality.

190  ## 2A. Overview of All Compared Platforms

191  We begin by providing some background on each of the resources selected for our comparison.
192  The first framework we selected was Galaxy, which is a scientific workflow, data integration, and
193  analysis platform that aims to make bioinformatics and computational biology workflows
194  accessible to researchers without computer programming or systems administration experience.
195  Like KBase, Galaxy is an open-source project. The project began in 2005, and it currently exists
196  as a collaborative effort involving primarily Pennsylvania State University, Johns Hopkins
197  University, and Oregon Health & Science University. Although Galaxy is available in many forms
198  (described in detail at https://galaxyproject.org/citing-galaxy/), we focus our comparison on the
199  version of Galaxy that most resembles KBase, which is the main public usegalaxy.org site.
200
201  The second framework we selected for comparison was CyVerse. CyVerse was launched by
202  NSF in 2008 initially as the iPlant Collaborative, but in 2013,  its scope was extended to all non-
203  human life sciences, with a renaming to CyVerse in 2016. The primary mission of CyVerse is to
204  transform broader science through data discovery. While CyVerse offers customizable and
205  programmatic interfaces, its two primary interfaces are a point and click web interface called the

206  Discovery Environment, and a cloud-hosted VM app based interface called Atmosphere. As with
207  Galaxy, we will focus most of our comparison on the web-based Discovery Environment. The
208  CyVerse Discovery Environment provides functionality to manage data, add new algorithms and
209  tools, and run analyses on appropriate computational resources. It also provides access to
210  storage (user's own, shared, and public data).
211
212  The third framework we selected for comparison was the BaseSpace Sequence Hub.
213  BaseSpace is a cloud-based, closed-source platform that provides tools to manage and analyze
214  Illumina sequence data.  Signing up for a free trial account allows users to store up to 1TB of
215  sequence data for an indefinite period of time, as well as giving them a limited number of
216  "compute credits". Many apps require "compute credits" to run, but free apps are also available.
217  Like CyVerse and Galaxy, BaseSpace offers a web interface for storing and analyzing biological
218  data, although this framework focuses primarily on sequencing data.
219
220  The fourth framework we selected for comparison was GenePattern. GenePattern provides
221  hundreds of analytical tools for gene expression (RNA-seq and microarray), sequence variation,
222  proteomic data, flow cytometry, and network analysis. These tools are all made available
223  through the online Gene Pattern Notebook environment with no programming experience
224  required. Like KBase, this notebook environment extends the Jupyter Notebook[22, 23] system,
225  allowing researchers to create documents that interleave formatted text, graphics and other
226  multimedia, executable code, and GenePattern analyses. GenePattern was developed at the
227  Broad Institute starting in 2006, and it is primarily funded by the National Cancer Institute.
228
229  Given the extensive metabolic modeling tools available in KBase, we wanted to include a
230  metabolic modeling framework in our comparison. For this, we chose Pathway Tools. Pathway
231  Tools was developed as a means of curating and visualizing biochemical pathways and
232  associated data for various organisms, dating all the way back to 1993. It is one of the longest
233  lived software/database/web-server suites available. Pathway Tools supports the use of all
234  known metabolism, in a wide taxonomical range, allowing users to create databases for
235  bacteria, fungi, and plants. Pathway Tools integrates a web server, enabling laboratories to host
236  their own instance of the database and publish their own data. Pathway Tools is somewhat
237  different from the other frameworks selected for comparison in that most functionalities are only
238  available offline through its installable software suite. Still, it is useful to include it in our
239  comparison as it embodies many distinctive design patterns.


240  ## 2B. Comparison of User Experience

241  KBase and all five of the platforms we selected for comparison offer a similar overall user
242  experience. All platforms have a graphical user interface, which enables users to view data and
243  run apps in point-and-click fashion. All platforms offer a centralized website from which users
244  can access a canonical version of the platform. Galaxy, GenePattern, and Pathway Tools also
245  offer a downloadable version, which can be installed and run on a user's own hardware. The
246  KBase SDK also allows for this, but in a much less user-friendly and more limited manner.
247

248  KBase's Narrative Interface is distinct from any other analysis platforms available today,
249  although it shares some common features with GenePattern[21]. KBase and GenePattern are
250  both built on the Jupyter platform, allowing users to fashion multi-step analyses within online
251  notebooks which they can then share. Both KBase and GenePattern extend Jupyter by offering
252  users a point-and-click menu of apps that can be run within the Jupyter notebook. However,
253  KBase is unique in also wrapping a data-layer around Jupyter, enabling users to browse and
254  view data objects imported or generated by apps within the notebook. Pathway Tools, CyVerse,
255  and BaseSpace are all app-centric interfaces with limited workflow support. Users select a
256  single app, then select data to run the app on. Then the user views and shares the output of the
257  analysis. Other point-and-click computational platforms do exist that enable users to
258  dynamically construct workflows, including Taverna[24], XSEDE[25], myExperiment[26], Kepler[27],
259  Pegasus[28], and Globus[29], but many of these platforms lack KBase's tight integration of tools and
260  data, and none of these platforms offers the "story-telling" capacity of a Jupyter notebook-based
261  interface.
262
263  Thus in the area of user experience, Galaxy and GenePattern are the most similar to KBase.


264  ## 2C. Comparison of Data Model

265  Data model is one area where the platforms we selected for comparison vary the most. In our
266  comparison, we found three distinct design patterns: (i) a file-based design (Galaxy, CyVerse,
267  BaseSpace, GenePattern); (ii) a structured object-based design (KBase); and (iii) an entirely
268  relational design (Pathway Tools).
269
270  Galaxy, CyVerse, GenePattern and BaseSpace all share very similar data frameworks in which
271  all data is stored in the native file format uploaded by the user (e.g., FASTA, FASTQ, BAM). In
272  these frameworks, the files are always augmented with associated metadata. This data model
273  has several advantages: (i) upload and download are easy because there is no need to
274  transform files into another form; (ii) sharing, provenance, and versioning are all easy because
275  individual objects are self-contained and these features can operate on the level of each
276  individual object; (iii) integration of files with tools is simple because tools typically operate on
277  native file formats directly. However, this approach also has disadvantages: (i) files are only as
278  consistent as their standards force them to be, and many files types (e.g., SBML, FASTA, GFF,
279  GenBank) actually involve extensive variability in how they represent data; (ii) there are often
280  many different file formats representing a single entity (e.g., FASTQ, FASTA, SRA for reads),
281  meaning many file format conversion utilities are required and the user spends significant time
282  transforming file formats; (iii) complex files (e.g., GenBank, FASTQ, SRA) are often treated as
283  black boxes, meaning they lack introspection and the files are not indexed in detail; and (iv) the
284  data that can be stored in a data type is limited by what is accommodated by its associated file
285  formats (e.g., in COBRA, extensive data is stashed in the typed object used to represent a
286  model, but when that typed object is converted to SBML, some of that data is lost due to
287  limitations in the SBML file format).
288

289  The KBase data model is similar to the file-based data model in that it is also object-based. This
290  means the KBase data model shares the largest benefit of the file-based system, which is that
291  individual objects are self-contained and can be independently shared, versioned, and
292  provenanced. However, instead of representing objects in their original file formats, KBase
293  represents objects in a single standardized, typed, rigorously specified, versioned, and validated
294  JSON-based format. The downside of this approach is that all uploaded and downloaded files
295  must be converted to and from their associated KBase object type, which makes supporting
296  upload and download a challenge for tool developers (the conversion process is transparent
297  and thus has very limited impact on end-users). This conversion can be lossy if the input file and
298  output data-type are not completely synchronized. This also makes tool integration more of a
299  challenge, as the tool developers must add some additional code to handle the conversion to
300  and from any KBase data type that their tool operates against. A developer may also need to
301  add a new data type if the output of the tool must be persistent (and reused) and the type
302  doesn't already exist. However, this approach eliminates nearly all the downsides of a file-based
303  system: (i) types are totally standardized and consistent even if their associated files are not; (ii)
304  there is only one representation for each fundamental type (e.g., "reads" vs SRA, FASTQ,
305  FASTA); (iii) types can be summarized, viewed, introspected and interconnected, although this
306  still isn't as easy and performant as a fully relational data model; and (iv) data types can be
307  rapidly expanded as required to meet the demands of new analyses being added to the
308  framework (e.g., adding atom mappings to metabolic models).
309
310  The data model in Pathway Tools is completely distinct from the other platforms, maintaining
311  data internally within a relational database, while using primarily custom flat-file-formats that are
312  specific to Pathway Tools for data exchange. This has the advantage of maintaining data in a
313  highly interconnected and queryable format within the Pathway Tools framework. Introspection,
314  search, standardization, and internal consistency are all great strengths of this approach.
315  However, this comes at the cost of granularity in support for versioning, sharing, and
316  provenance in Pathway Tools.


## 2D. Comparison of Provenance, Data Sharing, and Data Versioning

318  Our data model comparison reveals how the provenance, data sharing, and data versioning
319  features of our selected frameworks depend significantly on the data model of the framework.
320  Galaxy, CyVerse, BaseSpace, GenePattern and KBase, all of which have file-based or object-
321  based data models, have similar support for provenance and sharing. All of these frameworks
322  maintain information about the input parameters and apps used to produce each object stored
323  in the system (limited in the case of GenePattern). All of these frameworks also support sharing
324  at the level of individual objects. However, this always involves copying objects into "libraries"
325  (Galaxy), "folders" (CyVerse), "projects" (BaseSpace), "notebooks" (GenePattern), or
326  "Narratives" (KBase) and sharing those containers with other users. Galaxy, CyVerse,
327  BaseSpace, GenePattern, and KBase do vary in how they handle versioning. Only KBase
328  maintains an explicit version number of every overwritten object, and only KBase allows an
329  overwritten object to be reverted to a previous version. BaseSpace and GenePattern do not
330  support data versioning, instead relying on users not to overwrite data that they want to keep

331  (although this can result in loss of the data needed to repeat a particular downstream result).
332  CyVerse and Galaxy follow a very different approach in that objects are namespaced and
333  timestamped according to the workflow and the previous objects from which they were derived.
334  Thus, it is essentially impossible to overwrite an object, and all versions of all data are
335  preserved unless explicitly deleted by the user.
336
337  In this area, Pathway Tools is quite different from the other frameworks, which is not surprising
338  given its very different data model. Versioning in Pathway Tools only occurs at the level of an
339  entire PathwayDB, and sharing is only supported by exporting data into files, which can then be
340  shared offline and imported elsewhere. Pathway Tools only preserves provenance in the
341  annotations and curations made to its underlying relational database, although in this case,
342  provenance can be maintained with a granularity and detail that exceeds the other systems
343  discussed.


344  ## 2E. Comparison of Built-in Reference Data

345  Reference data is another area where KBase stands out among the platforms we selected for
346  comparison. Reference data is vital in that it serves to place user data and analysis results into
347  the broader context of all other known data of the same type. It is often more useful to
348  understand how a genome is different from its phylogenetically close neighbors than to
349  understand every single detail of the genome itself.
350
351  Among all our platforms selected for comparison, only KBase and Pathway Tools offer their own
352  internally managed, organized, and curated reference data collections. KBase offers a reference
353  database of over 90K microbial and eukaryotic genomes maintained and periodically
354  synchronized with RefSeq and Phytozome. To facilitate comparison, all of the microbial
355  genomes have been annotated using the RAST genome annotation app in KBase. KBase also
356  offers a reference database of biochemistry, including 27K compounds, 34K reactions, and 522
357  media formulations. The reference genomes in KBase are readily available for copying into any
358  KBase Narrative for analysis and comparison, but they also form the basis for some apps that
359  analyze user data in the context of this reference data (e.g., the *Insert Genomes into Species*
360  *Tree* app). The reference biochemistry in KBase forms the basis for the metabolic model
361  reconstruction, standardization, and gapfilling tools in KBase.
362
363  Pathway Tools offers ready access to a database of 9318 Pathway Genome Databases, each
364  of which represents a single genome, a metabolic reconstruction, and a basic model. Pathway
365  Tools also deeply integrates a reference biochemistry database comprised of 14K reactions and
366  13K compounds. Even more than KBase, this reference data in Pathway Tools is at the core of
367  every analysis a user does in the platform.
368
369  While Galaxy, GenePattern, CyVerse, and BaseSpace do not presently have their own
370  internally managed and organized reference data, it is important to note that: (i) all four of these
371  platforms include apps that facilitate the download of reference data from other existing
372  databases (e.g., NCBI); (ii) all four have public data available for download and access to

373   varying degree (e.g., CyVerse has a large and diverse *data commons* while BaseSpace has
374   example datasets); and (iii) many apps integrated into these platforms maintain their own
375   sizable reference databases. This last point is critical, and in fact, provides one reason why
376   maintaining an internal reference database is useful. Users will run many different apps on their
377   data, which internally may utilize their own wide range of internal reference data. Problems may
378   arise when mixing and matching apps that rely on different and inconsistent reference data.
379
380   Overall, KBase stands out in this category for its breadth of genomic and biochemical reference
381   data.


## 2F. Comparison of Third-party Development and Custom Code Support

384   In terms of support for third-party development, KBase, BaseSpace, GenePattern, CyVerse,
385   and Galaxy all offer a similar experience, while Pathway Tools is quite different. In large part,
386   this is due to the fact that KBase, BaseSpace, GenePattern, CyVerse, and Galaxy were all
387   designed in part to serve as platforms for the deployment of third-party apps. All of these
388   platforms have their own equivalent of an app catalog (called the Discovery Environment in
389   CyVerse and the ToolShed in Galaxy). These app catalogs enable users to discover apps, read
390   documentation related to the apps, rate the apps, and view apps created and shared by other
391   users. All of these platforms offer an SDK of some form, and use virtualization technology to
392   simplify deployment (Galaxy, CyVerse, and BaseSpace use Docker just like KBase does).
393   Galaxy, KBase, CyVerse, GenePattern, and BaseSpace enable users to create their own
394   custom UIs for their apps using general spec files encoded in either XML (Galaxy) or JSON
395   (KBase, CyVerse, BaseSpace). The complexity associated with the addition of an app is roughly
396   equivalent in all five platforms. The object-based data model in KBase does create added
397   complexity for developers as they need to convert between objects and files when wrapping a
398   tool, but conversely, these developers also often benefit from the greater standardization of
399   objects fed into their tool on KBase. Galaxy has a lower learning curve and simpler app
400   development process, but app registration system-wide is more difficult. CyVerse, GenePattern,
401   BaseSpace and KBase have steeper learning curves to create a new app, but KBase has the
402   simplest interface for registering, sharing, and maintaining an app in the App Catalog.
403
404   One unique aspect of the KBase SDK is the ease with which a developer can programmatically
405   call any KBase app from within another app (e.g., the metabolic model reconstruction app can
406   invoke the species tree building app internally to place a genome in a specific phylogenetic
407   neighborhood). This capability will grow in power and importance as the number of apps
408   available in the KBase app catalog increases.
409
410   A second unique aspect of the KBase SDK is the support for binding in a wide range of
411   programming languages, as well as the ability to construct a new module from a standard
412   template in any of the same programming languages. CyVerse, GenePattern, BaseSpace, and
413   Galaxy all offer only REST web services for interacting with the broader platforms (e.g.,
414   accessing data).

415
416    Pathway Tools also supports third party development, but in a much less formal, flexible, or
417    automated fashion. Pathway Tools offers a LISP command interface, which exposes an API to
418    access Pathway Tools data and run Pathway Tools applications interactively. Users can use
419    this API to create new tools, and they can work with the Pathway Tools developers to get their
420    tool integrated into the platform, but there is no way for a user to do this independently without
421    interacting with core developers. Additionally, the tools must be written, at least in part, in LISP
422    in order to interact with the rest of the Pathway Tools platform. KBase, CyVerse, BaseSpace,
423    and Galaxy all offer language-agnostic REST or JSON RPC interfaces for platform interaction.

424    ## 2G. Comparison of Custom Code Support within Platform Workflows

425    One key differentiator of KBase is that it's built on the Jupyter framework, so it enables users to
426    seamlessly integrate IPython code cells into their workflows, either to run KBase apps in bulk, or
427    to implement custom analysis steps that are not yet implemented within an app. Both of these
428    capabilities are used to good effect in the exemplar Narratives described above. We explored in
429    our comparison whether any other selected platforms offer a similar capability.
430
431    GenePattern is the most similar platform to KBase in terms of this capability since it is also built
432    on top of the Jupyter framework. GenePattern also allows seamless integration of custom code
433    and Markdown cells into users' workflows, and GenePattern offers a programmatic interface for
434    running apps in the platform.
435
436    Neither CyVerse nor BaseSpace exhibit this capacity. The only way to integrate custom code
437    into a workflow in these environments is to create a new app using the SDK. However, CyVerse
438    does allow users to run local scripts in Shell, Perl, Python or R and run basic utilities based on
439    these languages in the Discovery Environment for data/file processing. Users can write custom
440    workflows using the LISP command interface in Pathway Tools, but rapidly sharing these
441    custom workflows with others is difficult, as generally sharing workflows in Pathway Tools is not
442    supported.
443
444    Galaxy does not natively support the integration of custom code either. However, there are
445    deploys of Galaxy available within Jupyter, where Python code can be used to run Galaxy apps
446    within Jupyter code cells. This functionality is not nearly as integrated as the code cells in
447    KBase, but it does offer a similar capability.
448
449    Overall, a strong argument can be made that this important feature, at least in a fully integrated
450    form, is unique to the KBase and GenePattern platforms. One of the challenges in using any
451    online workflow system for scientific analysis is that virtually every scientific workflow is distinct
452    (given the importance of novel discovery in science). Thus, it is very easy to run into a step in
453    one's workflow for which no app conveniently exists in the platform one is using. When this
454    happens, the user is forced to pull their data out of the platform they are using and load it into
455    another environment where the needed analysis is available. This disrupts the continuity,
456    provenance, and containment of the analysis. The analysis can no longer be shared as a single

457  self-contained entity for others to run. Thus, support for custom code is a truly important function
458  for this type of platform, and it is distinct to KBase and GenePattern.


## 2H. Comparison of Available Scientific Functionality

460  Scientific functionality is one of the areas with the greatest variability among the platforms
461  included in our comparison. Although there are plenty of examples of specific apps that are
462  available in multiple platforms, no two platforms offered the same range of functionality, and all
463  platforms had distinct areas of strength. Table 1 shows approximate support for different types
464  of functionality in each platform. Note that this comparison focuses on functionality relating to
465  the analysis of microbial and plant genomes, as this reflects the mission-space of KBase as a
466  DOE resource.
467  The app functionality in KBase differs from existing systems in several ways. The seamless
468  integration of code cells that KBase offers is distinctive, but not entirely unique--GenePattern[21]
469  and Synapse offer a similar capability. Galaxy[18], Taverna[24], CyVerse[19], XSEDE[25],
470  myExperiment[26], and GenePattern[21] overlap with many of the bioinformatics workflows in KBase
471  but lack the metabolic modeling capabilities. COBRA Toolbox[30], Pathway Tools[31], and RAVEN
472  Toolbox[32] support metabolic modeling but offer only minimal support for genome sequence
473  analysis. In terms of science functionality, each platform has its own set of strengths and
474  weaknesses. There are many categories where nearly all platforms have at least something to
475  offer (although some have more than others). These include genome assembly, genome
476  annotation, RNA-seq, comparative genomics, expression analysis, and assembly. However,
477  among these categories, some platforms are clear gold standards in certain areas: CyVerse,
478  GenePattern, and Galaxy for RNA-seq, variation, and comparative genomics; and KBase for
479  assembly and annotation.
480
481  Finally, there are areas of functionality that really distinguish between platforms because very
482  few platforms offer any functionality at all. In metabolic modeling, only Pathway Tools and
483  KBase offer functionality; in metagenome annotation, only CyVerse and Galaxy offer
484  functionality; and in metabolomics and chemistry, only Galaxy offers functionality. Currently,
485  KBase's capabilities for community model reconstruction, plant model reconstruction,
486  community model gapfilling, and expression data model integration are completely unique to the
487  KBase platform. Unsurprisingly, an examination of the app run counts in KBase reveals that the
488  most runs are applied to the apps where KBase is strongest: annotation, modeling, and
489  assembly


## 2I. Summary of Platform Comparison

491  Overall, from this comparison we see that KBase is extremely distinct from other bioinformatics
492  and computational biology platforms that exist today. The data model is a key area of
493  distinction: KBase uses biological types as objects, while most other platforms use files. While
494  this makes file conversion more challenging for developers, it brings the benefit of generally
495  making all tools much more interoperable and integrated, and requiring far fewer apps simply to
496  convert data from one format to another. The KBase data model also permits introspection,

497  enabling viewers, summary statistics, dropdowns and search utilities that offer views of
498  subobjects within an entity (reactions in a model, genes in a genome).

499  User experience is another key differentiator in KBase. With its integration on top of the Jupyter
500  framework, KBase offers users the ability to run their analysis workflows within a notebook
501  environment that also supports the seamless integration of custom text, graphics, and even
502  executable code. Users can organize their apps, text, and graphics into rich, reproducible,
503  scientific stories that may be shared and extended by others. Only GenePattern shares this
504  capability with KBase, and GenePattern lacks the tight integration of data into this notebook
505  environment that KBase offers.

506  Custom code support is another strength for KBase. KBase and GenePattern are the only
507  platforms that offer this capacity natively within their primary interface. This capability is
508  essential, as it enables a user to fill gaps in a workflow by writing custom code directly within the
509  Narrative notebook. This enables a user to maintain the continuity of a Narrative as much as
510  possible by avoiding the need to export data and analyze offline or in another platform in order
511  to complete custom analyses. There is an enormous benefit to reproducibility if all work for a
512  single complex study is performed in a single environment.

513  Reference data is also a key distinguishing component of the KBase platform. Only Pathway
514  Tools and KBase maintain their own internal curated reference data. This data is crucial for
515  placing user data into context. While other platforms do support reference data used by
516  individual tools, there is great benefit in consolidating reference data for tools at much as
517  possible to ensure that all tools are using a common reference data. The reference data in
518  KBase is also instrumental to planned functionality like the Knowledge Engine.

519  Finally, scientific functionality is another key area of distinction for KBase. KBase offers diverse
520  functionality that is nearly as broad as Galaxy (in terms of the number of categories of apps
521  available) if not always as deep as Galaxy (in terms of the number of apps in each category).
522  Additionally, there are some areas of functionality where KBase is a gold standard compared
523  with the other platforms in our comparison, including genome assembly, annotation, and
524  metabolic modeling.

525  This comparison focused on the most distinguishing features of the KBase platform, exploring
526  how these features contrast with other platforms. As a result, this comparison ignores many
527  important features that differentiate these other platforms from KBase. For example,
528  BaseSpace, Galaxy, and CyVerse all excel at the annotation and analysis of human and other
529  eukaryotic genome data. In contrast, the KBase user agreement explicitly prohibits the upload
530  and analysis of human data in KBase. As a further example, Galaxy and Pathway Tools both
531  excel at portability, meaning users can easily install and run their own instance of these
532  systems. In contrast, platforms like KBase, GenePattern, and BaseSpace are fully centralized.

533

## Supplementary Note 3. Code and Data Availability

### 3A. Code Availability

The KBase code, available at github.com/kbase, is open source and freely distributed under the MIT License. The web-accessible KBase system (narrative.kbase.us) is run on DOE computing infrastructure and is freely available for anyone to use. KBase adheres to the FAIR (Findable, Accessible, Interoperable, Re-usable) data principles endorsed by many funding agencies and scientific organizations[33]

### 3B. Data Availability

All data generated or analyzed during this study are included in this published article and Supplementary Note 1 as links to the original work, or in the associated KBase Narratives linked here. An earlier version of this paper was published as a preprint[34].

## Supplementary Figures

**Supplementary Figure 1.** Code cell for batch processing. A custom code cell created within the Narrative Interface that constructs CMMs for thousands of genomes.



```
Construction of Core Metabolic Models in Bulk using a Code Cell # Janaka

import time
import pprint
import json
import sys
from biokbase.narrative.jobs.appmanager import AppManager
version = "release"

#Construction of Core Metabolic Models in Bulk using a Code Cell # Janaka Edirisinghe v.o1 03/16/2017

ws = biokbase.narrative.clients.get("workspace")

my_list = ws.list_objects({'ids':['358'],'type':'KBaseGenomes.Genome','minObjectID':0,'maxObjectID':50
count = 0
for genome in my_list:
    #print 'Genome :', genome[1]
    count += 1
    if count < 501:
        #print 'Now processing genome id :', genome[1]

        print 'Starting Core Metabolic construction for the genome ',genome[1]
        genome_ob = ws.copy_object({'from':{'objid':genome[0],'wsid': 358},'to':{'wsid': 20186,'name':
```
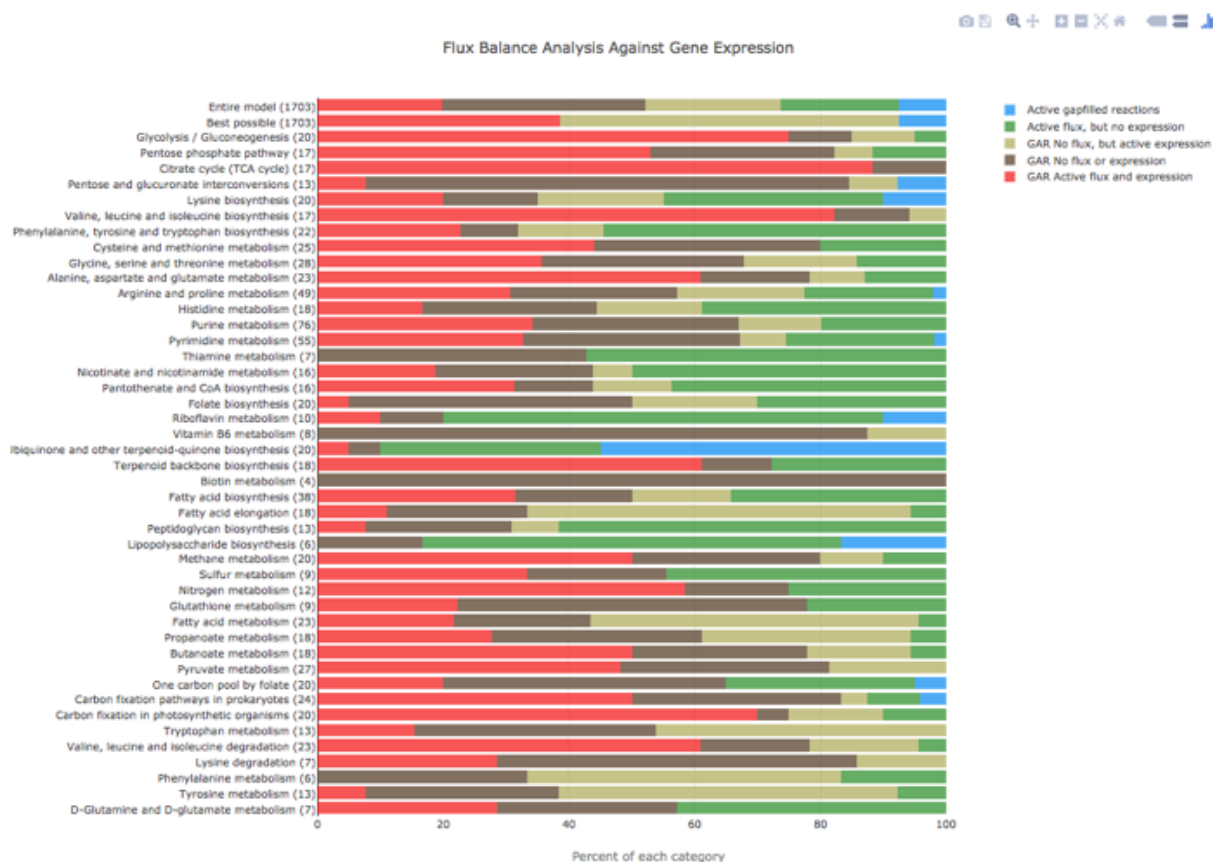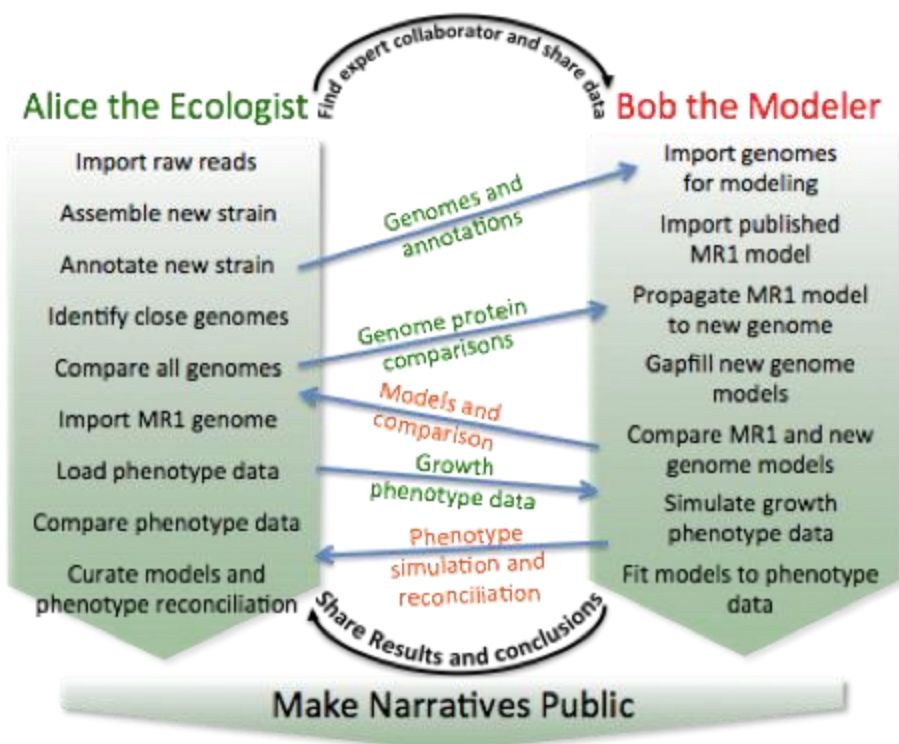
552 **Supplementary Figure 2.** "Compare Flux with Expression" output. Bar chart produced by the "Compare
553 Flux with Expression" app in KBase, which evaluates the reconciliation of metabolic model predictions
554 against expression data. All of the data points are categorized into biochemical pathways.



555
556

557 **Supplementary Figure 3.** Sharing of data, commentary and workflows using KBase Narratives. Two
558 researchers (Alice and Bob) create analysis workflows that complement each other's research resulting in
559 more intuitive and complete scientific conclusions.
560

561
562

## Supplementary Tables

564

**Supplementary Table 1.** Comparison of functionality available across evaluated platforms

| Category | KBase | CyVerse | BaseSpace | Galaxy | GenePattern | Pathway Tools |
|---|---|---|---|---|---|---|
| **Genome assembly** | High | Medium | Low | Medium | None | None |
| **Microbial genome annotation** | High | Medium | Low | Low | Low | Low |
| **Metagenome assembly and contig binning** | Medium | Medium | None | Low | None | None |
| **Metagenome annotation** | None | Medium | None | Medium | None | None |
| **Variation and GWAS** | None | High | Medium | High | High | None |

| RNA-seq | Medium | High | Low | High | High | None |
|---|---|---|---|---|---|---|
| **Metabolic modeling and chemistry** | High | None | None | None | Low | High |
| **Metabolomics and chemistry** | Low | None | None | Medium | None | Low |
| **Comparative genomics** | Low | High | Low | High | Low | Low |
| **Expression analysis** | Low | Low | Low | High | High | Low |

566


## Supplementary References

568

1.  Edirisinghe, J.N. et al. Modeling central metabolism and energy biosynthesis across microbial life. *BMC Genomics* **17**, 568 (2016).
2.  Wu, X. et al. Draft Genome Sequences of Two Janthinobacteriumlividum Strains, Isolated from Pristine Groundwater Collected from the Oak Ridge Field Research Center. *Genome Announc* **5** (2017).
3.  Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068-2069 (2014).
4.  Brettin, T. et al. RASTtk: a modular and extensible implementation of the RAST algorithm for building custom annotation pipelines and annotating batches of genomes. *Sci Rep* **5**, 8365 (2015).
5.  Li, L., Stoeckert, C.J., Jr. & Roos, D.S. OrthoMCL: identification of ortholog groups for eukaryotic genomes. *Genome Res* **13**, 2178-2189 (2003).
6.  Henry, C.S. et al. Microbial community metabolic modeling: A community data-driven network reconstruction. *J Cell Physiol* (2016).
7.  Weston, D.J. et al. Sphagnum physiology in the context of changing climate: emergent influences of genomics, modelling and host-microbiome interactions on understanding ecosystem function. *Plant Cell Environ* **38**, 1737-1751 (2015).
8.  Wilkinson, M.D. et al. The FAIR Guiding Principles for scientific data management and stewardship. *Sci Data* **3**, 160018 (2016).
9.  Wu, Y.W., Tang, Y.H., Tringe, S.G., Simmons, B.A. & Singer, S.W. MaxBin: an automated binning method to recover individual genomes from metagenomes using an expectation-maximization algorithm. *Microbiome* **2**, 26 (2014).
10. Aziz, R.K. et al. The RAST Server: rapid annotations using subsystems technology. *BMC Genomics* **9**, 75 (2008).
11. Price, M.N., Dehal, P.S. & Arkin, A.P. FastTree 2--approximately maximum-likelihood trees for large alignments. *PLoS One* **5**, e9490 (2010).
12. Henry, C.S. et al. High-throughput generation, optimization, and analysis of genome-scale metabolic models. *Nature Biotechnology* **Nbt.1672**, 1-6 (2010).

597    13.    Faria, J.P. et al. Constructing and Analyzing Metabolic Flux Models of Microbial
598           Communities. *Hydrocarbon and Lipid Microbiology Protocols* (2016).
599    14.    Latendresse, M. Efficiently gap-filling reaction networks. *BMC Bioinformatics* **15**, 225
600           (2014).
601    15.    Dreyfuss, J.M. et al. Reconstruction and validation of a genome-scale metabolic model
602           for the filamentous fungus Neurospora crassa using FARM. *PLoS Comput Biol* **9**,
603           e1003126 (2013).
604    16.    Ghosh, S. & Chan, C.K. Analysis of RNA-Seq Data Using TopHat and Cufflinks.
605           *Methods Mol Biol* **1374**, 339-361 (2016).
606    17.    Faria, J.P. et al. Computing and Applying Atomic Regulons to Understand Gene
607           Expression and Regulation. *Front Microbiol* **7**, 1819 (2016).
608    18.    Goecks, J., Nekrutenko, A., Taylor, J. & Team, G. Galaxy: a comprehensive approach
609           for supporting accessible, reproducible, and transparent computational research in the
610           life sciences. *Genome Biology* **11** (2010).
611    19.    Goff, S.A. et al. The iPlant Collaborative: Cyberinfrastructure for Plant Biology. *Front
612           Plant Sci* **2**, 34 (2011).
613    20.    Karp, P.D. et al. Pathway Tools version 19.0 update: software for pathway/genome
614           informatics and systems biology. *Brief Bioinform* **17**, 877-890 (2016).
615    21.    Reich, M. et al. GenePattern 2.0. *Nat Genet* **38**, 500-501 (2006).
616    22.    Perez, F. & Granger, B.E. IPython: A system for interactive scientific computing.
617           *Computing in Science & Engineering* **9**, 21-29 (2007).
618    23.    Kluyver, T. et al. Jupyter Notebooks—a publishing format for reproducible computational
619           workflows. *Positioning and Power in Academic Publishing: Players, Agents and
620           Agendas*, 87-90 (2016).
621    24.    Oinn, T. et al. Taverna: a tool for the composition and enactment of bioinformatics
622           workflows. *Bioinformatics* **20**, 3045-3054 (2004).
623    25.    Towns, J. et al. XSEDE: Accelerating Scientific Discovery. *Computing in Science &
624           Engineering* **16**, 62-74 (2014).
625    26.    Goble, C.A. et al. myExperiment: a repository and social network for the sharing of
626           bioinformatics workflows. *Nucleic Acids Research* **38**, W677-W682 (2010).
627    27.    Altintas, I. et al. Kepler: An extensible system for design and execution of scientific
628           workflows. *16th International Conference on Scientific and Statistical Database
629           Management, Proceedings*, 423-424 (2004).
630    28.    Deelman, E. et al. Pegasus: A framework for mapping complex scientific workflows onto
631           distributed systems. *Sci Programming-Neth* **13**, 219-237 (2005).
632    29.    Ananthakrishnan, R., Chard, K., Foster, I. & Tuecke, S. Globus platform-as-a-service for
633           collaborative science applications. *Concurrency and Computation-Practice & Experience*
634           **27**, 290-305 (2015).
635    30.    Schellenberger, J. et al. Quantitative prediction of cellular metabolism with constraint-
636           based models: the COBRA Toolbox v2.0. *Nat Protoc* **6**, 1290-1307 (2011).
637    31.    Karp, P.D. et al. The EcoCyc and MetaCyc databases. *Nucleic Acids Research* **28**, 56-
638           59 (2000).
639    32.    Agren, R. et al. The RAVEN toolbox and its use for generating a genome-scale
640           metabolic model for Penicillium chrysogenum. *PLoS Comput Biol* **9**, e1002980 (2013).
641    33.    Wilkinson, M.D. et al. The FAIR Guiding Principles for scientific data management and
642           stewardship. *Sci Data* **3**, 160018 (2016)
643    34.    Arkin, A.P. et al. The DOE Systems Biology Knowledgebase (KBase). *bioRxiv* **preprint
644           first posted online Dec. 22, 2016** (2016).
645