

# UC Berkeley

## UC Berkeley Previously Published Works

### Title

The Experiment Data Depot: A Web-Based Software Tool for Biological Experimental Data Storage, Sharing, and Visualization

### Permalink

<https://escholarship.org/uc/item/7nq2t7zt>

### Journal

ACS Synthetic Biology, 6(12)

### ISSN

2161-5063

### Authors

Morrell, William C

Birkel, Garrett W

Forrer, Mark

et al.

### Publication Date

2017-12-15

### DOI

10.1021/acssynbio.7b00204

Peer reviewed

## The Experiment Data Depot: a web-based software tool for biological experimental data storage, sharing, and visualization

William Morrell, Garrett Birkel, Mark Forrer, Teresa Lopez, Tyler Backman, Michael Dussault, Christopher J. Petzold, Edward E.K. Baidoo, Zak Costello, David Ando, Jorge Alonso Gutierrez, Kevin George, Aindrila Mukhopadhyay, Ian Vaino, Jay D Keasling, Paul D. Adams, Nathan J Hillson, and Hector Garcia Martin

*ACS Synth. Biol.*, **Just Accepted Manuscript** • DOI: 10.1021/acssynbio.7b00204 • Publication Date (Web): 21 Aug 2017

Downloaded from <http://pubs.acs.org> on August 23, 2017

### Just Accepted

“Just Accepted” manuscripts have been peer-reviewed and accepted for publication. They are posted online prior to technical editing, formatting for publication and author proofing. The American Chemical Society provides “Just Accepted” as a free service to the research community to expedite the dissemination of scientific material as soon as possible after acceptance. “Just Accepted” manuscripts appear in full in PDF format accompanied by an HTML abstract. “Just Accepted” manuscripts have been fully peer reviewed, but should not be considered the official version of record. They are accessible to all readers and citable by the Digital Object Identifier (DOI®). “Just Accepted” is an optional service offered to authors. Therefore, the “Just Accepted” Web site may not include all articles that will be published in the journal. After a manuscript is technically edited and formatted, it will be removed from the “Just Accepted” Web site and published as an ASAP article. Note that technical editing may introduce minor changes to the manuscript text and/or graphics which could affect content, and all legal disclaimers and ethical guidelines that apply to the journal pertain. ACS cannot be held responsible for errors or consequences arising from the use of information contained in these “Just Accepted” manuscripts.



1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

	Division Hillson, Nathan; Joint BioEnergy Institute, Fuels Synthesis Division Garcia Martin, Hector; Lawrence Berkeley National Laboratory, Joint BioEnergy Institute

SCHOLARONE™  
Manuscripts

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

# The Experiment Data Depot: a web-based software tool for biological experimental data storage, sharing, and visualization.

William C. Morrell,<sup>‡,§,∇</sup> Garrett W. Birkel,<sup>‡,||,⊥,∇</sup> Mark Forrer,<sup>‡,§,||</sup> Teresa Lopez,<sup>‡,§,||</sup> Tyler W. H. Backman,<sup>‡,||,⊥</sup> Michael Dussault,<sup>‡</sup> Christopher J. Petzold,<sup>‡,||,⊥</sup> Edward E. K. Baidoo,<sup>‡,||,⊥</sup> Zak Costello,<sup>‡,||,⊥</sup> David Ando,<sup>‡,⊥</sup> Jorge Alonso-Gutierrez,<sup>‡,⊥</sup> Kevin W. George,<sup>‡,⊥</sup> Aindrila Mukhopadhyay,<sup>‡,⊥</sup> Ian Vaino,<sup>‡</sup> Jay D. Keasling,<sup>‡,⊥,¶,||,⊥</sup> Paul D. Adams,<sup>‡,||,#</sup> Nathan J. Hillson,<sup>\*,‡,||,⊥,@</sup> and Hector Garcia Martin<sup>\*,‡,||,⊥,△</sup>

<sup>‡</sup>*DOE Joint BioEnergy Institute, Emeryville, CA, USA.*

<sup>§</sup>*Biotechnology and Bioengineering and Biomass Science and Conversion Department, Sandia National Laboratories, Livermore, CA, USA.*

<sup>||</sup>*DOE Agile BioFoundry, Emeryville, CA, USA.*

<sup>⊥</sup>*Biological Systems and Engineering Division, Lawrence Berkeley National Laboratory, Berkeley, CA, USA.*

<sup>¶</sup>*Department of Chemical and Biomolecular Engineering, University of California, Berkeley, CA, USA.*

<sup>||</sup>*Department of Bioengineering, University of California, Berkeley, CA, USA.*

<sup>⊥</sup>*Novo Nordisk Foundation Center for Biosustainability, Technical University Denmark, DK2970-Horsholm, Denmark.*

<sup>#</sup>*Molecular Biophysics and Integrated Bioimaging Division, Lawrence Berkeley National Laboratory, Berkeley, CA, USA.*

<sup>@</sup>*DNA Synthesis Science Program, DOE Joint Genome Institute, Walnut Creek, CA.*

<sup>△</sup>*BCAM, Basque Center for Applied Mathematics, Bilbao, Spain.*

<sup>∇</sup>*This article is distributed under a Creative Commons Attribution 4.0 International License.*

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

2 **Running header**

3 The Experiment Data Depot

## Abstract

Although recent advances in synthetic biology allow us to produce biological designs more efficiently than ever, our ability to predict the end result of these designs is still nascent. Predictive models require large amounts of high-quality data to be parametrized and tested, which are not generally available. Here, we present the Experiment Data Depot (EDD), an online tool designed as a repository of experimental data and metadata. EDD provides a convenient way to upload a variety of data types, visualize these data, and export them in a standardized fashion for use with predictive algorithms. In this paper, we describe EDD and showcase its utility for three different use cases: storage of characterized synthetic biology parts, leveraging proteomics data to improve biofuel yield, and the use of extracellular metabolite concentrations to predict intracellular metabolic fluxes.

**Keywords:** Database, -omics data, data standards, data mining, flux analysis, synthetic biology

The field of biology has undergone a radical transformation in the 20th and 21st centuries: whereas biology had previously been a descriptive science, focused on classifying and explaining biological behavior, the advent of genetic engineering and synthetic biology provides the possibility of changing the instruction set of biological entities and modifying their behavior.<sup>1</sup> The ensuing anticipated industrialization of biology in the 21st century<sup>2</sup> is expected to significantly impact society in several ways: a biobased economy has the potential to address key environmental challenges, transform manufacturing processes, increase the productivity and scope of the agricultural sector, reduce the economy's dependence on oil, improve human health, and grow new jobs and industries.<sup>3</sup>

However, while our capability to create new biological designs is advancing quickly, our ability to predict the outcome of engineered biological systems remains nascent. DNA synthe-

1  
2  
3  
4 30 sis productivity improves as fast as Moore's law,<sup>4</sup> and new tools for facile genome engineering  
5  
6 31 have revolutionized our capabilities to introduce site-specific modifications in the genomes  
7  
8 32 of cells and organisms.<sup>5</sup> Nonetheless, while it is increasingly more manageable to make the  
9  
10 33 DNA changes we intend, the end result on cell biology is generally unforeseen.<sup>6</sup>

11  
12 34 One of the main obstacles in predicting the behavior of biological systems is a concerning  
13  
14 35 lack of repeatability in bioengineering, as compared to other engineering disciplines. While  
15  
16 36 it is possible to produce a blueprint and specific instructions to construct (*e.g.*) a cell phone  
17  
18 37 in China that will satisfy the same specifications as the same phone built in the U.S., the  
19  
20 38 same is not the case for bioengineered systems.<sup>7</sup> Recent studies by Amgen and Bayer were  
21  
22 39 able to reproduce only 10-30% of biotech findings published in top-tiered journals,<sup>6,8,9</sup> and  
23  
24 40 there is a growing concern regarding lack of reproducibility.<sup>10</sup> This lack of reproducibility  
25  
26 41 not only hampers predictability, but also significantly limits investment in the field: the rule  
27  
28 42 of thumb that has been reported to be applied among venture capitalists is that 50% of  
29  
30 43 studies in top-journals are irreproducible.<sup>8</sup>

31  
32 44 Greater predictability and reproducibility requires efficient data, metadata, and proto-  
33  
34 45 col collection and sharing.<sup>7</sup> New computational biology approaches for predicting biological  
35  
36 46 behavior are becoming available, ranging from machine learning techniques to mechanistic  
37  
38 47 models.<sup>11-14</sup> However, the large amounts of standardized high-quality data that are needed  
39  
40 48 to rigorously validate or improve these models are lacking. Concurrently, the post-genomic  
41  
42 49 revolution has provided experimentalists with large-scale data sets of -omics data that are  
43  
44 50 orders of magnitude larger than they are typically trained to analyze. Hence, the collab-  
45  
46 51 oration between experimentalists and computational specialists could become much more  
47  
48 52 fruitful and frequent through a more robust exchange of data. In the field of synthetic bi-  
49  
50 53 ology, for example, it has been shown that careful characterization of synthetic biological  
51  
52 54 parts enables accurate prediction of full pathway behavior.<sup>15</sup>

55  
56 55 However, description of the experimental details is typically only reported in the materials  
57  
58 56 and methods of papers in non-standard, often incomplete, ways.<sup>6</sup> New frameworks such as  
59  
60

1  
2  
3  
4 57 the ISA (Investigation/Study/Assay) software<sup>16</sup> are appearing which provide a standardized  
5  
6 58 description of experiment design and metadata that have become the standard way to report  
7  
8 59 results to journals like Nature Scientific Data. This framework and others<sup>17,18</sup> have facili-  
9  
10 60 tated the appearance of tools for sharing transcriptomics,<sup>19,20</sup> proteomics,<sup>21–23</sup> metabolomics  
11  
12 61 data,<sup>24</sup> and even combinations of different -omics data types.<sup>25</sup> In parallel, data collection  
13  
14 62 and storage systems based on standards developed for medical purposes (DICOM, Digital  
15  
16 63 Imaging and COmmunication in Medicine<sup>26</sup>) are being applied to synthetic biology part  
17  
18 64 characterization (DICOM-SB<sup>27</sup>). However, none of these tools provides a single data repos-  
19  
20 65 itory for all -omics data types that is able to extract data straight from instrument output,  
21  
22 66 visualize this data, and export the data in formats that are readily applicable to modeling  
23  
24 67 tools and libraries.

25  
26 68 Here, we present the Experiment Data Depot (EDD), an online tool designed as a repos-  
27  
28 69 itory of experimental data and metadata (Fig. 1). EDD can uptake experimental data,  
29  
30 70 provide visualization of these data, and produce downloadable data in several standard out-  
31  
32 71 put formats. The input of data to EDD is performed through automated data streams:  
33  
34 72 each of these input streams automatically parses the standard outputs of the instruments  
35  
36 73 most commonly used for bioengineering. New input streams can be easily added to adapt  
37  
38 74 to local data production. The current version of EDD handles transcriptomics, proteomics,  
39  
40 75 metabolomics, HPLC, and Biolector<sup>®</sup> fermentation data. EDD provides a quick visualiza-  
41  
42 76 tion of imported data that allows for a quality check by showing whether the imported data  
43  
44 77 are within the expected range or not. Since data are stored internally in a relational database,  
45  
46 78 all data output is consistent. Outputs can be provided in terms of different standardized  
47  
48 79 files (Systems Biology Markup Language, SBML,<sup>28,29</sup> or CSV) or through a representa-  
49  
50 80 tional state transfer (RESTful) Application Programming Interface (API, in development).  
51  
52 81 Since the most common complaint of data scientists<sup>30</sup> is that they spend most of their time  
53  
54 82 preparing data for analysis rather than doing the analysis itself, the ability to obtain data in  
55  
56 83 standardized formats should be of great utility. SBML and CSV files can be used in conjunc-



tion with libraries such as COBRApy<sup>31</sup> or Scikit-learn<sup>32</sup> to generate actionable results for metabolic engineering. We showcase this capability by using HPLC data to predict internal metabolic fluxes of cells, and by leveraging proteomic data to improve biofuel yield. We also demonstrate EDD's capability to store information on characterized synthetic biology parts.

EDD is not a LIMS (Laboratory Information Management System): it is not meant to store raw data (*e.g.*, mass spectrometry traces). Rather, it only stores processed *biologically interpretable* data (*e.g.*, metabolite concentrations, protein expression levels, oxygen input rates, *etc.*), *i.e.* data that can be immediately interpreted by a biologist without requiring detailed knowledge of the analytical measurement technique.

## Methods

### Experiment description terms (EDD ontology)

EDD describes experiments in terms of studies, lines, strains, protocols, assays, measurements, and values (see Fig. 2 for an illustrative example).

- **Study** is used to describe a single continuous experiment meant to answer a single question. For example, an experiment characterizing the properties of a library of promoters in *Escherichia coli* would be a Study. Another example would be screening a panel of mutant enzymes for specificity to a molecule of interest.
- **Line** describes a single culture or line of enquiry within a Study. A single flask with a *E. coli* strain culture, or a well of *Saccharomyces cerevisiae* in a plate, are examples of Lines in EDD. Lines are grouped together under a Study in the EDD hierarchy, therefore a Study contains a set of Lines.
- **Strain** describes the biological entity used in a Line. A Line entry includes information about the strain or enzyme being used, making it possible to search for any Line or Study that uses a specific strain. Multiple Lines within a Study can use the same Strain,

1  
2  
3  
4 108 either as biological replicates, or under differing conditions. Additional information  
5  
6 109 concerning the strain(s) and/or plasmid(s) used in a Line is made available through  
7  
8 110 links to the Inventory of Composable Elements (ICE),<sup>33</sup> which serves as a repository  
9  
10 111 for DNA sequences, the physical location in the laboratory freezer, and other strain  
11  
12 112 metadata.

- 13  
14  
15 113 • **Protocol** denotes the method used to obtain information from a Line (*e.g.*, pro-  
16  
17 114 teomics). A Protocol is not tied to any particular Study; it is any repeatable process  
18  
19 115 meant to be used across many Studies. The description of a Protocol can be anything  
20  
21 116 from a simple list of written instructions, to a reference to a document or manual, or  
22  
23 117 a robot program.
- 24  
25  
26 118 • **Assay** is the application of a Protocol on a specific Line (*e.g.*, using proteomics to  
27  
28 119 study protein expression of Line C1, an *E. coli* culture). Assays are grouped under  
29  
30 120 Lines in the EDD hierarchy, thus a Line contains a set of Assays (see C1-PROT-1 and  
31  
32 121 C1-PROT-2 in Fig. 2).
- 33  
34  
35 122 • **Measurement** describes a quantity measured by an Assay (*e.g.*, the count of phos-  
36  
37 123 phoglucose isomerase proteins per cell found using the proteomics protocol on line  
38  
39 124 C1). Some Protocols will measure only one quantity (*e.g.*, optical density at 600 nm),  
40  
41 125 while others could measure multiple quantities (*e.g.*, several proteins for proteomics or  
42  
43 126 several extracellular metabolites for HPLC).
- 44  
45  
46 127 • **Values** are individual points of data for a Measurement. A Measurement could contain  
47  
48 128 only a single value, or several of them.
- 49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## 129 **Key capabilities**

### 130 **Data input**

131 Data input into EDD has been streamlined (Fig. 3 and Screencast 1 in Supporting Infor-  
132 mation). The data input menu consists of a set of prescribed import modules, plus a more  
133 general import option (Fig. 4). The assumption in this design is that typically the same types  
134 of data are imported, and new data types are only rarely added. The prescribed data inputs  
135 include options for HPLC data, targeted proteomics data, metabolomics concentration data,  
136 metabolite labeling patterns (such as those used in  $^{13}\text{C}$  Metabolic Flux Analysis<sup>34,35</sup>), tran-  
137 scriptomics data, and data obtained from the m2p-labs Biolector<sup>®</sup> automated fermentation  
138 platform.<sup>36</sup> A specific input format is expected for each data type depending on the data  
139 source (e.g. HPLC or Biolector) to standardize and facilitate data input. An example of  
140 the data format is shown in the input form as a guidance (see Screencast 1). New data type  
141 inputs can be easily added by including a new import module conforming to the interface  
142 for import/export modules (see Supporting Information).

143 Data lacking a specified format or type can be uploaded through a general import option.  
144 This option attempts to allow greater flexibility in defining rows and columns of an input  
145 table. A large variety of spreadsheet layouts may be handled by the general import, but  
146 this requires the user of EDD to define mappings of spreadsheet rows and columns to EDD  
147 datatypes.

### 148 **Visualization**

149 EDD provides visualization of experimental data through interactive tables and graphs (see  
150 Fig. 5 and Screencast 2 in Supporting Information). The guiding principle of visualization  
151 in EDD is that it is not meant to solve all visualization needs, but rather provide a general  
152 overview of datasets via visualization of the most common needs, while the rest can be  
153 tackled through data downloads and more sophisticated visualization tools (*e.g.*, Spotfire<sup>37</sup>

1  
2  
3  
4 154 or Plot.ly).

5  
6 155 The EDD study detail view contains several sections to present different facets of data  
7  
8 156 contained in the study: an overview part (“Overview”), a table describing lines and metadata  
9  
10 157 (“Experiment Description”) and an interactive graph displaying all collected data (“Data”).  
11  
12 158 The “Data” section (Fig. 5) allows the user to see different measurements for each line (*e.g.*,  
13  
14 159 acetate concentration for *E. coli* wild type strain or the number of copies of fumC protein in  
15  
16 160 engineered strain p3BB4) via different graph types: line, or bar graphs where data is grouped  
17  
18 161 by varying criteria. An interactive menu allows the user to toggle among different data types  
19  
20 162 or lines, in order to compare them. In this way, one can, for example, compare glucose  
21  
22 163 consumption for several strains, or lactate vs acetate production of a single strain. This  
23  
24 164 visualization gives the researcher a quick data quality check by testing whether the gathered  
25  
26 165 data matches intuitive expectations. The toggling is enabled through progressive filtering  
27  
28 166 of metadata criteria: Line, Strain, Protocol, Assay, Measurement (plus other metadata  
29  
30 167 customized for the Study). The filtering draws one column for each metadata type that has  
31  
32 168 more than one unique value in the Study, then lists the unique values in the column. When a  
33  
34 169 value in the column is checked the overview plot is updated to show only the records related  
35  
36 170 to the checked value. Also, the contents of all the columns to the right of the modified column  
37  
38 171 are updated to show which values remain in the currently visibly subset of records. In this  
39  
40 172 way, the user can progressively drill down into arbitrary groups of their data efficiently (see  
41  
42 173 Screencast 2 in Supporting Information for a demonstration).

43  
44 174 The “Experiment Description” section of the Study detail view collects the metadata and  
45  
46 175 descriptors of Lines into a searchable, filterable, and sortable table. Lines can be searched  
47  
48 176 through a box which filters out all lines not meeting the search criteria, and sorted by clicking  
49  
50 177 on headers, as in spreadsheets. The relevant metadata fields can be shown or hidden through  
51  
52 178 an options menu.  
53  
54  
55  
56  
57  
58  
59  
60

## 179 Data standardization

180 EDD provides a single repository of data and a set of unified workflows for data input which  
181 facilitate standardized data collection and storage. This standardization facilitates compar-  
182 ison of experiments accumulated over time and provides a unified input for data analysis.  
183 Furthermore, detailed protocols and metadata parameters for each type of measurement are  
184 stored within EDD. Including this additional context in data standards is important, so the  
185 researcher analyzing the data does not need to be the same individual or team who exe-  
186 cuted the experiments. This decoupling enables effective division of labor and helps improve  
187 productivity.<sup>7</sup>

188 EDD uses PubChem Compound Identifiers (cids) as the primary identifier for track-  
189 ing metabolites.<sup>38</sup> Common genome-scale models are supported by a pre-generated map-  
190 ping that connects BiGG<sup>39</sup> identifiers to cids by using ChEBI<sup>40</sup> as an intermediate, as  
191 there are BiGG<->ChEBI and ChEBI<->PubChem cross-references, but no direct BiGG<-  
192 >PubChem cross-references available. For databases other than BiGG, identifier mappings  
193 are not automatically resolved to PubChem cids. Novel metabolites not yet included in Pub-  
194 Chem can be added to the database via the administration interface, which stores chemical  
195 structures as a SMILES<sup>41</sup> string.

196 Proteins are tracked using the UniProt unique identifier (UPI,<sup>42</sup>), and *E. coli* genes are  
197 currently tracked using Blattner numbers (b-numbers<sup>43</sup>). Support for NCBI GenBank<sup>44</sup>  
198 accession numbers, a more standard and universal identifier than b-numbers, will be added  
199 in the very near future. Novel proteins and genes are also supported by adding them directly  
200 via the administration interface.

## 201 Data output

202 EDD provides access to all the data pertaining to an experiment in the form of standardized  
203 output files and a RESTful API in order to access data programmatically (in development).  
204 See Screencasts 3 and 4 in Supporting Information for a demonstration.

1  
2  
3  
4 205 Two output formats are provided at this time: comma separated values (CSV) and SBML.  
5  
6 206 The CSV format is a general spreadsheet format providing selected information for a given  
7  
8 207 experiment. Options on the CSV export can customize the output to include a subset of the  
9  
10 208 data of interest. There are three basic options for spreadsheet layout (illustrated in Fig. S1  
11  
12 209 in Supporting Information):

- 14  
15 210 • Rows of samples, columns of metadata and points; "short and wide". Suited for  
16  
17 211 researchers reading data across lots of samples.
- 18  
19 212 • Rows of data points, columns of metadata; "tall and skinny". Suited for loading into  
20  
21 213 analysis packages like Spotfire or R.
- 22  
23 214 • Rows of metadata and points, columns of samples; a transpose view of "short and  
24  
25 215 wide". Suited for researchers reading lots of points across a few samples.

26  
27  
28  
29 216 The SBML format is tailored to enable and facilitate flux analysis through COBRA  
30  
31 217 methods<sup>45</sup> or <sup>13</sup>C MFA.<sup>35</sup> The SBML output contains exchange fluxes and growth rates  
32  
33 218 calculated from the data stored in EDD as explained in the Supporting Information. In  
34  
35 219 order to make the SBML output useful for <sup>13</sup>C MFA,<sup>46</sup> it was necessary to supplement the  
36  
37 220 SBML standard with ways to include <sup>13</sup>C labeling patterns for different metabolites (see  
38  
39 221 Supporting Information). New standards for different outputs can be added as explained in  
40  
41 222 detail in the Supporting Information.

42  
43 223 The RESTful API is structured along the hierarchies illustrated in Figs. 2 and 7 (see  
44  
45 224 <https://github.com/JBEI/edd/tree/master/docs/Interface.md>). Accessing a Study will list  
46  
47 225 all the Lines in the study, accessing a Line will list all the Assays on the line, and so on, until  
48  
49 226 a script or program can access individual data points. When completed, the RESTful API  
50  
51 227 will allow access to the data in EDD with more complex query criteria than a straightforward  
52  
53 228 export can accommodate.  
54  
55  
56  
57  
58  
59  
60

## 229 **Read/edit permissions**

230 EDD includes a permissions model for Studies. A Study will be created by default with only  
231 permissions for the creator to view or edit. Without adding alternate permissions, a Study  
232 will be private, visible only to the individual creating the Study. Additional permissions may  
233 be granted to individual users, to groups of users, or to all users with accounts on the EDD  
234 server. There are two types of permissions available: the Read permission allows for viewing,  
235 searching, and exporting data from a Study; and the Write permission allows for adding,  
236 modifying, importing, or deleting data from a Study, as well as modifying permissions on  
237 the Study.

## 238 **Implementation**

239 The EDD code is open source under a Berkeley Software Distribution (BSD) license. The  
240 front-end of EDD is written in TypeScript, JavaScript, and HTML/CSS. EDD runs in any  
241 modern web browser, but Chrome is recommended (<https://www.google.com/chrome/>). The  
242 back-end is coded in Python and built on the Django platform (see Fig. 6). The code and  
243 documentation are available on Github (<https://github.com/JBEI/EDD>) and is divided into  
244 the following modules:

## 245 **Templates and Views**

246 The Django template framework is used to handle the layout and structure of EDD pages.  
247 Templates enforce a separation between how data in EDD are processed and how the same  
248 data are presented. By separating processing and presentation, the code for both is easier to  
249 generalize and re-use. A base template defines the overall look-and-feel of application pages  
250 and consistent navigation across the application. Additional templates referencing the base  
251 template define the structure for the major pages within EDD (*e.g.*, show study details; or,  
252 import instrument data).

253 Individual requests to EDD are handled with view functions. EDD directs requests to

1  
2  
3  
4 254 view functions based on the contents of the request URL. Then, the view function processes  
5  
6 255 the data in the request, loads and updates data from the database, and builds a response  
7  
8 256 using the view's template.  
9

## 10 11 257 **Front-end and visualization**

12  
13  
14 258 The lines, bars, axes, and labels in the overview plot are rendered in SVG via the D3  
15  
16 259 JavaScript library (d3js.org). Hovering over any line or bar triggers a CSS-based visual  
17  
18 260 effect to make it stand out from the others, and provides more details on the data behind  
19  
20 261 the visualization.

21  
22 262 The progressive filtering of metadata criteria is accomplished by creating a Typescript  
23  
24 263 class for a filtering column that accepts and then emits a set of records, and then subclassing  
25  
26 264 it for each of the base five kinds of metadata (Line, Strain, Protocol, Assay, Measurement),  
27  
28 265 plus a sixth subclass for all the customized metadata types that can appear in a Study. The  
29  
30 266 Measurement subclass is itself further subclassed for Metabolites, Proteins, and Transcripts.  
31  
32 267 When a Study page loads, each of these classes is instantiated once, and the resulting filtering  
33  
34 268 object is placed in an ordered list. Then, when a Study begins receiving data records from  
35  
36 269 the server, additional instantiations of the customized metadata subclass are made, one for  
37  
38 270 each new custom type detected. These objects are added to the beginning of the list.

39  
40 271 Each object is responsible for a column in the filtering section, and for accumulating and  
41  
42 272 then managing its list of unique values. To achieve progressive filtering, a set of all the data  
43  
44 273 records in the Study is fed into the first object in the list, which then emits another set,  
45  
46 274 possibly shortened by removing all the records that do not match any checked values in the  
47  
48 275 column. That set is passed to the next object, and further reduced, and so on, until the final  
49  
50 276 set is fed into the overview plot for display.  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



## 277 Database

278 Access to the EDD database is provided through the Django Object Relational Manager  
279 (ORM, Fig. 7). The ORM offers an interface to interact with entities in the database  
280 directly with Python code. This abstraction layer allows for EDD code to generally work  
281 with higher-level concepts of Studies, Assays, or Measurements instead of the underlying data  
282 models (*i.e.*, no need for SQL queries). Code execution can be triggered upon specified events  
283 through signal handlers in the ORM system. For example, a signal handler is responsible  
284 for updating Study information in the search index whenever a Study changes.

285 The data model for EDD centers on a few abstract concepts, tied together into the  
286 nested hierarchy of Study, Line, Assay, Measurement (Figs. 2 and 7). EDDObject defines  
287 the base for these parts of EDD. Each EDDObject has a unique machine-readable identifier,  
288 a human-readable name and description, update history, comments, files, and arbitrary  
289 metadata. Metadata, in turn, is defined by a MetadataType object. Each metadata value  
290 on an EDDObject references a MetadataType, containing the information needed for other  
291 code to interpret the value.

292 As an example, a Line is an EDDObject that has metadata describing the conditions of  
293 a biological sample. The specific metadata types used are customizable for each Line. The  
294 metadata that needs to be captured will differ between an experiment concerning cultures  
295 grown in flasks, compared to an experiment concerning corn growing in a field. Some meta-  
296 data values, like Strain, are in turn EDDObjects, containing additional metadata. Lines  
297 concerning strains link to corresponding strain entries in a strain repository, such as ICE.

298 The definitions of metadata are fully configurable, and can leverage existing specifications  
299 of metadata, such as those included in the DICOM-SB standard<sup>27</sup> or ISA-Tab.<sup>16</sup>

## 300 Importers and Exporters

301 EDD defines an interface for generalized import and export of data in various formats. There  
302 are two types of inputs: a protocol-specific input from a particular instrument (*e.g.*, HPLC

1  
2  
3  
4 303 or transcriptomics data), and a general import for data types not otherwise covered. Import  
5  
6 304 modules transform the data into structures of the EDD database (Fig. 4). Export modules  
7  
8 305 do the reverse process transforming selected data from the EDD database into other useful  
9  
10 306 output formats. These modules are the primary way to move data into and out of EDD.  
11  
12 307 Structuring the code as modules interfacing to and from the EDD database allows for the  
13  
14 308 input of complex workflows through the flexible combination of these modules. Hence, an  
15  
16 309 experiment that produced HPLC, transcriptomics, and proteomics data can have its data  
17  
18 310 introduced in EDD through a successive application of the respective modules (Fig. 3).

## 21 **Services**

22  
23  
24 312 EDD makes use of several open-source systems to provide services to the main application.  
25  
26 313 Each service is run using Docker containers ([www.docker.com](http://www.docker.com)), allowing for standard instal-  
27  
28 314 lation and deployment across servers. Installing and running a service only requires having a  
29  
30 315 Docker host and the name of a service image. Docker handles downloading all the packages  
31  
32 316 and code needed to run the service in an image. No separate installation is required, and  
33  
34 317 most service images will have a reasonable default configuration included.

35  
36 318 Code for EDD is itself collected into an image that will run in a Docker container. A  
37  
38 319 Dockerfile included in the source code describes all the required setup and install for the core  
39  
40 320 EDD service, and can be built into an image that is run just like any other service. Building  
41  
42 321 this image once will allow the same image to be copied to any Docker host and launch a new  
43  
44 322 instance.

45  
46 323 A simple overview of the services driving EDD is included in (Fig. 8). All services are  
47  
48 324 contained within the Docker Host. EDD connects to the Internet and outside world at two  
49  
50 325 points: with the Nginx web server ([www.nginx.com](http://www.nginx.com)) to handle web requests, and with the  
51  
52 326 Exim mail server (<http://www.exim.org/>) to send email notifications. Incoming requests to  
53  
54 327 Nginx get routed to the core EDD image running a Django website in the Gunicorn WSGI  
55  
56 328 application server, or to a backend file storage service. The core EDD service connects  
57  
58  
59  
60

1  
2  
3  
4 329 to several other services to implement specific features. Text search and faceting uses a  
5  
6 330 Solr document index service ([lucene.apache.org/solr/](http://lucene.apache.org/solr/)). A Redis cache ([redis.io](http://redis.io)) stores login  
7  
8 331 session information and copies of the latest versions of static web resources like images  
9  
10 332 and scripts. The core data model of EDD is implemented with a SQL schema running in  
11  
12 333 a PostgreSQL service ([www.postgresql.org](http://www.postgresql.org)). Any tasks that would take longer than the  
13  
14 334 duration of a typical web request are handled by a Celery service ([www.celeryproject.org](http://www.celeryproject.org))  
15  
16 335 running a copy of the EDD Docker image. Communication between the EDD application  
17  
18 336 service and the EDD worker service is mediated by a RabbitMQ message queue service  
19  
20 337 ([www.rabbitmq.com](http://www.rabbitmq.com)). Management of the message queue is handled by an optional Flower  
21  
22 338 service, which can also be connected to the Nginx service to enable management of the task  
23  
24 339 queue from outside of the Docker host.

25  
26 340 This microservice architecture of the EDD application ecosystem is intended to simplify  
27  
28 341 the process of expanding an installation of EDD. All of the services represented by rectangular  
29  
30 342 boxes in Fig. 8 are stateless services, meaning capacity can be added by replacing the service  
31  
32 343 box with a simple load balancer dividing the workload among multiple container copies.  
33  
34 344 The three stateful services: Solr, Redis, and Postgres; represented by upright cylinders, all  
35  
36 345 offer their own clustering solutions to scale beyond a single node. The file storage service,  
37  
38 346 represented by an overturned cylinder, can use any standard data storage strategy; from  
39  
40 347 local disks, to large RAID arrays, to large cloud storage providers like Amazon AWS S3  
41  
42 348 buckets.

## 46 47 349 **Results and discussion**

48  
49  
50 350 In this section, we present two example workflows that use experimental data contained  
51  
52 351 within EDD to produce actionable items for metabolic engineering. Another possible use  
53  
54 352 of EDD is to store synthetic biology parts characterization data, as is demonstrated by the  
55  
56 353 public version of EDD (<https://public-edd.jbei.org>). This instance of EDD holds the data  
57  
58  
59  
60

1  
2  
3  
4 354 for all the synthetic biology parts characterized in a recent publication concerning a Cas9-  
5  
6 355 based toolkit for instituting genetic changes in *S. cerevisiae* to optimize heterologous gene  
7  
8 356 expression.<sup>47</sup>

9  
10 357 The first workflow will show how to upload time-resolved HPLC data into EDD. We will  
11  
12 358 demonstrate the visualization capabilities and then download the data as a SBML file. We  
13  
14 359 will then show how to use this SBML file in conjunction with the COBRApy<sup>31</sup> library to  
15  
16 360 predict intracellular metabolic fluxes (which provide a comprehensive description of cellular  
17  
18 361 metabolism) through FBA (Flux Balance Analysis). FBA has important applications in  
19  
20 362 bioengineering,<sup>48,49</sup> microbial ecology<sup>50</sup> and biomedicine.<sup>51</sup>

21  
22 363 The second workflow will show how to upload targeted proteomics data into EDD, how  
23  
24 364 to view these data and how to download them for further analysis. We provide an example of  
25  
26 365 this further analysis by using the proteomics data obtained from a bioengineered *E. coli* strain  
27  
28 366 to increase production of limonene, repeating an analysis done in a previous publication.<sup>52</sup>

29  
30 367 Both of these workflows (and their input files) are demonstrated through Screencasts 4  
31  
32 368 and 5 in the Supporting Information, or at <https://public-edd.jbei.org/pages/tutorials/>.

## 36 369 **Using metabolite concentration data to derive internal metabolic** 37 38 370 **fluxes through Flux Balance Analysis (FBA)**

39  
40  
41 371 This workflow demonstrates how to upload time-resolved HPLC data into EDD, visualize  
42  
43 372 them and download them in the SBML format so internal metabolic fluxes can be calculated  
44  
45 373 through FBA.<sup>53</sup> The full workflow is showcased in Screencast 4. We will first introduce the  
46  
47 374 data in EDD in two steps (Fig. 3).

48  
49 375 We start at the main page and click on "Add New Study" on the upper right. The initial  
50  
51 376 step involves providing basic metadata information such as the study name, a brief descrip-  
52  
53 377 tion of the study and a contact person. This action prompts for an experiment description,  
54  
55 378 which can be introduced by dragging and dropping the file "FBA\_Experiment\_Description.xlsx"  
56  
57 379 (available as Supporting information and <https://public-edd.jbei.org/pages/tutorials/>). This  
58  
59  
60

1  
2  
3  
4 380 excel file contains a description of the experimental design on the basis of lines, as well as  
5  
6 381 the protocols applied and the corresponding assays (Fig. 2). Line information includes links  
7  
8 382 to detailed strain and plasmids information in ICE, as well as carbon source and media. In  
9  
10 383 this case, this minimal example describes two shaking flask cultures (line BW1 and ArcA)  
11  
12 384 of *E. coli* for which HPLC measurements of glucose and acetate are available at times 0,  
13  
14 385 7.5, 9.5, 11, 13, 15, and 17 hours. This template can be modified as desired to describe  
15  
16 386 different experiments. As soon as the experiment description is uploaded, the user can view  
17  
18 387 the corresponding lines and other experimental details.

19  
20 388 The next step is to upload data by clicking on "Import Data" on the upper right corner.  
21  
22 389 This action takes us to a data import page where the desired input format (the general  
23  
24 390 import in this case) and corresponding protocol ("HPLC" in this case) are chosen. The  
25  
26 391 HPLC data can be found in the "FBA\_HPLC.xlsx" file. Dragging and dropping this file in  
27  
28 392 the import page will make EDD parse the data and show an initial visualization, where the  
29  
30 393 user can discard undesired time points (e.g. having resulted from experimental mistakes).  
31  
32 394 EDD automatically matches the metabolite names to the database of standard metabolite  
33  
34 395 names included, and the user can correct this assignment if needed. Once "Submit Import"  
35  
36 396 is pressed, the data are now available on the main page of EDD for visualization. OD data  
37  
38 397 is uploaded in an analagous manner.

39  
40 398 The filtering section below the data graph provides the means to only look at certain parts  
41  
42 399 of the data set. For example, clicking on 'arcA' below 'Strain' only shows the HPLC data  
43  
44 400 corresponding to the arcA strain. Clicking on 'D-Glucose' below 'Metabolite' only shows the  
45  
46 401 HPLC data corresponding to the glucose measurement. Clicking on both, only shows the  
47  
48 402 acetate curves for the arcA strain (see Screencast 2 in the Supporting Information).

49  
50 403 Data can be downloaded in a standardized format for later analysis. In this case we will  
51  
52 404 download them in the SBML format. Exchange fluxes are automatically calculated from  
53  
54 405 the extracellular metabolite concentrations described in the HPLC data (see Supporting  
55  
56 406 Information). This file can be obtained by clicking on 'BW1' line, then selecting "Export  
57  
58  
59  
60

1  
2  
3  
4 407 Data" and then selecting "to SBML" and "Take Action". This procedure will take the  
5  
6 408 user to an export page that will determine the export parameters. The first one is which  
7  
8 409 genome-scale model to use as a base (*i.e.*, which genome-scale model to apply the previously  
9  
10 410 calculated exchange fluxes to). We will choose the *E. coli* iJO1366 model in this case, for  
11  
12 411 the sake of example. The second step will involve selecting which OD measurement values  
13  
14 412 will be used to constrain the biomass (biomass is assumed to be proportional to OD through  
15  
16 413 a constant value that is explicitly provided in this section and can be changed as needed).  
17  
18 414 These values are already preselected, so we only need to check that they are not obviously  
19  
20 415 wrong (*e.g.* set to zero). Step three involves pairing the calculated exchange fluxes with the  
21  
22 416 corresponding reactions in the genome-scale model. Finally, we can download the SBML file  
23  
24 417 for the desired time point by clicking on "Download SBML".

25  
26 418 The final step involves using the COBRApy library<sup>31</sup> and the SBML file downloaded  
27  
28 419 to predict internal metabolic flux profiles through Flux Balance Analysis.<sup>53</sup> We can predict  
29  
30 420 fluxes in five lines of code:

```
31  
32  
33 421 import cobra  
34  
35 422 model=cobra.io.read_sbml_model('EciJR904at20hrs.xml')  
36  
37 423 solution=model.optimize()  
38  
39 424 solution.objective_value  
40  
41 425 solution.fluxes[['PGI','GND']]
```

42  
43 426 which shows a value of predicted flux of 2.61 mmol/gdw/hr for PGI (glucose-6-phosphate  
44  
45 427 isomerase) and 0.91 mmol/gdw/hr for GND (hosphogluconate dehydrogenase). This code  
46  
47 428 along with the expected results are shown in Jupyter notebook A in the Supporting Infor-  
48  
49 429 mation.  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## 430 Using targeted proteomics data to improve biofuel production through 431 Principal Component Analysis (PCAP)

432 This workflow shows how to use EDD and the Scikit-learn library to leverage targeted pro-  
433 teomics data to improve biofuel production (limonene) by bioengineered *E. coli*, as demon-  
434 strated in Alonso-Gutierrez *et al.*<sup>52</sup> This workflow is showcased in Screencast 5 in the Sup-  
435 porting Information.

436 This example provides a demonstration of how to add several types of data using the two  
437 step process in Fig. 3. The initial steps of how to create a study are the same as for the pre-  
438 vious example, in terms of providing the basic metadata. The description of the experiment  
439 can be found in 'PCAP\_Experiment\_Description.xls': in this case there are thirty shake  
440 flask cultures (lines 2X-Mh to 2X-Hm) of *E. coli* for which targeted proteomics data samples  
441 are taken at 24 hrs. Dragging and dropping the file into the page obtained by clicking on  
442 "Add New Study" creates a new study reflecting all these details. The proteomics data can  
443 be found in the "PCAP\_Proteomics.csv" file. We can add these data to the study by clicking  
444 on "Import data" and following the instructions in the input page as shown in the previous  
445 example. This example has two additional data types associated besides the targeted pro-  
446 teomics data: limonene production measured through GC-MS ("PCAP\_GCMS.csv" file)  
447 and optical density measured through spectroscopy ("PCAP\_OD.xlsx" file). Adding the  
448 limonene measurements is as straightforward as pressing again "Import data" and follow-  
449 ing the instructions in the input page. Adding the optical density data follows the same  
450 procedure.

451 EDD offers several ways to visualize the data we previously loaded. In this example,  
452 the line graphs displaying the dependence with time are of limited use, since all data are  
453 collected at a single time. By clicking on "Bar Graphs" at the top of the "Data" tab, we  
454 can see this data in bar form grouped by measurement, line or time, as indicated by the  
455 different buttons. Hovering over each bar or data point gives further information. As before,  
456 we can filter certain types of data by clicking on "Filtering" and using the ensuing menu.

1  
2  
3  
4 457 By clicking on a line, protocol, or protein, we only see the data corresponding to that line,  
5  
6 458 protocol, or protein. The assays applied to each line and the sampling times are available  
7  
8 459 by clicking on the "Table" tab.

9  
10 460 We will now download the data from EDD for further analysis using Principal Component  
11  
12 461 Analysis of Proteomics (PCAP<sup>52</sup>). First, we select the lines we would like to download and  
13  
14 462 we click on "Export Data" and select "as CSV/etc" from the download menu options. This  
15  
16 463 provides a CSV file with a defined format that can be used as input for Jupyter notebook B  
17  
18 464 (see Supporting Information).

19  
20 465 The next steps involve taking the proteomics and production data and use Principal  
21  
22 466 Component Analysis to find which proteins need to have their expression changed in order  
23  
24 467 to improve biofuel production. This procedure is carried out using the Scikit-learn library,<sup>32</sup>  
25  
26 468 and is demonstrated in Jupyter notebook B. The input is the CSV file obtained from EDD,  
27  
28 469 and the output is Fig. 4 from Alonso-Gutierrez *et al.*,<sup>52</sup> which predicts which part of the  
29  
30 470 proteomics phase space is associated to improved limonene production (see publication for  
31  
32 471 further details).

## 33 34 35 36 37 472 **Conclusion**

38  
39  
40 473 We have presented in this manuscript EDD, an interactive online open-source tool that  
41  
42 474 serves as a repository of experimental data. Linked with ICE, EDD provides a standardized  
43  
44 475 description of experiments: from the strains and plasmids involved, to the protocols used,  
45  
46 476 the experimental design for sampling, and the data extracted. While the initial use cases and  
47  
48 477 the examples provided here are geared towards microorganism cultivation and phenotyping,  
49  
50 478 the data schema and different functionalities can be adapted to other uses (*e.g.*, enzyme  
51  
52 479 characterization or plant bioengineering).

53  
54 480 Data input can be done either manually through a web interface or through automated  
55  
56 481 workflows for typical data types. The latter includes input for: HPLC data, transcriptomics,  
57  
58  
59  
60



1  
2  
3  
4 482 proteomics data, metabolomics data, and Biolector data. These workflows provide a drag-  
5  
6 483 and-drop interface that parses data into the database automatically. These workflows are  
7  
8 484 modular, and new modules can be written for additional data types (*e.g.*, chip-Seq, etc).  
9  
10 485 Once the API in development is finished, it will provide the possibility of automating data  
11  
12 486 input, and hence ease the integration of data from other databases and publications.

13  
14 487 Data visualization is provided for each study through an interactive window where dif-  
15  
16 488 ferent data types can be seen simultaneously (Fig. 5). Different data types and strains can  
17  
18 489 be interactively filtered in or out to facilitate comparisons. Data for each protocol can be  
19  
20 490 found at the bottom of each study, along with sampling details.

21  
22 491 Data standardization is enabled by forcing all data into an ontology and using stan-  
23  
24 492 dardized ontologies for data (for example, all metabolomics data uses the same metabolite  
25  
26 493 names). Furthermore, the user is forced to include a minimum of metadata as a description  
27  
28 494 of metadata. A flexible use of metadata means that, beyond that minimum obligatory core,  
29  
30 495 extra metadata can be included, if desired, by the experimentalist.

31  
32 496 Data output can be done using a variety of formats, including CSV or SBML files. These  
33  
34 497 output streams are modular and new modules can be added for different output formats. By  
35  
36 498 virtue of the internal organization of EDD, all data output is consistent and can be used to  
37  
38 499 feed a variety of modeling or data mining approaches.

40 500 EDD improves on single -omics type databases such as PRIDE,<sup>54</sup> MOPED<sup>55</sup> and PAXdb<sup>56</sup>  
41  
42 501 (for *e.g.* proteomics) because it is able to integrate multiple types of -omics data (*e.g.* tran-  
43  
44 502 scriptomics, proteomics and metabolomics). Furthermore, the metadata typically stored in  
45  
46 503 these systems (*e.g.* PRIDE) focuses on data acquisition and sample preparation metadata  
47  
48 504 (*i.e.* trypsin amount, digestion length..), whereas experiment metadata (*e.g.* shaking speed,  
49  
50 505 culture volume, growth temperature) is typically lacking in these databases but is captured  
51  
52 506 on EDD. However, while some of these databases provide data analysis capabilities (*e.g.*  
53  
54 507 MOPED or PaxDb), EDD was not meant to perform complex data analysis. There are  
55  
56 508 many available tools available for data analysis (*e.g.* through Kbase or Jupyter notebooks)

1  
2  
3  
4 509 and we believe EDD's mission is not to choose those tools for the user but, rather, feed those  
5  
6 510 tools the standardized data they need, in order to streamline their use (see for example the  
7  
8 511 multi-omics data viewer Arrowland, <https://public-arrowland.jbei.org/>).

9  
10 512 In this manuscript, we have described two use cases for EDD in metabolic engineering (all  
11  
12 513 data available in the Supporting Information and <https://public-edd.jbei.org/pages/tutorials/>):  
13  
14 514 1) using extracellular metabolite concentrations to predict internal metabolic fluxes for an  
15  
16 515 *E. coli* strain using FBA, and 2) using proteomics data to increase biofuel production in a  
17  
18 516 bioengineered strain. These use cases are presented as tutorials and showcase the utility of  
19  
20 517 EDD for metabolic engineering and synthetic biology applications. EDD is, however, a tool  
21  
22 518 in continuous development. We present here a tool that addresses some of our current needs,  
23  
24 519 but the code is available to be modified and adapted to fit other future needs that require  
25  
26 520 collection and storage of large amounts of experimental data.

27  
28 521 EDD also provides a platform to disseminate the data produced at one institution to other  
29  
30 522 institutions, hence becoming a repository of data of use for testing and parametrizing models.  
31  
32 523 For example, JBEI's<sup>57</sup> public instance of EDD (<https://public-edd.jbei.org>) holds the infor-  
33  
34 524 mation for all the synthetic biology parts characterized in a recent JBEI publication which  
35  
36 525 provides the largest, most comprehensive Cas9-based toolkit to quickly institute genetic  
37  
38 526 changes in *S. cerevisiae* to optimize heterologous gene expression.<sup>47</sup> We expect to continue  
39  
40 527 to seed JBEI's public instance of EDD with data related to future publications from LBNL  
41  
42 528 (e.g. associated to JBEI or the Agile BioFoundry: <http://agilebio.lbl.gov/>), and very soon  
43  
44 529 open the possibility to other external researchers of uploading their own data. An alternative  
45  
46 530 is for external researchers to set their own instances of EDD (as explained in detail in the  
47  
48 531 github repository, [https://github.com/JBEI/edd/blob/master/docs/Developer\\_Setup.md](https://github.com/JBEI/edd/blob/master/docs/Developer_Setup.md)).  
49  
50 532 We also welcome contributions and joint development (see <https://github.com/JBEI/edd/blob/master/Con>  
51  
52 533 to fit other user's needs. Our final goal is to create a web of EDDs for different institutions  
53  
54 534 able to efficiently exchange data, as is the case for the web of registries (<https://www.jbei.org/jbeis->  
55  
56 535 [inventory-of-composable-elements-ice-tutorial-now-available/](https://www.jbei.org/jbeis-inventory-of-composable-elements-ice-tutorial-now-available/)).

1  
2  
3  
4 536 In the current world, where there is an increasingly strong trend to disclose algorithms  
5  
6 537 as open source code,<sup>58</sup> but training data is viewed as extremely valuable,<sup>59</sup> EDD will pro-  
7  
8 538 vide significant value as more experiments are available. We hope EDD will help enabling  
9  
10 539 reproducibility and predictability in the fields of metabolic engineering and synthetic biology.  
11  
12

## 13 14 540 **Competing interests**

15  
16  
17 541 The authors declare that they have no competing interests.  
18  
19

## 20 21 22 542 **Author's contributions**

23  
24  
25 543 HGM, GWB, WCM, MF, NJH, TL designed the software. JAG and KWB helped conceive  
26  
27 544 the tool functionalities. HGM, WCM, PDA, JDK, NJH, IV, EEKB, CP, ZC, DA, GWB,  
28  
29 545 TWHB, JAG, KWB, and AM wrote the paper. WCM, GWB, MF, TL, TWHB, MD wrote  
30  
31 546 the code.  
32  
33

## 34 35 36 547 **Acknowledgement**

37  
38  
39 548 This work was part of the DOE Joint BioEnergy Institute ([http:// www.jbei.org](http://www.jbei.org)) and  
40  
41 549 part of the Agile BioFoundry (<http://agilebiofoundry.org>) supported by the U.S. Depart-  
42  
43 550 ment of Energy, Energy Efficiency and Renewable Energy, Bioenergy Technologies Office,  
44  
45 551 through contract DE-AC02-05CH11231 between Lawrence Berkeley National Laboratory  
46  
47 552 and the U. S. Department of Energy. The United States Government retains and the  
48  
49 553 publisher, by accepting the article for publication, acknowledges that the United States  
50  
51 554 Government retains a nonexclusive, paid-up, irrevocable, worldwide license to publish or  
52  
53 555 reproduce the published form of this manuscript, or allow others to do so, for United  
54  
55 556 States Government purposes. The Department of Energy will provide public access to  
56  
57 557 these results of federally sponsored research in accordance with the DOE Public Access  
58  
59  
60

1  
2  
3  
4 558 Plan (<http://energy.gov/downloads/doe-public-access-plan>). NJH was also supported by the  
5  
6 559 DOE Joint Genome Institute (<https://jgi.doe.gov>) by the U.S. Department of Energy, Office  
7  
8 560 of Science, Office of Biological and Environmental Research, through contract DE-AC02-  
9  
10 561 05CH11231 between Lawrence Berkeley National Laboratory and the U.S. Department of  
11  
12 562 Energy. This research is also supported by the Basque Government through the BERC 2014-  
13  
14 563 2017 program and by Spanish Ministry of Economy and Competitiveness MINECO: BCAM  
15  
16 564 Severo Ochoa excellence accreditation SEV-2013-0323.

17  
18 565 We acknowledge and thank Daniel Lopez for contributing EDD's logo. We also thank Jacob  
19  
20 566 Coble, Sarah LaFrance, Xianwei "John" Meng, Oge Nnadi, Hector Plahar, Lisa Simirenko  
21  
22 567 and Ernst Oberortner for reviewing EDD's source code and their helpful suggestions.  
23  
24  
25

## 26 568 Supporting Information Available

27  
28  
29  
30 569 The following supporting information is available:  
31

- 32  
33 570 • Supporting text including (SupportingText.pdf):  
34  
35  
36 571 – A detailed explanation of how exchange fluxes are calculated for export from  
37  
38 572 EDD using SBML.  
39  
40 573 – An explanation of the extensions of SBML that were required in order to store  
41  
42 574 <sup>13</sup>C labeling data, transcriptomics, proteomics, metabolomics and fluxomics data.  
43  
44  
45 575 – Instructions on how to add new output standards for EDD.  
46  
47 576 – A supplementary figure explaining the different layouts of exported spreadsheets.  
48  
49  
50 577 • Five screencasts as tutorials that show five fundamental functionalities in EDD:  
51  
52  
53 578 – Screencast 1: Data upload (1-Data Input.mp4).  
54  
55  
56 579 – Screencast 2: Data visualization (2-Data Visualization.mp4).  
57  
58 580 – Screencast 3: Data download (3-Data Download.mp4).  
59  
60

- 1  
2  
3  
4 581 – Screencast 4: Using metabolite data for flux analysis (4-FBA.mp4).  
5  
6 582 – Screencast 5: Using proteomics data to increase biofuel production (5-PCAP.mp4).  
7  
8  
9 583 • A zip file (Examples.zip) containing Jupyter notebooks and input files to recreate:  
10  
11  
12 584 – Using metabolite data for flux analysis (Screencast 4):  
13  
14 585 \* Jupyter Notebook A.ipynb (+ corresponding html version)  
15  
16 586 \* FBA\_Experiment\_Description.xlsx  
17  
18 587 \* FBA\_HPLC.xlsx  
19  
20  
21 588 \* FBA\_OD.xlsx.  
22  
23 589 – Using proteomics data to increase biofuel production (Screencast 5):  
24  
25  
26 590 \* Jupyter Notebook B.ipynb (+ corresponding html version)  
27  
28 591 \* PCAP\_Experiment\_Description.xlsx  
29  
30 592 \* PCAP\_GCMS.csv  
31  
32 593 \* PCAP\_OD.xlsx  
33  
34  
35 594 \* PCAP\_Proteomics.csv  
36  
37  
38

## 595 References

- 39  
40  
41  
42 596 (1) Russo, E. (2003) Special Report: The birth of biotechnology. *Nature* 421, 456–457.  
43  
44  
45 597 (2) *Industrialization of Biology*; The National Academies Press, 2015.  
46  
47  
48 598 (3) House, T. W. (2012) National Bioeconomy Blueprint April 2012. *Industrial Biotechnol-*  
49  
50 599 *ogy* 8, 97–102.  
51  
52  
53 600 (4) Tang, N., Ma, S., and Tian, J. *Synthetic Biology*; Elsevier BV, 2013; pp 3–21.  
54  
55  
56 601 (5) Doudna, J. A., and Charpentier, E. (2014) The new frontier of genome engineering with  
57  
58 602 CRISPR-Cas9. *Science* 346, 1258096.  
59  
60

- 1  
2  
3  
4 603 (6) Gardner, T. S. (2013) Synthetic biology: from hype to impact. *Trends Biotechnol.* *31*,  
5  
6 604 123–125.  
7  
8  
9 605 (7) Chubukov, V., Mukhopadhyay, A., Petzold, C. J., Keasling, J. D., and Martín, H. G.  
10  
11 606 (2016) Synthetic and systems biology for microbial production of commodity chemicals.  
12  
13 607 *NPJ Syst. Biol. Appl.* *2*, 16009.  
14  
15  
16 608 (8) Prinz, F., Schlange, T., and Asadullah, K. (2011) Believe it or not: how much can  
17  
18 609 we rely on published data on potential drug targets? *Nat. Rev. Drug Discovery* *10*,  
19  
20 610 712–712.  
21  
22  
23 611 (9) Begley, C. G., and Ellis, L. M. (2012) Drug development: Raise standards for preclinical  
24  
25 612 cancer research. *Nature* *483*, 531–533.  
26  
27  
28 613 (10) Baker, M. (2016) 1,500 scientists lift the lid on reproducibility. *Nature* *533*, 452–454.  
29  
30  
31 614 (11) Karr, J. R., Sanghvi, J. C., Macklin, D. N., Gutschow, M. V., Jacobs, J. M., Bolival, B.,  
32  
33 615 Assad-Garcia, N., Glass, J. I., and Covert, M. W. (2012) A Whole-Cell Computational  
34  
35 616 Model Predicts Phenotype from Genotype. *Cell* *150*, 389–401.  
36  
37  
38 617 (12) Hyduke, D. R., Lewis, N. E., and Palsson, B. Ø. (2013) Analysis of omics data with  
39  
40 618 genome-scale models of metabolism. *Mol. BioSyst.* *9*, 167–174.  
41  
42  
43 619 (13) Nelli, F. *Python Data Analytics*; Springer Science Business Media, 2015; pp 237–264.  
44  
45  
46 620 (14) Gill, R. T., Halweg-Edwards, A. L., Clauset, A., and Way, S. F. (2015) Synthesis aided  
47  
48 621 design: The biological design-build-test engineering paradigm? *Biotechnol. Bioeng.*  
49  
50 622 *113*, 7–10.  
51  
52  
53 623 (15) Davidsohn, N., Beal, J., Kiani, S., Adler, A., Yaman, F., Li, Y., Xie, Z., and Weiss, R.  
54  
55 624 (2015) Accurate Predictions of Genetic Circuit Behavior from Part Characterization  
56  
57 625 and Modular Composition. *ACS Synth. Biol.* *4*, 673–681.  
58  
59  
60

- 1  
2  
3  
4 626 (16) Rocca-Serra, P., Brandizi, M., Maguire, E., Sklyar, N., Taylor, C., Begley, K., Field, D.,  
5  
6 627 Harris, S., Hide, W., Hofmann, O., Neumann, S., Sterk, P., Tong, W., and Sansone, S.-  
7  
8 628 A. (2010) ISA software suite: supporting standards-compliant experimental annotation  
9  
10 629 and enabling curation at the community level. *Bioinformatics* 26, 2354–2356.
- 11  
12  
13 630 (17) Brazma, A. et al. (2001) Minimum information about a microarray experiment  
14  
15 631 (MIAME)-toward standards for microarray data. *Nat Genet* 29, 365–71.
- 16  
17  
18 632 (18) Taylor, C. F. (2006) Minimum Reporting Requirements for Proteomics: A MIAPE  
19  
20 633 Primer. *PROTEOMICS* 6, 39–44.
- 21  
22  
23 634 (19) Clough, E., and Barrett, T. *Methods in Molecular Biology*; Springer Science Business  
24  
25 635 Media, 2016; pp 93–110.
- 26  
27  
28 636 (20) Brazma, A. (2003) ArrayExpress—a public repository for microarray gene expression  
29  
30 637 data at the EBI. *Nucleic Acids Res.* 31, 68–71.
- 31  
32  
33 638 (21) Jones, P. (2006) PRIDE: a public repository of protein and peptide identifications for  
34  
35 639 the proteomics community. *Nucleic Acids Res.* 34, D659–D663.
- 36  
37  
38 640 (22) Vizcaíno, J. A. et al. (2014) ProteomeXchange provides globally coordinated proteomics  
39  
40 641 data submission and dissemination. *Nat Biotechnol* 32, 223–226.
- 41  
42  
43 642 (23) Desiere, F. (2006) The PeptideAtlas project. *Nucleic Acids Res.* 34, D655–D658.
- 44  
45  
46 643 (24) Haug, K., Salek, R. M., Conesa, P., Hastings, J., de Matos, P., Rijnbeek, M., Ma-  
47  
48 644 hendraker, T., Williams, M., Neumann, S., Rocca-Serra, P., Maguire, E., Gonzalez-  
49  
50 645 Beltran, A., Sansone, S.-A., Griffin, J. L., and Steinbeck, C. (2012) MetaboLights—an  
51  
52 646 open-access general-purpose repository for metabolomics studies and associated meta-  
53  
54 647 data. *Nucleic Acids Res.* 41, D781–D786.
- 55  
56  
57 648 (25) Gonzalez-Beltran, A., Maguire, E., Georgiou, P., Sansone, S.-A., and Rocca-Serra, P.

- 1  
2  
3  
4 649 (2013) Bio-GraphIIIn: a graph-based integrative and semantically-enabled repository  
5  
6 650 for life science experimental data. *EMBnet.journal* 19, 46.  
7  
8  
9 651 (26) Pianykh, O. S. *Digital Imaging and Communications in Medicine (DICOM)*; Springer  
10  
11 652 Nature, 2011; pp 3–5.  
12  
13 653 (27) de Murieta, I. S., Bultelle, M., and Kitney, R. I. (2016) Toward the First Data Acqui-  
14  
15 654 sition Standard in Synthetic Biology. *ACS Synth. Biol.*  
16  
17  
18 655 (28) Finney, A., and Hucka, M. (2003) Systems biology markup language: Level 2 and  
19  
20 656 beyond. *Biochim. Soc. Trans.* 31, 1472–1473.  
21  
22  
23 657 (29) Hucka, M. *Encyclopedia of Systems Biology*; Springer Science Business Media, 2013; pp  
24  
25 658 2057–2063.  
26  
27  
28 659 (30) 2015 Data Science Report. 2015; [https://visit.crowdflower.com/](https://visit.crowdflower.com/2015-data-scientist-report)  
29  
30 660 [2015-data-scientist-report](https://visit.crowdflower.com/2015-data-scientist-report).  
31  
32  
33 661 (31) Ebrahim, A., Lerman, J. A., Palsson, B. O., and Hyduke, D. R. (2013) COBRAPy:  
34  
35 662 CONstraints-Based Reconstruction and Analysis for Python. *BMC Syst. Biol.* 7, 74.  
36  
37  
38 663 (32) others,, et al. (2011) Scikit-learn: Machine learning in Python. *J. Mach. Learn. Res.*  
39  
40 664 *12*, 2825–2830.  
41  
42  
43 665 (33) Ham, T. S., Dmytriv, Z., Plahar, H., Chen, J., Hillson, N. J., and Keasling, J. D.  
44  
45 666 (2012) Design implementation and practice of JBEI-ICE: an open source biological  
46  
47 667 part registry platform and tools. *Nucleic Acids Res.* 40, e141–e141.  
48  
49  
50 668 (34) Wiechert, W. (2001) 13C Metabolic Flux Analysis. *Metab. Eng.* 3, 195–206.  
51  
52  
53 669 (35) Martín, H. G., Kumar, V. S., Weaver, D., Ghosh, A., Chubukov, V., Mukhopadhyay, A.,  
54  
55 670 Arkin, A., and Keasling, J. D. (2015) A Method to Constrain Genome-Scale Models  
56  
57 671 with 13C Labeling Data. *PLoS Comput. Biol.* 11, e1004363.  
58  
59  
60



- 1  
2  
3  
4 672 (36) Funke, M., Buchenauer, A., Schnakenberg, U., Mokwa, W., Diederichs, S., Mertens, A.,  
5  
6 673 Müller, C., Kensy, F., and Büchs, J. (2010) Microfluidic biolector-microfluidic biopro-  
7  
8 674 cess control in microtiter plates. *Biotechnol. Bioeng.* *107*, 497–505.
- 9  
10  
11 675 (37) Wilkins, C. L. (2000) Books and Software: Data mining with Spotfire Pro 4.0. *Anal.*  
12  
13 676 *Chem.* *72*, 550 A–550 A.
- 14  
15  
16 677 (38) Bolton, E. E., Wang, Y., Thiessen, P. A., and Bryant, S. H. (2008) PubChem: integrated  
17  
18 678 platform of small molecules and biological activities. *Annu. Rep. Comput. Chem.* *4*,  
19  
20 679 217–241.
- 21  
22  
23 680 (39) Schellenberger, J., Park, J. O., Conrad, T. M., and Palsson, B. Ø. (2010) BiGG: a Bio-  
24  
25 681 chemical Genetic and Genomic knowledgebase of large scale metabolic reconstructions.  
26  
27 682 *BMC bioinf.* *11*, 213.
- 28  
29  
30 683 (40) Degtyarenko, K., De Matos, P., Ennis, M., Hastings, J., Zbinden, M., McNaught, A.,  
31  
32 684 Alcántara, R., Darsow, M., Guedj, M., and Ashburner, M. (2007) ChEBI: a database  
33  
34 685 and ontology for chemical entities of biological interest. *Nucleic Acids Res.* *36*, D344–  
35  
36 686 D350.
- 37  
38  
39 687 (41) Weininger, D., Weininger, A., and Weininger, J. L. (1989) SMILES. 2. Algorithm for  
40  
41 688 generation of unique SMILES notation. *J. Chem. Inf. Model.* *29*, 97–101.
- 42  
43  
44 689 (42) Consortium, U. (2015) UniProt: a hub for protein information. *Nucleic Acids Res.* *43*,  
45  
46 690 D204–12.
- 47  
48  
49 691 (43) Blattner, F. R., Plunkett, G., Bloch, C. A., Perna, N. T., Burland, V., Riley, M.,  
50  
51 692 Collado-Vides, J., Glasner, J. D., Rode, C. K., and Mayhew, G. F. (1997) The complete  
52  
53 693 genome sequence of *Escherichia coli* K-12. *Science* *277*, 1453–1462.
- 54  
55  
56 694 (44) Benson, D. A., Cavanaugh, M., Clark, K., Karsch-Mizrachi, I., Lipman, D. J., Ostell, J.,  
57  
58 695 and Sayers, E. W. (2012) GenBank. *Nucleic Acids Res.* *41*, D36–D42.

- 1  
2  
3  
4 696 (45) Schellenberger, J., Que, R., Fleming, R. M. T., Thiele, I., Orth, J. D., Feist, A. M.,  
5  
6 697 Zielinski, D. C., Bordbar, A., Lewis, N. E., Rahmanian, S., Kang, J., Hyduke, D. R., and  
7  
8 698 Palsson, B. Ø. (2011) Quantitative prediction of cellular metabolism with constraint-  
9  
10 699 based models: the COBRA Toolbox v2.0. *Nat Protoc* 6, 1290–1307.
- 11  
12  
13 700 (46) Martín, H., Kumar, V., Weaver, D., Ghosh, A., Chubukov, V., Mukhopadhyay, A.,  
14  
15 701 Arkin, A., and Keasling, J. (2015) A Method to Constrain Genome-Scale Models with  
16  
17 702 <sup>13</sup>C Labeling Data. *PLoS Comput Biol* 11, e1004363.
- 18  
19  
20 703 (47) Apel, A. R., d Espaux, L., Wehrs, M., Sachs, D., Li, R. A., Tong, G. J., Garber, M.,  
21  
22 704 Nnadi, O., Zhuang, W., Hillson, N. J., Keasling, J. D., and Mukhopadhyay, A. (2016)  
23  
24 705 A Cas9-based toolkit to program gene expression in *Saccharomyces cerevisiae*. *Nucleic*  
25  
26 706 *Acids Res.* 45, 496–508.
- 27  
28  
29 707 (48) Yim, H. et al. (2011) Metabolic engineering of *Escherichia coli* for direct production of  
30  
31 708 1,4-butanediol. *Nat. Chem. Biol.* 7, 445–452.
- 32  
33  
34 709 (49) Park, J. H., Lee, K. H., Kim, T. Y., and Lee, S. Y. (2007) Metabolic engineering of  
35  
36 710 *Escherichia coli* for the production of L-valine based on transcriptome analysis and in  
37  
38 711 silico gene knockout simulation. *Proc. Natl. Acad. Sci. U. S. A.* 104, 7797–7802.
- 39  
40  
41 712 (50) Stolyar, S., Dien, S. V., Hillesland, K. L., Pinel, N., Lie, T. J., Leigh, J. A., and  
42  
43 713 Stahl, D. A. (2007) Metabolic modeling of a mutualistic microbial community. *Mol.*  
44  
45 714 *Syst. Biol.* 3.
- 46  
47  
48 715 (51) Frezza, C. et al. (2011) Haem oxygenase is synthetically lethal with the tumour sup-  
49  
50 716 pressor fumarate hydratase. *Nature* 477, 225–228.
- 51  
52  
53 717 (52) Alonso-Gutierrez, J., Kim, E.-M., Batth, T. S., Cho, N., Hu, Q., Chan, L. J. G.,  
54  
55 718 Petzold, C. J., Hillson, N. J., Adams, P. D., Keasling, J. D., Martin, H. G., and Lee, T. S.  
56  
57 719 (2015) Principal component analysis of proteomics (PCAP) as a tool to direct metabolic  
58  
59 720 engineering. *Metabol. Eng.* 28, 123–133.
- 60

- 1  
2  
3  
4 721 (53) Lewis, N. E., Nagarajan, H., and Palsson, B. O. (2012) Constraining the metabolic  
5  
6 722 genotype–phenotype relationship using a phylogeny of in silico methods. *Nat. Rev.*  
7  
8 723 *Microbiol.*
- 9  
10  
11 724 (54) Vizcaíno, J. A., Csordas, A., Del-Toro, N., Dianes, J. A., Griss, J., Lavidas, I.,  
12  
13 725 Mayer, G., Perez-Riverol, Y., Reisinger, F., and Ternent, T. (2015) 2016 update of  
14  
15 726 the PRIDE database and its related tools. *Nucleic Acids Res.* *44*, D447–D456.
- 16  
17  
18 727 (55) Kolker, E., Higdon, R., Haynes, W., Welch, D., Broomall, W., Lancet, D., Stanberry, L.,  
19  
20 728 and Kolker, N. (2011) MOPED: model organism protein expression database. *Nucleic*  
21  
22 729 *Acids Res.* *40*, D1093–D1099.
- 23  
24  
25 730 (56) Wang, M., Weiss, M., Simonovic, M., Haertinger, G., Schrimpf, S. P., Hengart-  
26  
27 731 ner, M. O., and von Mering, C. (2012) PaxDb, a database of protein abundance averages  
28  
29 732 across all three domains of life. *Mol. Cell. Proteomics* *11*, 492–500.
- 30  
31  
32 733 (57) Scheller, H. V., Singh, S., Blanch, H., and Keasling, J. D. (2010) The Joint BioEn-  
33  
34 734 ergy Institute (JBEI): Developing New Biofuels by Overcoming Biomass Recalcitrance.  
35  
36 735 *BioEnergy Res.* *3*, 105–107.
- 37  
38  
39 736 (58) Metz, C. Google Just Open Sourced TensorFlow, Its Artifi-  
40  
41 737 cial Intelligence Engine. 2015; [https://www.wired.com/2015/11/  
42  
43 738 google-open-sources-its-artificial-intelligence-engine/](https://www.wired.com/2015/11/google-open-sources-its-artificial-intelligence-engine/).
- 44  
45  
46 739 (59) Vanian, J. IBM bought The Weather Company because weather af-  
47  
48 740 fects nearly everything. 2015; [http://fortune.com/2015/10/28/  
49  
50 741 ibm-weather-company-acquisition-data/](http://fortune.com/2015/10/28/ibm-weather-company-acquisition-data/).
- 51  
52  
53  
54  
55  
56  
57  
58  
59  
60

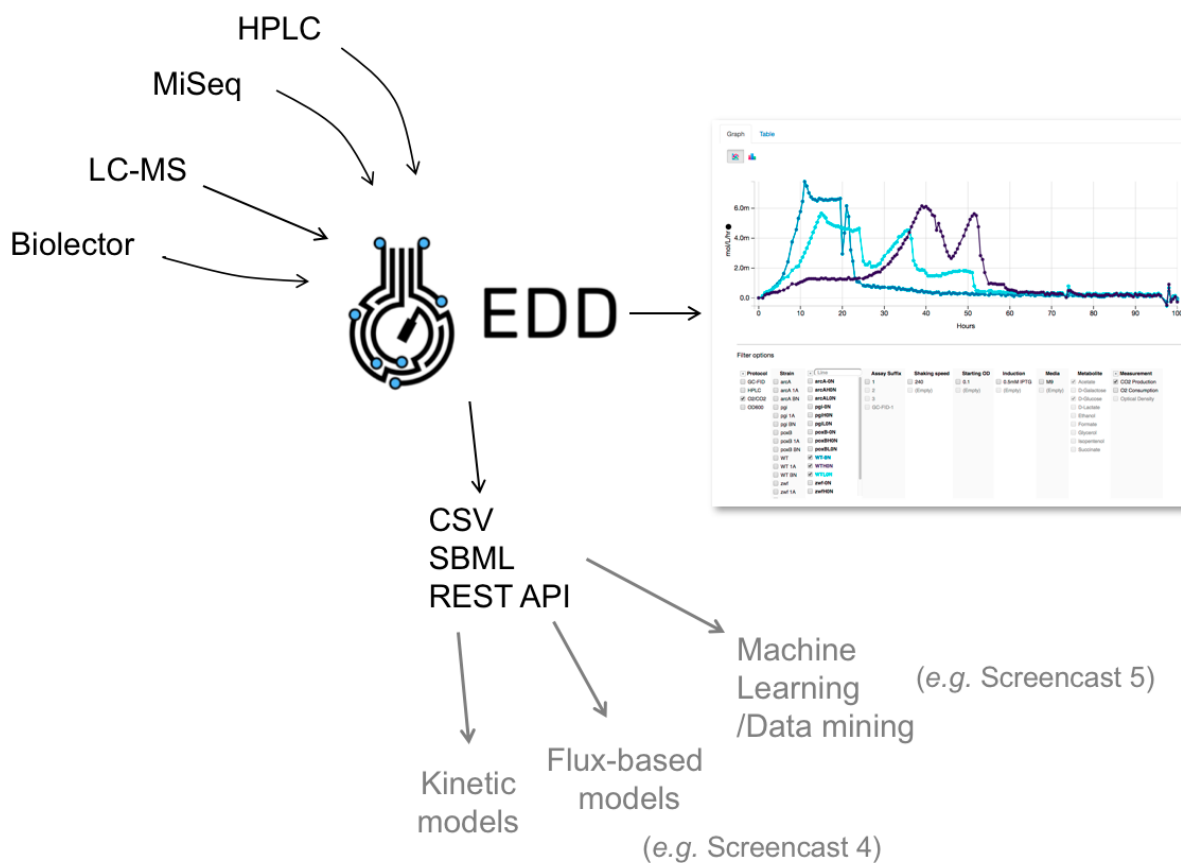
742 **Figures**

Figure 1: **Overview and key capabilities of EDD.** EDD collects data from different instruments, stores and visualizes them in an interactive way, and enables downloading them in a standardized format for use with a variety of modeling and analysis techniques. Screencasts 4 and 5, available in the Supporting Information (or <https://public-edd.jbei.org/pages/tutorials/>), provide step by step example tutorials to calculate internal metabolic fluxes, or to use proteomics data to improve biofuel production through data mining.

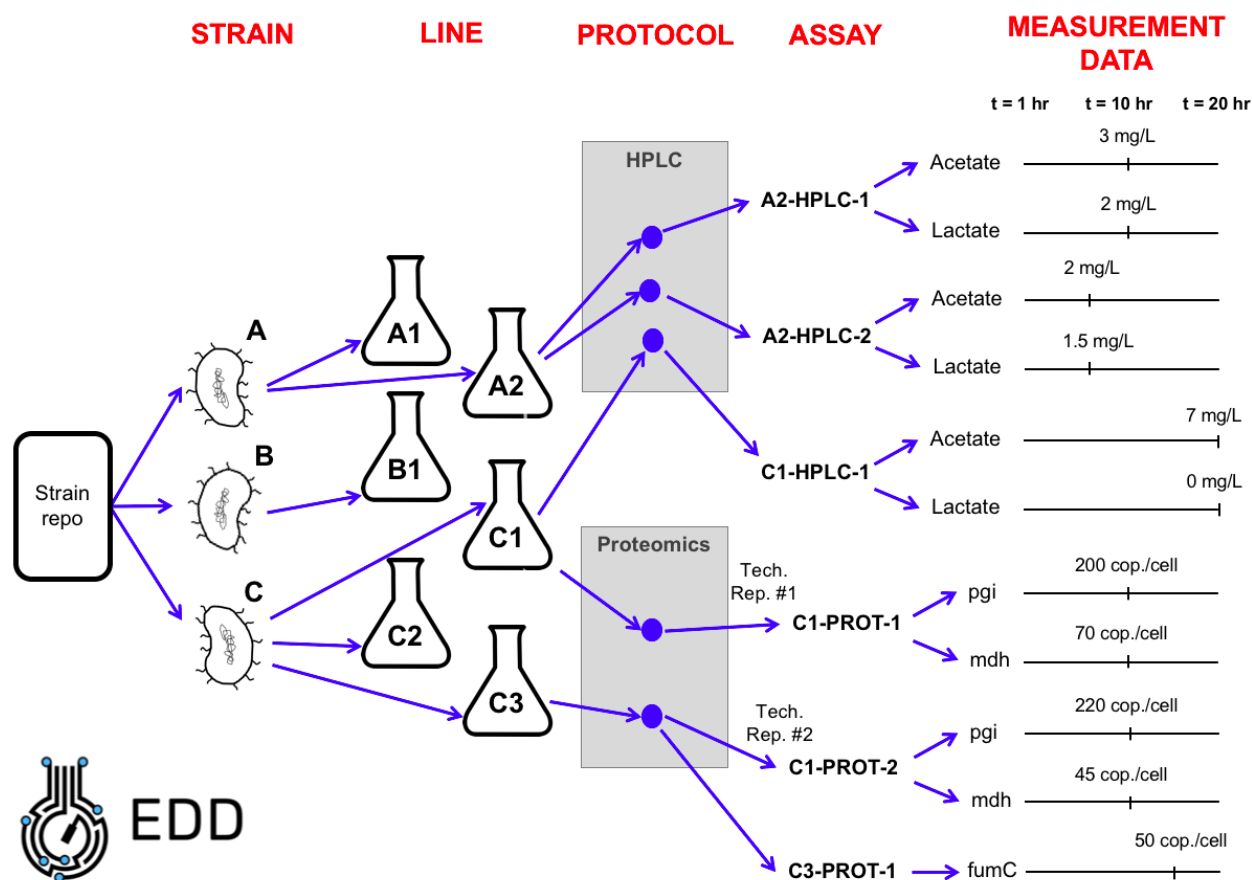


Figure 2: **Experiment description on EDD.** Example of how a common experiment would be described in EDD. This *study* involves culturing three *strains* (A, B and C) from a strain repository in several shaking flasks. Strain A is cultured in two flasks giving rise to two *lines* (A1 and A2). Strain B is cultured in a single flask (line B1) and strain C is cultured in three different flasks (lines C1, C2, and C3). The HPLC (High Pressure Liquid Chromatography measuring extracellular metabolite concentrations) *protocol* is applied to line A2 at t=10 hr giving rise to *assay* A2-HPLC-1. For assay A2-HPLC-1 the *measurement data* for acetate and lactate were 3 and 2 mg/L, respectively, at t=10 hr. Line C1 is subject to two different protocols: HPLC (t= 20 hr) and proteomics (quantitative measurement of expressed proteins, t=10 hr). Proteomics assay C1-PROT-1 on line C1 yields a measurement of 200 copies of pgi (phosphoglucose isomerase) per cell, and 70 copies of mdh (malate dehydrogenase) per cell at t=10 hr. A technical replicate of this measurement, coming from a different line (flask), constitutes a different assay C1-PROT-2.

# 1 Load experiment description

Line	Media	Vol.
BW1	M9	50mL
arcA	M9	50mL
.....		

# 2 Add data from instrument

HPLC data		
Time	BW1	BW1
0	Gluc.	Ac.
7.5	22	0.1
9.5	15	0.6
11	10	0.9
.....		



EDD

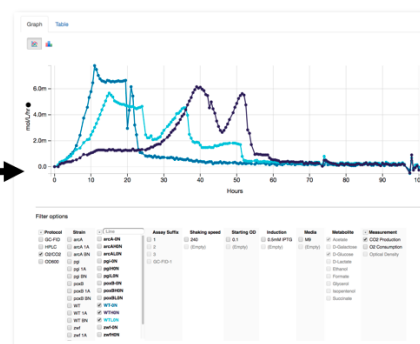


Figure 3: **Data input into EDD** has been streamlined. Users can input data in two steps. The first step involves adding description of the experiment describing lines and metadata for the study, as exemplified in Fig. 2. The second step involves uploading the data: (*e.g.*, HPLC data with metabolite concentrations). The input is modular, so additional data (*e.g.*, proteomics, transcriptomics, *etc.*) can be added later using the same import protocols. See Screencast 1 in the Supporting Information, or at <https://public-edd.jbei.org/pages/tutorials/>, for a demonstration.

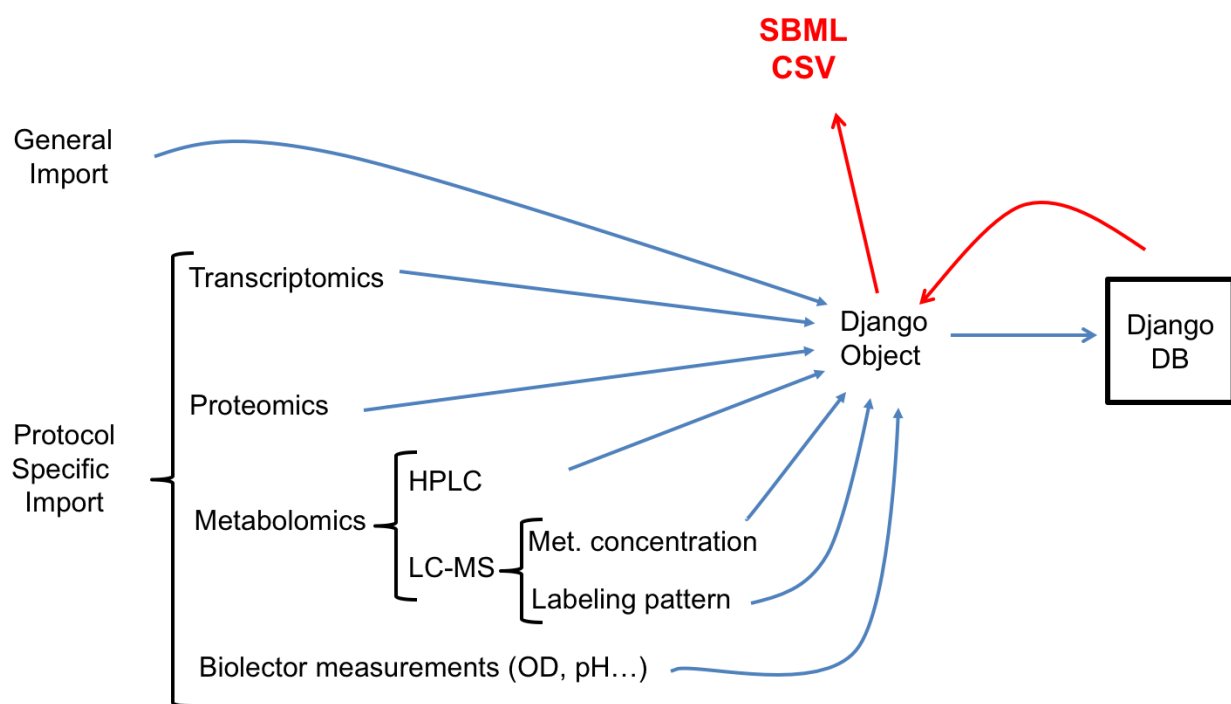
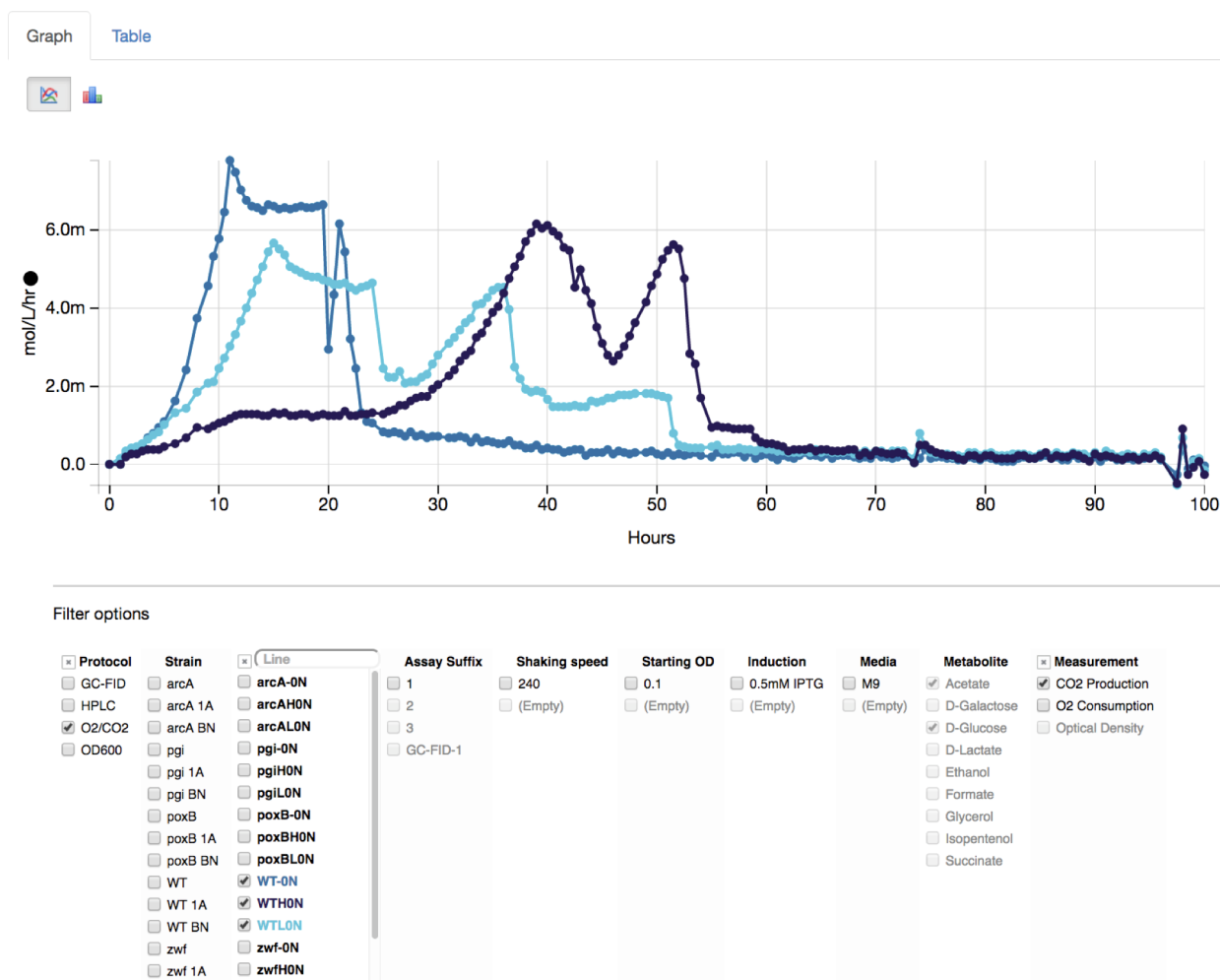


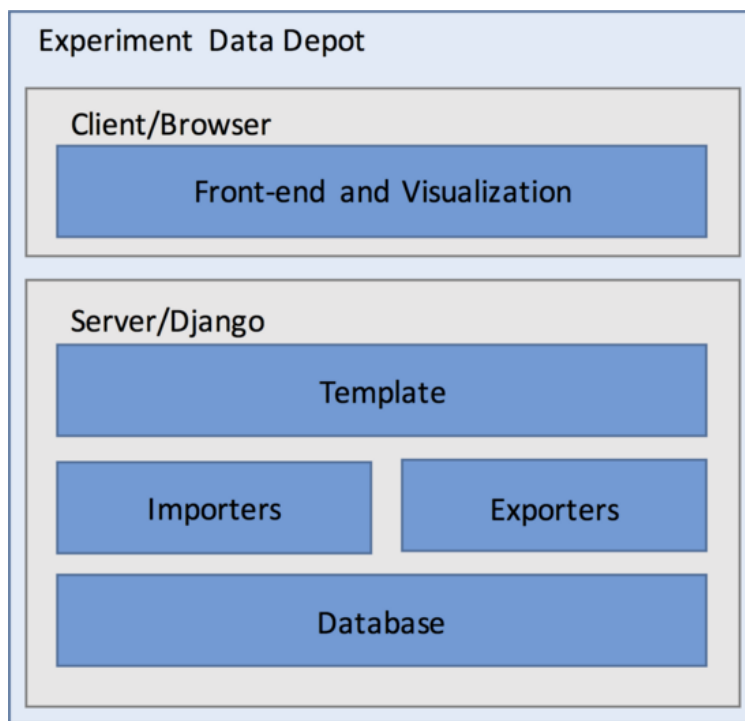
Figure 4: **Export and import modules.** Inputs are divided into two groups: protocol specific import that comes from a specific machine with a predetermined format, and a general import. Inputs are written so as to produce a Django object that is then stored in the database. The same modules are used for data export in SBML and CSV format. See Screenshot 3 in the Supporting Information, or at <https://public-edd.jbei.org/pages/tutorials/>, for a demonstration of data export.



43  
44  
45  
46  
47  
48  
49  
50  
51  
52

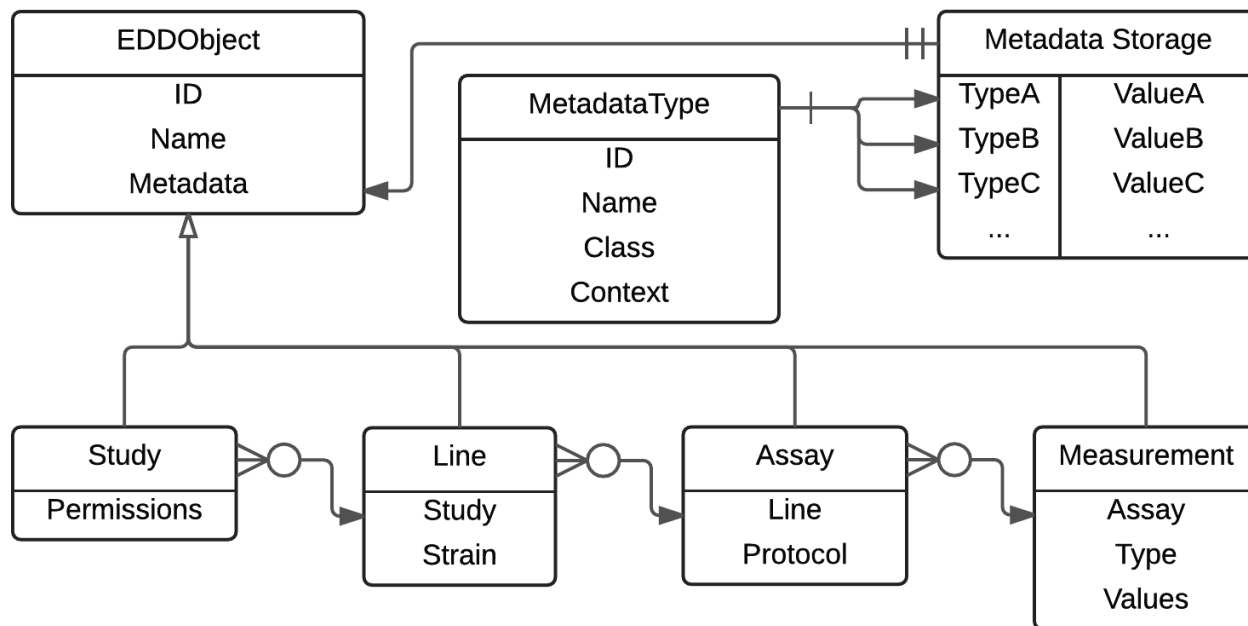
Figure 5: **Interactive data visualization.** The "Data" tab provides an interactive visualization of all data contained in a single study. The "Filter options" menu contains a classification of data and metadata. By clicking on each of the buttons in the menu one can choose to view *e.g.*, only the acetate, D-glucose, and O<sub>2</sub> consumption data for the 'WT BN' line. The user can also compare lines by checking them (*e.g.*, 'WT BN' vs 'WT 1A'). See Screenshot 2 in the Supporting Information, or at <https://public-edd.jbei.org/pages/tutorials/>, for a demonstration.





27  
28  
29  
30  
31  
32  
33

Figure 6: **High-level diagram of EDD code structure.** The front-end and visualization run on the client (internet browser) and are coded in TypeScript. The backend involves the database, importer, exporters and the templates and is coded in python using the Django framework.



55  
56  
57  
58  
59  
60

Figure 7: **Database schema for EDD data.** The database is accessed through the Django Object Relational Manager (ORM) and encodes the experiment descriptors shown in Fig. 2.

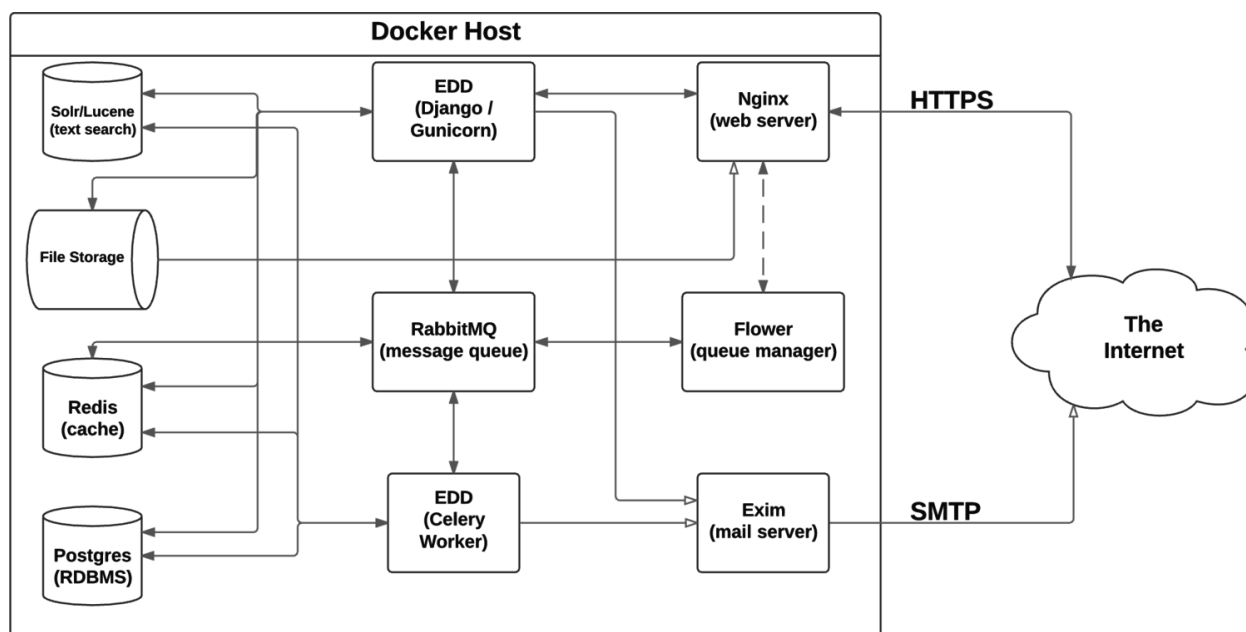


Figure 8: **Service diagram for EDD.** Multiple services combine together to create EDD. This microservice architecture simplifies the process of expanding an installation of EDD.

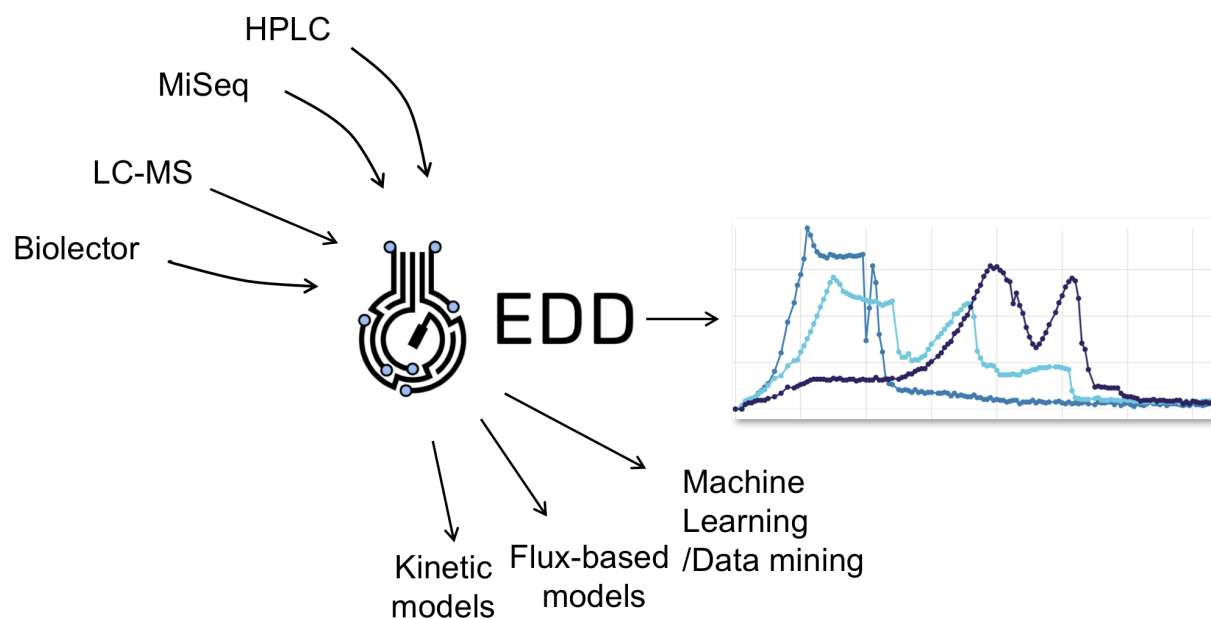


TABLE OF CONTENTS (ToC) graphic.