

# DBgen: A Python Library for Defining Scalable, Maintainable, Accessible, Reconfigurable, Transparent (SMART) Data Pipelines

Michael J. Statt<sup>1</sup>, Kristopher S. Brown<sup>1</sup>, Santosh Suram<sup>2</sup>, Linda Hung<sup>2</sup>,  
Daniel Schweigert<sup>2</sup>, John Gregoire<sup>3</sup>, Brian A. Rohr<sup>1,\*</sup>

---

## Abstract

In this work, we present DBgen, a Python library that provides a framework for defining extract-transform-load (ETL) pipelines to create and populate SQL databases. DBgen is most useful when the underlying data has complex relationships, requires multi-step analysis, is large-scale, and the type of data being collected changes frequently. Scientific data often fits this description. With current tooling, defining ETL pipelines for this particularly difficult-to-manage data is so onerous that a great deal of it does not end up being stored in a database and is opaque. DBgen is designed to fill the gap in the current tooling and reduce the barrier to defining ETL pipelines such data.

*Keywords:* Database, ETL, Python, Data management

---

---

\*Correspondance to:

*Email address:* `brian.rohr@modelyst.io` (Brian A. Rohr)

<sup>1</sup>Modelyst LLC, Palo Alto, CA 94306, USA

<sup>2</sup>Toyota Research Institute, Los Altos, CA 94022, USA

<sup>3</sup>California Institute of Technology, Pasadena, CA, USA

## Required Metadata

### Current code version

Nr.	Code metadata description	Please fill in this column
C1	Current code version	1.0.0
C2	Permanent link to code/repository used for this code version	<a href="https://github.com/modelyst/dbgen">https://github.com/modelyst/dbgen</a> :
C4	Legal Code License	Apache 2.0
C5	Code versioning system used	git
C6	Software code languages, tools, and services used	Python 3 on Linux, OSX or Windows. Dependencies listed in requirements.txt in code repository.
C7	Compilation requirements, operating environments & dependencies	N/A
C8	Link to developer documentation/manual	<a href="https://www.dbgen.modelyst.com">https://www.dbgen.modelyst.com</a>
C9	Support email for questions	brian.rohr@modelyst.io

Table 1: Code metadata

## 1. Motivation and Significance

DBgen is designed to reduce the barrier to creating databases for complicated data sources in accordance with the FAIR data principle[1]. The FAIR principle, which states that data should be Findable, Accessible, Interoperable, and Reusable, is widely accepted and has shown utility in a variety of fields of scientific research including medicine[2, 3, 4], meteorology & oceanography[5, 6, 7, 8], oral speech studies[9], botany[10], mass spectrometry[11], and many others. Although many agree that it is important to make data FAIR, a great deal of scientific data remains stored in a way that does not meet these principles. We contend that this is in part because defining extract-transform-load (ETL) pipelines to get scientific data into SQL databases is a particularly challenging task. Scientific data is complex, analysis-heavy, frequently-changing, large-scale, and the people who understand this particularly difficult-to-manage data best are rarely experts in SQL and data engineering. These challenges regarding making scientific data fair have been highlighted previously [12, 13, 14, 15, 16] as has the importance of directing resources toward solving these infrastructure problems [17]. In order to handle this challenging case, scientific ETL pipelines need to be scalable, maintainable, accessible, reconfigurable, and transparent

20 (SMART). In section 1.1, we describe how the SMART principles address the  
21 challenges in defining scientific ETL pipelines. Then, in section 1.2, we iden-  
22 tify a gap in the current tooling which makes it difficult to define SMART  
23 ETL pipelines. In section 2, we describe what DBgen is and how it helps to  
24 fill this gap in the current tooling.

### 25 1.1. SMART ETL Pipelines

26 *Scalable:* As high-throughput experimentation becomes more widely used,  
27 the scale of scientific data continues to increase. This demands that any sci-  
28 entific ETL pipeline can handle large amounts of data efficiently.

29 *Maintainable:* The fundamental relationships between the entities of in-  
30 terest in scientific data are very complex. There are many of types of entities  
31 that need to be tracked, and the relations between these entities are often-  
32 times many-to-many. There is also a great deal of meta-data that needs to  
33 be stored and formally linked to the data. This makes the FAIR principles  
34 of rich meta-data and strong provenance challenging to achieve. A database  
35 architecture capable of capturing this complexity requires many tables and  
36 foreign keys. Accordingly, the ETL pipeline that populates such a database  
37 is a fairly complex piece of software. As with any complex piece of software,  
38 maintainability is crucial, and making the code modular and easy to debug  
39 is imperative for maintainability.

40 *Accessible:* The number of scientific researchers who are comfortable with  
41 higher-level languages like Python far exceeds the number who are comfort-  
42 able with SQL, so in order to make ETL pipelines editable by scientists, it  
43 is important to abstract away the SQL code.

44 *Reconfigurable:* Scientific ETL pipelines must be easily reconfigurable  
45 because scientific research is inherently frequently-changing. Scientists fre-  
46 quently decide to do new types of experiments and analyses as the very goal  
47 of doing the research is to learn information that changes their understanding  
48 of the subject of research. As they do new types of experiments, they start  
49 to track new types of information and carry out new analyses. Changing  
50 the entities, attributes, and relations that are tracked in a database is the  
51 definition of a schema change. Therefore, schema changes occur far more  
52 often in scientific research than in other fields, and the ETL pipeline needs  
53 to be easy to reconfigure when these inevitable changes occur.

54 *Transparent:* In many cases, researchers are interested in a high-level  
55 analysis of the data, not just the raw data itself. The process of converting  
56 raw data to high-level results usually requires many small steps. Each data  
57 processing step yields an intermediate result, which may be of interest in its  
58 own right or may be useful in other analyses later. When the data pipeline  
59 is complete, there can be a large, complicated web of data processing steps

60 between the raw data and the final results. It is imperative that the full flow  
61 of data from the original, raw, file all the way through to the final analysis  
62 is transparent.

### 63 *1.2. Current Tooling*

64 There are many tools, including DBeaver, TablePlus, and MySQL Work-  
65 bench, that make it easy to define complex empty database schemas; how-  
66 ever, in the case of scientific data, defining the ETL pipeline is much more  
67 difficult than defining the empty database schema. There are several existing  
68 Python packages, including psycpg2, pymysql, and sqlalchemy that expose  
69 SQL functionality to a python user, but they do not provide any framework  
70 for defining SMART ETL pipelines. DBgen fills this gap in the current tool-  
71 ing. If one were to implement 50 complex, SMART ETL pipelines, each for a  
72 different use case, yes, some of the code would specific to each use case, but a  
73 large amount of the code would be common to all of the use cases. DBgen is  
74 that code that is common code. In other words, DBgen is a Python package  
75 that provides a framework for the definition of SMART ETL pipelines, just  
76 as pytorch [18] and tensorflow [19] provide a framework for the definition of  
77 GPU-accelerated, complex deep learning models.

## 78 **2. Software Description**

### 79 *2.1. Software Architecture*

80 Each DBgen Model defines a complete build procedure for a database:  
81 both instantiating the empty database and populating it with data. The  
82 DBgen Model consists of two graphs, one for each of those steps. The schema  
83 graph defines the empty database architecture, and the ETL graph defines  
84 the ETL process, which populates the database with data.

85 In DBgen, there are three key classes associated with defining the schema  
86 graph (Entities, Attributes, and Relations), and four key classes associated  
87 with defining the ETL graph (Generators, Queries, PyBlocks, and Loads).  
88 Both pieces are put together in one, large object, called a Model, which  
89 defines the entire database build procedure, both instantiating the database  
90 and populating it.

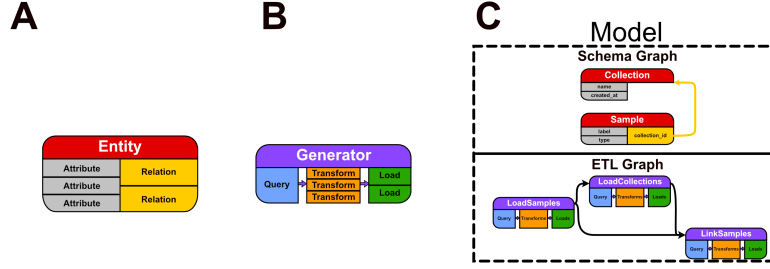


Figure 1: A depiction of the relationships between the key classes in DBgen. (A) shows how the classes that comprise the schema graph relate to each other. (B) shows how the classes that comprise the ETL graph relate to each other. (C) shows how many instances of these classes can be composed to create a DBgen Model.

## 2.2. The Schema Graph: Defining the Empty Database Schema

*Entity:* In DBgen, each Entity fully defines an empty database table. It consists of a name and any number of Attributes and Relations, which are described next. Each entity is a node in the schema graph.

*Attributes:* Attributes, define the columns of each database table. They have a name, a description, and a data type. Attributes in DBgen can be either "identifying" or "non-identifying." The information stored in the identifying columns are necessary and sufficient to identify exactly one row in a given table. DBgen disallows the existence of two rows in the same table with the same identifying data. For example, in a table of movies, one could decide to make the title and release date Attributes identifying. In DBgen, this would guarantee that a query for a specific title and release date would return no more than one row, and DBgen would also require a title and a release date to create a row in the table. There may be many other non-identifying Attributes, like duration and average critic rating. Although the back-end details are different, this concept is analogous to a composite primary key.

*Relations:* Relations define the relationships between the tables and therefore the edges of the schema graph. For those who are familiar with database terminology, a Relation defines a foreign key.

## 2.3. Populating the Database

The ETL graph consists of Generators, which are the nodes in the ETL graph. The information specified in the Query and Loads objects, also described below, allow DBgen to automatically compute the appropriate edges for the ETL graph.

116     *Generators:* Each Generator, defines a single extract-transform-load (ETL)  
117 step. It consists of a Query, a PyBlock, and Loads, which represent the ex-  
118 tract, transform, and load steps, respectively. A common pattern in DBgen  
119 data pipelines is to query the database, use that information as inputs to a  
120 function, and insert the result back into the database. This allows for com-  
121 plicated, multi-step data analyses to be carried out in a flexible, modular,  
122 maintainable, and transparent manner.

123     *Query:* The Query object in DBgen defines a SQL query. By using a  
124 Query object rather than a raw SQL string, the Generator knows which  
125 database columns need to be queried in order for that ETL step to run.  
126 DBgen will later use this information to compute dependencies among the  
127 ETL steps and automatically run them in the correct order.

128     *PyBlock:* The PyBlock object represents the transform portion of the  
129 ETL step. It is any arbitrary python function, which enables sophisticated  
130 analyses, including predictions from machine learning models, to be run. The  
131 inputs to the function come from the Query object, and the outputs from  
132 the function are inserted back into the database in the load step, which is  
133 described next. Pyblocks can also read data in from the file system, which  
134 is commonly used for the early ETL steps in the data pipeline.

135     *Loads:* As the name suggests, Loads represent the load portion of a given  
136 ETL step, which is the step in which data is inserted into the database. In  
137 an instance of the Loads object, the user specifies which outputs from the  
138 PyBlock are inserted into which columns in the database. Now, each Gen  
139 knows which database columns must be populated before it is run (from the  
140 Query object) and also which database columns are populated by that Gen  
141 (from the Loads object). DBgen needs this information to determine the  
142 correct order to run the ETL steps.

143     *Model:* Finally, the Model object defines the entire database build pro-  
144 cedure. It is comprised of a set of Entities and a set of Generators. When  
145 a Model is run, it uses the information in the Entities to create the schema  
146 graph and to instantiate the empty database tables. Then, it uses the infor-  
147 mation from each Generator’s Query and Load objects to compute the order  
148 in which the Gens need to be run, thereby creating the ETL graph. Finally,  
149 it executes each Generator’s Query, Pyblock, and Loads steps to actually  
150 populate the database.

#### 151 2.4. Software Functionalities

152     *Automatic Ordering:* Complicated data pipelines are comprised of a large  
153 number of small, simple processing steps. Without DBgen, the user must  
154 make sure that the ETL steps are run in the correct order. If one ETL step  
155 needs to query a given column in the database, the ETL step that populates

156 that column needs to be run first. This process is laborious, especially in  
157 the case of scientific research, where frequent changes to the data pipeline  
158 are expected. DBgen automatically determines the correct order to run the  
159 ETL steps by using the structured information in each Generator object to  
160 create a directed acyclic graph (DAG). This feature is significant because it  
161 allows the user to not think about the whole data pipeline when adding or  
162 editing an ETL step, even if the change completely disrupts the dependencies  
163 among the ETL steps. This feature contributes to the reconfigurability and  
164 maintainability of DBgen ETL pipelines. Furthermore, this computational  
165 graph can be visualized to show the full flow of data from its source, through  
166 all processing steps, to the final destination in the database. This adds  
167 transparency to all DBgen ETL pipelines.

168 *Primary Key and Foreign Key Handling:* In large, complex database  
169 schemas with many foreign keys, querying tables to populate these foreign  
170 keys properly is both computationally expensive and laborious for the user to  
171 write. DBgen obviates the need for this altogether. The technical details of  
172 how DBgen accomplishes this are described in the supplemental information.

173 *Automatic Detection of Changing Inputs and Functions:* Every time a  
174 DBgen ETL step is run, DBgen stores a hash of the inputs it received and  
175 the function that processed the data. Then, when a pipeline is re-run, DBgen  
176 automatically knows which steps need to be re-run. This avoids the compu-  
177 tational expense of re-running functions, and perhaps more importantly, it  
178 allows the user to add data and make edits to the pipeline without thinking  
179 about the execution of the pipeline at all. If the user wants to change a  
180 function, they just change the function and re-run the DBgen model. DBgen  
181 will automatically detect that only that function was changed, and that ETL  
182 step and all of its dependents will be re-run, and the rest of the steps will be  
183 skipped, as their results are unchanged. This is a recurring theme in DBgen:  
184 separate the definition of the pipeline from the execution of the pipeline, and  
185 abstract away the execution portion, thereby enabling the user to zoom in  
186 and make edits on any small portion of the pipeline without needing to think  
187 at all about how that may impact the broader data pipeline. This contributes  
188 to the reconfigurability and scalability of DBgen ETL pipelines.

189 *DBgen Log Database - A Dashboard for the ETL Process:* The log database  
190 is a separate, small database that is designed to be a dashboard for the data  
191 pipeline that helps the user debug complicated data pipelines with ease. The  
192 log database has one row per ETL step per attempt at running the data  
193 pipeline. The most common use case is to query this table for the most  
194 recent run attempt. For each row, it contains the following information:

- 195
- The status of the ETL step (not started, failed, completed, etc.)

- 196 • The runtime for the ETL step
- 197 • The full traceback for any errors that were encountered when trying to
- 198 run the ETL step
- 199 • Which columns the ETL step queries
- 200 • Which columns the ETL step populates

201 So, if there is an error in one of the functions, and a user tries to run  
 202 the data pipeline, the user can go to the dashboard and immediately see:  
 203 these steps finished, this step encountered an error, all of its dependencies  
 204 did not run as a result, and here is the error message. This makes even the  
 205 most complicated ETL pipelines easy to debug and maintain. Even when  
 206 there are no errors, the user can query the runtime of each step to identify  
 207 bottlenecks in the process. The DBgen log database makes

208 *Efficient, Cloud-Optimized Insertion Methods:* DBgen copies data from  
 209 the machine that is executing the ETL pipeline to the machine that is host-  
 210 ing the database in large batches, which is more computationally efficient  
 211 than traditional methods. Additionally, DBgen makes it easy to run each  
 212 individual ETL step on a separate compute instance, so if one ETL step is  
 213 particularly compute-intensive, it can be run on a more powerful compute  
 214 instance. Together, these features ensure that all DBgen ETL pipelines are  
 215 scalable. The details are described in the supplemental information.

216 *Abstraction:* With DBgen, a user can implement an entire database and  
 217 ETL pipeline without writing one line of SQL code. This provides accessibil-  
 218 ity to ETL pipelines defined using DBgen. With the SQL abstracted away,  
 219 a much broader user base is able to define or at least understand and edit  
 220 existing ETL pipelines.

### 221 3. Illustrative Example

222 The example we will discuss is a database of materials science research  
 223 data created and populated using DBgen called ESAMP, which stands for  
 224 "Event-Sourced Architecture for Materials Provenance." To give a sense of  
 225 the inherent complexity of the problem, in materials science research, a ma-  
 226 terial sample of interest is usually derived from one or many other samples  
 227 through a series of procedures. Each procedure may produce one or many  
 228 data files, which may be used in one or many analyses to yield results that  
 229 are of interest to the researcher. Some analyses depend on other analyses  
 230 having already been carried out. A query against the database should be  
 231 able to answer complicated questions like, "show the top 10 best-performing



232 batteries as determined by a specific lifetime test that were derived from  
233 anodes that are at least 10% cobalt as determined by X-ray photo-electron  
234 spectroscopy (XPS), and exclude any batteries that have the solvent that  
235 arrived on May 18th anywhere in their processing history.”

236 ESAMP models this complexity completely without making any simplify-  
237 ing assumptions. This is significant because it adds transparency, flexibility,  
238 and provenance to the curation of datasets for machine learning (ML) or  
239 other data analysis techniques. When the underlying data is modeled with-  
240 out assumptions, the user is able to write a SQL query with a certain set  
241 of constraints and assumptions to generate a dataset. If the researcher finds  
242 that the data is imbalanced or comes across another problem with their first  
243 dataset, they can easily edit the query to generate a new, improved dataset.

244 Importantly, any future researcher who looks at that dataset and wants  
245 to know exactly where the data came from and what assumptions and con-  
246 straints were applied can simply look at the SQL query that was used to  
247 generate it. This provides much-needed transparency in materials science  
248 ML projects.

249 The underlying data that is now in the ESAMP database was originally  
250 stored in a large zip archive containing additional zip archives containing  
251 tens of thousands of automatically-generated, custom-structured text files.  
252 The data pipeline that extracts the data from that structure and populates  
253 the ESAMP database requires over 50 ETL steps with a complicated tree  
254 of dependencies. The process of designing the ESAMP architecture required  
255 a great deal of iteration, so the ETL pipeline had to be adjusted and redef-  
256 ined dozens of times. Accomplishing this without DBgen would have been  
257 prohibitively laborious.

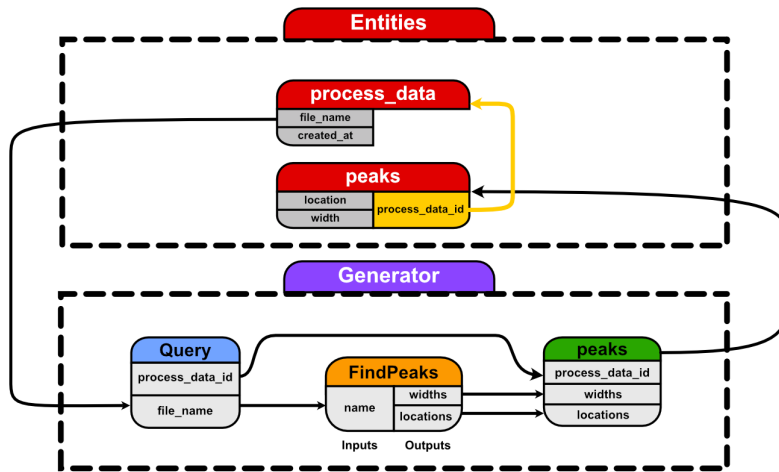


Figure 2: An example of the flow of information within a Generator

Figure 2 illustrates the flow of information within an example Generator. In this example, process data has been stored within the process\_data Entity using its file\_name and the goal of this Generator is to extract each file and find the locations and widths of the peaks within the underlying data. As this generator queries the file\_name from the process\_data Entity and loads the resulting peaks into the peaks entity, DBgen automatically places this Generator after any Generator that populates the process\_data entity and before any Generator that depends on the peaks Entity being populated. Therefore, the author of this Generator can narrow their field of vision to just the data flowing from the data extraction to the loading. As discussed above, DBgen is specifically designed to allow for the authors of the pipeline to easily and effectively change their mind. The three common types of changes made to a pipeline are:

1. New data entering the pipeline
2. New functions to process the data
3. New schema for storing the inputs and outputs

The severity of these changes typically increases from 1 to 3, with an updated schema being the most extreme. However, DBgen helps to greatly reduce the complexity of implementing each type of change.

Firstly, detecting whether a file entering the pipeline has been seen before is critical to reducing the computational strain on the overall pipeline. To avoid reprocessing of duplicate data, this generator will store a hash of the

281 generator and each input in the DBgen Log Database discussed above. This  
 282 enables DBgen to automatically detect whether it has seen an extracted  
 283 input before and, if so, it can skip reprocessing the duplicate row. While this  
 284 may seem trivial for a path to a file, this duplicate detection generalizes to  
 285 complex inputs that are aggregated and processed from many entities within  
 286 the schema.

287 Secondly, scientists regularly update and improve the functions they use  
 288 to transform data. This could be a change to the FindPeaks transform or the  
 289 addition of a pre-processing step to remove the outliers in the data before  
 290 the FindPeaks transform processes the data (as shown in Figure 3). In either  
 291 case, DBgen makes these changes easy to make by automatically ordering the  
 292 functions within a generator. This means that FindPeaks need only request  
 293 an output from the RemoveOutlier for DBgen to know that RemoveOutlier  
 294 needs to be run prior to FindPeaks. Additionally, the generator's hash will  
 295 change with the addition of a transform or the modification of any transforms  
 296 underlying code. This signals to DBgen that each input must be reprocessed  
 297 regardless of whether it had been processed by the previous version of the  
 298 generator.

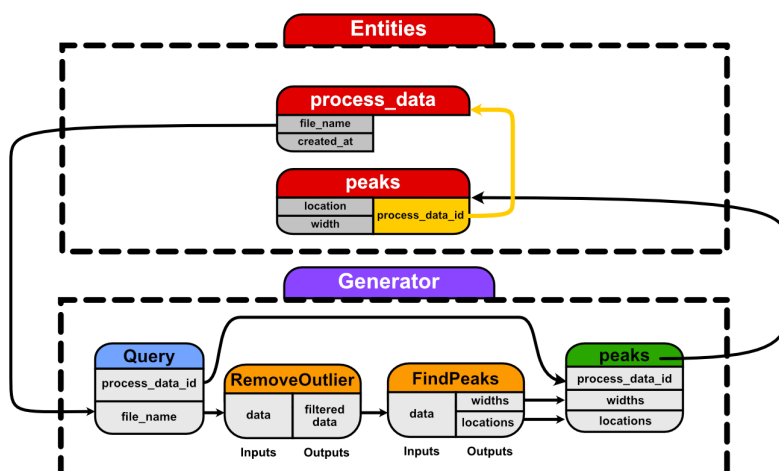


Figure 3: A generator with two PyBlocks, or transform steps.

299 Finally, schema changes can be facilitated by DBgen by isolating the  
 300 schema that stores the data from the generator that modifies it through  
 301 the use clear, well-defined interfaces of Loads and Extracts. When schema  
 302 changes are made, such as the modification of the relationship between pro-  
 303 cess\_data and peaks from a one-to-many to a many-to-many relationship,  
 304 each generator that depends upon or populates the modified entities in the

305 schema is clearly logged in the DBgen Log Database. Additionally, the ab-  
306 stractions of Query objects and Load objects allows for complex schema  
307 changes to be accommodated by only a few lines of code.

## 308 4. Impact

309 DBgen provides a framework that makes it easy to implement SMART  
310 ETL pipelines. We believe this will have the largest impact in the field of  
311 scientific research. Specifically, we believe that DBgen will result in more  
312 scientific data being stored in accordance with the FAIR principles of data  
313 management. Ultimately, this will accelerate innovation in computational  
314 methods applied to experimental scientific data. As a point of analogy, ImageNet [20] was published in 2009 and played a significant role in the develop-  
315 ment of new convolutional neural network architectures shortly thereafter, as  
316 evidenced by the publication of AlexNet in 2012[21], and Resnet in 2016[22].  
317 In many fields of scientific research, there is no analogous database that can  
318 be used to advance the development of computational methods in each field.  
319 This is not because it wouldn't be useful, but rather because it is difficult  
320 to achieve for the reasons mentioned in section 1. We believe that these  
321 databases do not exist today because the software tools needed to generate  
322 these complicated databases and corresponding ETL pipelines does not exist.  
323 DBgen aims to fill that need.

## 325 5. Conclusions

326 In this work, we present DBgen, a python library that adds useful ab-  
327 stractions to the process of defining complex databases and ETL pipelines  
328 and thereby reduces the barrier to storing data in accordance with the FAIR  
329 principle. We also present a set of principles (SMART) that ETL pipelines  
330 should ideally abide by, analogous to the FAIR principles for data storage.  
331 We use materials science R&D data as an example of an inherently complex  
332 data source. We show that modeling the data without making assumptions  
333 demands a complicated database architecture, which would be difficult to  
334 create without DBgen. We show that modeling the data in this way adds  
335 transparency and flexibility to dataset curation. Altogether, we provide evi-  
336 dence that DBgen is a useful tool that greatly reduces the barrier to storing  
337 scientific data in accordance with the widely accepted FAIR principles.

## 338 6. Conflict of Interest

339 Dr. Brian Rohr and Dr. Michael Statt regularly use DBgen in their  
340 materials science data engineering services work at Modelyst, LLC.

## References

- [1] M. D. Wilkinson, M. Dumontier, I. J. Aalbersberg, G. Appleton, M. Axton, A. Baak, N. Blomberg, J.-W. Boiten, L. B. da Silva Santos, P. E. Bourne, et al., The fair guiding principles for scientific data management and stewardship, *Scientific data* 3 (1) (2016) 1–9.
- [2] J. Wise, A. G. de Barron, A. Splendiani, B. Balali-Mood, D. Vasant, E. Little, G. Mellino, I. Harrow, I. Smith, J. Taubert, et al., Implementation and relevance of fair data principles in biopharmaceutical r&d, *Drug discovery today* 24 (4) (2019) 933–938.
- [3] P. Deshpande, A. Rasin, J. Furst, D. Raicu, S. Antani, Diis: a biomedical data access framework for aiding data driven research supporting fair principles, *Data* 4 (2) (2019) 54.
- [4] C. Vesteghem, R. F. Brøndum, M. Sønderkær, M. Sommer, A. Schmitz, J. S. Bødker, K. Dybkær, T. C. El-Galaly, M. Bøgsted, Implementing the fair data principles in precision oncology: review of supporting initiatives, *Briefings in bioinformatics* 21 (3) (2020) 936–945.
- [5] T. Tanhua, S. Pouliquen, J. Hausman, K. O’Brien, P. Bricher, T. de Bruin, J. J. Buck, E. F. Burger, T. Carval, K. S. Casey, et al., Ocean fair data services, *Frontiers in Marine Science* 6 (2019) 440.
- [6] R. Devarakonda, G. Prakash, K. Guntupally, J. Kumar, Big federal data centers implementing fair data principles: Arm data center example, in: 2019 IEEE International Conference on Big Data (Big Data), IEEE, 2019, pp. 6033–6036.
- [7] C. Gries, M. Servilla, M. O’Brien, K. Vanderbilt, C. Smith, D. Costa, S. Grossman-Clarke, Achieving fair data principles at the environmental data initiative, the us-iter data repository, *Biodiversity Information Science and Standards* 3 (2019) e37047.
- [8] L. Lannom, D. Koureas, A. R. Hardisty, Fair data and services in biodiversity science and geoscience, *Data Intelligence* 2 (1-2) (2020) 122–130.
- [9] S. Calamai, F. Frontini, Fair data principles and their application to speech and oral archives, *Journal of New Music Research* 47 (4) (2018) 339–354.
- [10] C. Pommier, C. Michotey, G. Cornut, P. Roumet, E. Duchêne, R. Flores, A. Lebreton, M. Alaux, S. Durand, E. Kimmel, et al., Applying

- 375 fair principles to plant phenotypic data management in gnpis, Plant  
376 Phenomics 2019 (2019).
- 377 [11] Y. Watanabe, K. F. Aoki-Kinoshita, Y. Ishihama, S. Okuda, Glycopost  
378 realizes fair principles for glycomics mass spectrometry data, Nucleic  
379 Acids Research 49 (D1) (2021) D1523–D1528.
- 380 [12] M. Boeckhout, G. A. Zielhuis, A. L. Bredenoord, The fair guiding prin-  
381 ciples for data stewardship: fair enough?, European journal of human  
382 genetics 26 (7) (2018) 931–936.
- 383 [13] C. Draxl, M. Scheffler, Big-data-driven materials science and its fair  
384 data infrastructure, arXiv preprint arXiv:1904.05859 (2019).
- 385 [14] R. David, L. Mabile, A. Specht, S. Stryeck, M. Thomsen, M. Yahia,  
386 C. Jonquet, L. Dollé, D. Jacob, D. Bailo, et al., Fairness literacy: the  
387 achilles’ heel of applying fair principles, CODATA Data Science Journal  
388 19 (32) (2020) 1–11.
- 389 [15] L. B. Da Silva Santos, M. D. Wilkinson, A. Kuzniar, R. Kaliyaperumal,  
390 M. Thompson, M. Dumontier, K. Burger, Fair data points supporting  
391 big data interoperability, Enterprise Interoperability in the Digitized and  
392 Networked Factory of the Future. ISTE, London (2016) 270–279.
- 393 [16] L. Garcia, J. Bolleman, S. Gehant, N. Redaschi, M. Martin, Fair adop-  
394 tion, assessment and challenges at uniprot, Scientific data 6 (1) (2019)  
395 1–4.
- 396 [17] S. Stall, L. Yarmey, J. Cutcher-Gershenfeld, B. Hanson, K. Lehnert,  
397 B. Nosek, M. Parsons, E. Robinson, L. Wyborn, Make scientific data  
398 fair (2019).
- 399 [18] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin,  
400 A. Desmaison, L. Antiga, A. Lerer, Automatic differentiation in pytorch  
401 (2017).
- 402 [19] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S.  
403 Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow,  
404 A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kud-  
405 lur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah,  
406 M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker,

- 407 V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wat-  
408 tenberg, M. Wicke, Y. Yu, X. Zheng, TensorFlow: Large-scale ma-  
409 chine learning on heterogeneous systems, software available from ten-  
410 sorflow.org (2015).  
411 URL <https://www.tensorflow.org/>
- 412 [20] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, L. Fei-Fei, Imagenet: A  
413 large-scale hierarchical image database, in: 2009 IEEE conference on  
414 computer vision and pattern recognition, Ieee, 2009, pp. 248–255.
- 415 [21] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with  
416 deep convolutional neural networks, Advances in neural information pro-  
417 cessing systems 25 (2012) 1097–1105.
- 418 [22] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image  
419 recognition, in: Proceedings of the IEEE conference on computer vision  
420 and pattern recognition, 2016, pp. 770–778.