

Data Modeling and Governance in the Life Sciences

Rancho BioSciences

December 2021



In recent years, through a confluence of factors, we have witnessed a transformation that has led to the exciting era of 'Open Life Sciences Data' and resulted in an explosion in the availability of R&D data public resources.

Some of the factors include:

- 1) impact of the Internet and cloud infrastructure in "leveling the playing field" and increasing accessibility to providers, data, and tools
- 2) cost efficiencies bettering Moore's Law, has made Next Generation Sequencing and other science platforms available at lower price points (we already reached the \$1,000 per genome price point^{1,2})
- 3) organizations have increasingly embraced open data principles over time.
 - a. The Human Genome Project contributed with proactive steps with what has become known as the Bermuda Principles leading to immediate release and publication of sequencing assemblies throughout the project³.
 - b. The Supreme Court Molecular Pathology v. Myriad Genetics, Inc, case ruling that human genes could not be patented additionally led to increased open data⁴.
 - c. Additionally, as access and cost barriers have come down and as data standards have emerged, it has become easier to share data across the world. Many of the largest scientific data repositories now exist in the public domain⁵. Organizations such as the [Broad Institute](#), the [European Bioinformatics Institute \(EMBL-EBI\)](#), the [US National Center for Biotechnology Information \(NCBI\)](#), the [Swiss Institute of Bioinformatics \(SIB\)](#), and many others have significantly contributed to open data principles through provision of excellent resources.
- 4) emerging analytical methods, including Artificial Intelligence and Machine & Deep Learning methods, require consistent high-quality data to deliver more reliable insights.
- 5) better understanding of biology and disease requires coordinated measurement and analysis of genomic, proteomic, cellular imaging, and many other of measures of disease and treatment.

Much of the data collected in the initial life science data wave were collected using FAIR-lite standards – standards that were mostly FAIR but left data scientists and informaticians wanting in various ways. While these resources provided a great leap forward for pan-omics platform data access and availability, data preparation time including collecting, cleaning, and organizing data still accounted for ~80% of a data scientist's total effort spent for data analysis efforts⁶.

¹ <https://www.forbes.com/sites/katiejennings/2020/10/28/how-human-genome-sequencing-went-from-1-billion-a-pop-to-under-1000/?sh=5ea851cd8cea>

² <https://www.genome.gov/about-genomics/fact-sheets/DNA-Sequencing-Costs-Data>

³ https://en.wikipedia.org/wiki/Bermuda_Principles

⁴ https://www.supremecourt.gov/opinions/12pdf/12-398_1b7d.pdf

⁵ https://en.wikipedia.org/wiki/List_of_biological_databases

⁶ <https://www.forbes.com/sites/gilpress/2016/03/23/data-preparation-most-time-consuming-least-enjoyable-data-science-task-survey-says/>

Opportunities: FAIR Data Principles

In 2016, the '[FAIR Guiding Principles for scientific data management and stewardship](#)' were published in *Scientific Data*⁷. The authors intended to provide guidelines to improve the **F**indability, **A**ccessibility, **I**nteroperability, and **R**euse of digital assets. The principles emphasize machine-actionability (i.e., the capacity of computational systems to find, access, interoperate, and reuse data with none or minimal human intervention) because humans increasingly rely on computational support to deal with data as a result of the increase in volume, complexity, and creation speed of data⁸.

FAIR data principles were designed to help people align data, metadata, and infrastructure in a way that is logically searchable and usable:

FAIR Principles

- **Findable**
 - (Meta)data are assigned a globally unique and persistent identifier
 - Data are described with rich metadata
 - Metadata clearly and explicitly include the identifier of the data they describe
 - (Meta)data are registered or indexed in a searchable resource
- **Accessible**
 - (Meta)data are retrievable by their identifier using a standardized communications protocol
 - The protocol is open, free, and universally implementable
 - The protocol allows for an authentication and authorization procedure, where necessary
 - Metadata are accessible, even when the data are no longer available
- **Interoperable**
 - (Meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation
 - (Meta)data use vocabularies that follow FAIR principles
 - (Meta)data include qualified reference to other (meta)data
- **Reusable**
 - (Meta)data are richly described with a plurality of accurate and relevant attributes
 - (Meta)data are associated with detailed provenance
 - (Meta)data meet domain-relevant community standards

 www.RanchoBioSciences.com

Figure 1: FAIR Data Guiding Principles

Catching up with FAIRification and FAIR Data Standards

While FAIR data principles are designed to streamline access to R&D data, fully incorporating these principles across an organization is challenging. Multiple legacy data systems have often been built to local requirements, often for identical data types. Data systems are not designed to support accessible and interoperable data. Data and information provenance across internal systems is often not available. Interoperability with public data is complicated at best and often requires manual assessment and curation. Finally, the cost to retrofit existing solutions or implement new solutions that fully enable FAIR data standards is high and often not figured into IT solution budgets.

In short, most of the challenges faced by organizations looking to FAIRify their data and solutions, revolve around the absence or lack of completeness of broader, principled data modeling and governance practices. Implementing FAIR data practices across an organization requires significant change management for both IT and business practices. It requires substantive investment in data stewardship and new solutions. Most of all, implementing FAIR data practices requires robust data modeling and governance.

⁷ <https://www.nature.com/articles/sdata201618>

⁸ <https://www.go-fair.org/fair-principles/>

Driving Towards a Solution

Master Conceptual Model

The Master Conceptual Model is developed with the client to provide an authoritative overview of their overall business processes. It is typically a top-down modeling approach meant to support maximal breadth of organizational operations and goals. The master conceptual model is high-level, attempting to balance the twin requirements of being general enough to encompass all data domains while specific enough to be useful.

In the example shown in Figure 2, the Biologics Drug Discovery, Development, and Deployment Map (4DM) business process map from the NIH National Center for Advancing Translational Sciences is presented^{9, 10}.

This business process map can, if developed through extensive business analysis efforts, be straightforwardly used to develop a master conceptual model. An example of key entities that support basic science research and target identification (Figure 2, top left) is presented in Figure 3. The master conceptual model shows key entities and relationships amongst them for capturing data and information related to animal model and cell lines studies. The same basic construct, because of the model's generality, can be used to describe human clinical trial studies required to support the clinical research and development business process shown at the middle left of Figure 2.

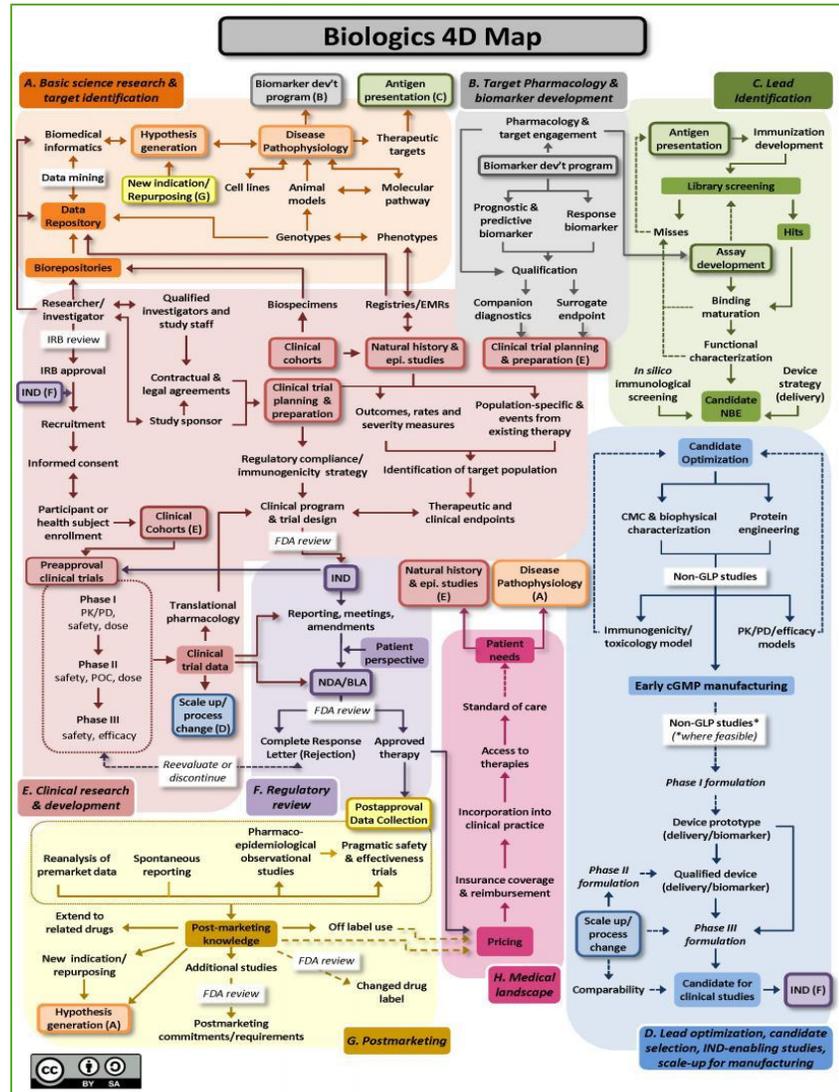


Figure 2: Biologics 4D (Drug Discovery, Development and Deployment) Map showing the basic entities and flows through the biologics therapeutic development process.

⁹ <https://ncats.nih.gov/translation/maps>

¹⁰ <https://pubmed.ncbi.nlm.nih.gov/29269942/>

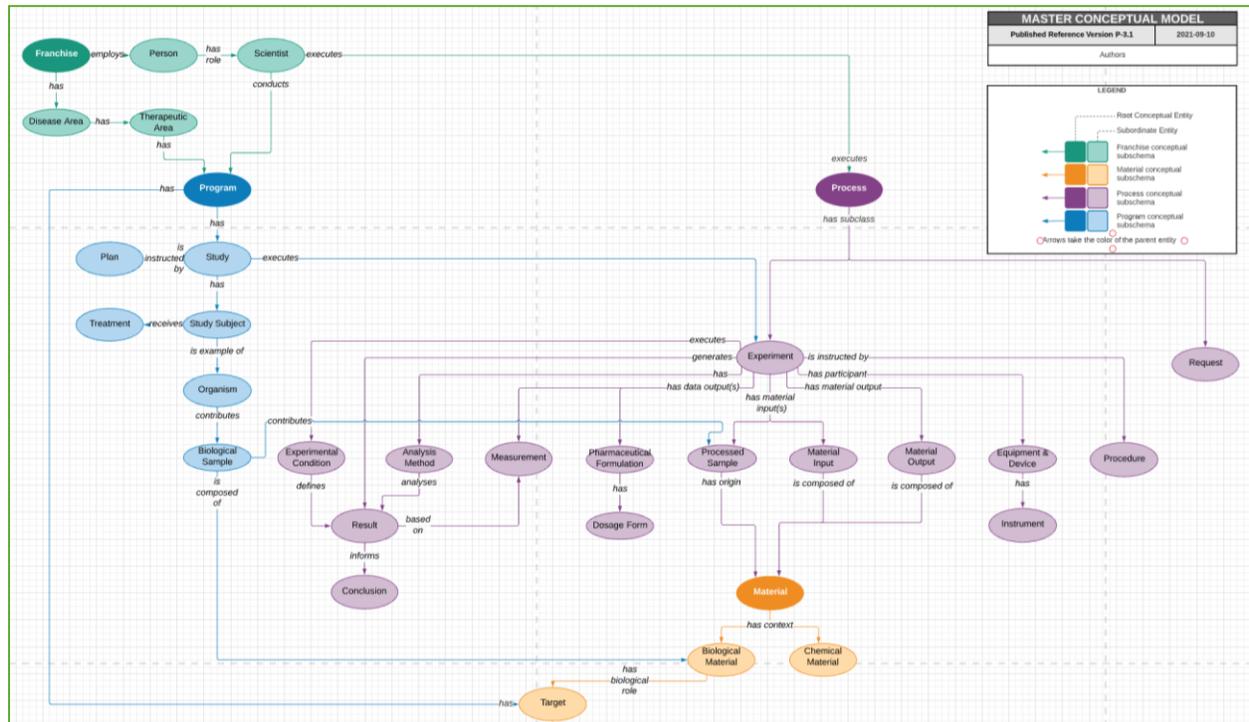


Figure 3: Master Conceptual Model for a basic Life Sciences Study and Experiment construct.

Contextual Models

Development and maintenance of Contextual Models follows the development of the master conceptual model. These contextual models provide clients with a framework for consistency through a detailed description of increasingly specific business practices. Contextual models inherit their base structure from the master conceptual model and allow for greater definition and differentiation of distinct domains, such as the research and clinical domains depicted in Figure 2. The entities in Figure 3 are more completely defined with entity subtypes and attribution. This continues through to the development of well-defined solution logical models.

The collection of the master conceptual model, business area contextual models, and solution logical models provides a roadmap of an organization’s data and information. Combined with external data models such as FDAs Electronic Common Technical Document (eCTD)¹¹, the Observational Medical Outcomes Partnership (OMOP) common data model¹², the NIH Genomics Data Commons (GDC) Data Model¹³, and others, an organization has not only their own data and information roadmaps, but also those for public resources. Alignment of internal and external models allows for simpler data exchange and direct data use operations. We are presently leveraging existing models, solutions, and application programming interfaces (APIs) and extending best practices where needed to enable our client’s solutions and their access to public data resources.

Ontologies, Taxonomies, and Master Dictionaries and Controlled Vocabularies

Fully documented model references, and controlled vocabularies, ontologies, and taxonomies that describe the models must be leveraged if already existing, developed where they are not, and maintained. These ontologies, taxonomies and vocabularies are what fill the instantiated solutions with metadata and data. Utilizing best practices for selection, use, and maintenance of ontologies such as those presented by Malone, *et al.*, with *Ten Simple Rules for Selecting a Bio-ontology*¹⁴ is critical for robust model definitions and data descriptions as described by FAIR.

In addition, utilizing public ontology resources enables use of important ontology class mapping resources. For example, SNOMED CT, ICD-10, and the NCI thesaurus (NCIt) all provide ontologies and taxonomies for important clinical entities and attributes. Each resource provides some support for cross-referencing or mapping of their terms to those within other ontologies. In this way, if one organization has chosen ICD-10 as a primary source and another organization has chosen NCIt, the cross-references provided by each, create a data translation layer that facilitates data management and interoperability requirements when data from disparate solutions are brought together.

```

1  {
2  "tables": [
3  {
4    "table_name": "master_model_sample",
5    "primary_key_field_name": "sample_identifier",
6    "table_fields": [
7    {
8      "field_name": "sample_identifier",
9      "ref_table": "master_model_subject",
10     "ref_table_key": "subject_identifier",
11     "ref_type": "Many-to-one",
12     "field_type": "NUMBER",
13     "field_length": 16,
14     "field_definition": "unique database identifier",
15     "field_source": "database sequence",
16     "unique": true,
17     "nullable": false
18     },
19     ],
20   },
21   {
22     "field_name": "sample_number",
23     "ref_table": "",
24     "ref_table_key": "",
25     "ref_type": "",
26     "field_type": "VARCHAR2",
27     "field_length": 512,
28     "field_definition": "sample identification number - typically from a sample barcode",
29     "field_source": "text",
30     "unique": false,
31     "nullable": false
32     },
33   ],
34   {
35     "field_name": "sample_mass",
36     "ref_table": "",
37     "ref_table_key": "",
38     "ref_type": "",
39     "field_type": "BINARY_DOUBLE",
40     "field_length": 4,
41     "field_definition": "mass numeric value - paired with sample_mass_units",
42     "field_source": "text",
43     "unique": false,
44     "nullable": true
45     },
46   ],
47   {
48     "field_name": "sample_mass_units",
49     "ref_table": "",
50     "ref_table_key": "",
51     "ref_type": "",
52     "field_type": "TEXT",
53     "field_length": 65535,
54     "field_definition": "mass numeric value - paired with sample_mass",
55     "field_source": "Units of measurement ontology (UO): http://purl.obolibrary.org/obo/0000000",
56     "unique": false,
57     "nullable": true
58     },
59   ],
60   ],
61   },
62   ],
63   },
64   ],
65   },
66   ],
67   },
68   ],
69   },
70   ],
71   },
72   ],
73   },
74   ],
75   },
76   ],
77   },
78   ],
79   },
80   ],
81   },
82   ],
83   },
84   ],
85   },
86   ],
87   },
88   ],
89   },
90   ],
91   },
92   ],
93   },
94   ],
95   },
96   ],
97   },
98   ],
99   },
100  ],
101  },
102  ],
103  },
104  ],
105  ],
106  ],
107  ],
108  ],
109  ],
110  ],
111  ],
112  ],
113  ],
114  ],
115  ],
116  ],
117  ],
118  ],
119  ],
120  ],
121  ],
122  ],
123  ],
124  ],
125  ],
126  ],
127  ],
128  ],
129  ],
130  ],
131  ],
132  ],
133  ],
134  ],
135  ],
136  ],
137  ],
138  ],
139  ],
140  ],
141  ],
142  ],
143  ],
144  ],
145  ],
146  ],
147  ],
148  ],
149  ],
150  ],
151  ],
152  ],
153  ],
154  ],
155  ],
156  ],
157  ],
158  ],
159  ],
160  ],
161  ],
162  ],
163  ],
164  ],
165  ],
166  ],
167  ],
168  ],
169  ],
170  ],
171  ],
172  ],
173  ],
174  ],
175  ],
176  ],
177  ],
178  ],
179  ],
180  ],
181  ],
182  ],
183  ],
184  ],
185  ],
186  ],
187  ],
188  ],
189  ],
190  ],
191  ],
192  ],
193  ],
194  ],
195  ],
196  ],
197  ],
198  ],
199  ],
200  ],
201  ],
202  ],
203  ],
204  ],
205  ],
206  ],
207  ],
208  ],
209  ],
210  ],
211  ],
212  ],
213  ],
214  ],
215  ],
216  ],
217  ],
218  ],
219  ],
220  ],
221  ],
222  ],
223  ],
224  ],
225  ],
226  ],
227  ],
228  ],
229  ],
230  ],
231  ],
232  ],
233  ],
234  ],
235  ],
236  ],
237  ],
238  ],
239  ],
240  ],
241  ],
242  ],
243  ],
244  ],
245  ],
246  ],
247  ],
248  ],
249  ],
250  ],
251  ],
252  ],
253  ],
254  ],
255  ],
256  ],
257  ],
258  ],
259  ],
260  ],
261  ],
262  ],
263  ],
264  ],
265  ],
266  ],
267  ],
268  ],
269  ],
270  ],
271  ],
272  ],
273  ],
274  ],
275  ],
276  ],
277  ],
278  ],
279  ],
280  ],
281  ],
282  ],
283  ],
284  ],
285  ],
286  ],
287  ],
288  ],
289  ],
290  ],
291  ],
292  ],
293  ],
294  ],
295  ],
296  ],
297  ],
298  ],
299  ],
300  ],
301  ],
302  ],
303  ],
304  ],
305  ],
306  ],
307  ],
308  ],
309  ],
310  ],
311  ],
312  ],
313  ],
314  ],
315  ],
316  ],
317  ],
318  ],
319  ],
320  ],
321  ],
322  ],
323  ],
324  ],
325  ],
326  ],
327  ],
328  ],
329  ],
330  ],
331  ],
332  ],
333  ],
334  ],
335  ],
336  ],
337  ],
338  ],
339  ],
340  ],
341  ],
342  ],
343  ],
344  ],
345  ],
346  ],
347  ],
348  ],
349  ],
350  ],
351  ],
352  ],
353  ],
354  ],
355  ],
356  ],
357  ],
358  ],
359  ],
360  ],
361  ],
362  ],
363  ],
364  ],
365  ],
366  ],
367  ],
368  ],
369  ],
370  ],
371  ],
372  ],
373  ],
374  ],
375  ],
376  ],
377  ],
378  ],
379  ],
380  ],
381  ],
382  ],
383  ],
384  ],
385  ],
386  ],
387  ],
388  ],
389  ],
390  ],
391  ],
392  ],
393  ],
394  ],
395  ],
396  ],
397  ],
398  ],
399  ],
400  ],
401  ],
402  ],
403  ],
404  ],
405  ],
406  ],
407  ],
408  ],
409  ],
410  ],
411  ],
412  ],
413  ],
414  ],
415  ],
416  ],
417  ],
418  ],
419  ],
420  ],
421  ],
422  ],
423  ],
424  ],
425  ],
426  ],
427  ],
428  ],
429  ],
430  ],
431  ],
432  ],
433  ],
434  ],
435  ],
436  ],
437  ],
438  ],
439  ],
440  ],
441  ],
442  ],
443  ],
444  ],
445  ],
446  ],
447  ],
448  ],
449  ],
450  ],
451  ],
452  ],
453  ],
454  ],
455  ],
456  ],
457  ],
458  ],
459  ],
460  ],
461  ],
462  ],
463  ],
464  ],
465  ],
466  ],
467  ],
468  ],
469  ],
470  ],
471  ],
472  ],
473  ],
474  ],
475  ],
476  ],
477  ],
478  ],
479  ],
480  ],
481  ],
482  ],
483  ],
484  ],
485  ],
486  ],
487  ],
488  ],
489  ],
490  ],
491  ],
492  ],
493  ],
494  ],
495  ],
496  ],
497  ],
498  ],
499  ],
500  ],
501  ],
502  ],
503  ],
504  ],
505  ],
506  ],
507  ],
508  ],
509  ],
510  ],
511  ],
512  ],
513  ],
514  ],
515  ],
516  ],
517  ],
518  ],
519  ],
520  ],
521  ],
522  ],
523  ],
524  ],
525  ],
526  ],
527  ],
528  ],
529  ],
530  ],
531  ],
532  ],
533  ],
534  ],
535  ],
536  ],
537  ],
538  ],
539  ],
540  ],
541  ],
542  ],
543  ],
544  ],
545  ],
546  ],
547  ],
548  ],
549  ],
550  ],
551  ],
552  ],
553  ],
554  ],
555  ],
556  ],
557  ],
558  ],
559  ],
560  ],
561  ],
562  ],
563  ],
564  ],
565  ],
566  ],
567  ],
568  ],
569  ],
570  ],
571  ],
572  ],
573  ],
574  ],
575  ],
576  ],
577  ],
578  ],
579  ],
580  ],
581  ],
582  ],
583  ],
584  ],
585  ],
586  ],
587  ],
588  ],
589  ],
590  ],
591  ],
592  ],
593  ],
594  ],
595  ],
596  ],
597  ],
598  ],
599  ],
600  ],
601  ],
602  ],
603  ],
604  ],
605  ],
606  ],
607  ],
608  ],
609  ],
610  ],
611  ],
612  ],
613  ],
614  ],
615  ],
616  ],
617  ],
618  ],
619  ],
620  ],
621  ],
622  ],
623  ],
624  ],
625  ],
626  ],
627  ],
628  ],
629  ],
630  ],
631  ],
632  ],
633  ],
634  ],
635  ],
636  ],
637  ],
638  ],
639  ],
640  ],
641  ],
642  ],
643  ],
644  ],
645  ],
646  ],
647  ],
648  ],
649  ],
650  ],
651  ],
652  ],
653  ],
654  ],
655  ],
656  ],
657  ],
658  ],
659  ],
660  ],
661  ],
662  ],
663  ],
664  ],
665  ],
666  ],
667  ],
668  ],
669  ],
670  ],
671  ],
672  ],
673  ],
674  ],
675  ],
676  ],
677  ],
678  ],
679  ],
680  ],
681  ],
682  ],
683  ],
684  ],
685  ],
686  ],
687  ],
688  ],
689  ],
690  ],
691  ],
692  ],
693  ],
694  ],
695  ],
696  ],
697  ],
698  ],
699  ],
700  ],
701  ],
702  ],
703  ],
704  ],
705  ],
706  ],
707  ],
708  ],
709  ],
710  ],
711  ],
712  ],
713  ],
714  ],
715  ],
716  ],
717  ],
718  ],
719  ],
720  ],
721  ],
722  ],
723  ],
724  ],
725  ],
726  ],
727  ],
728  ],
729  ],
730  ],
731  ],
732  ],
733  ],
734  ],
735  ],
736  ],
737  ],
738  ],
739  ],
740  ],
741  ],
742  ],
743  ],
744  ],
745  ],
746  ],
747  ],
748  ],
749  ],
750  ],
751  ],
752  ],
753  ],
754  ],
755  ],
756  ],
757  ],
758  ],
759  ],
760  ],
761  ],
762  ],
763  ],
764  ],
765  ],
766  ],
767  ],
768  ],
769  ],
770  ],
771  ],
772  ],
773  ],
774  ],
775  ],
776  ],
777  ],
778  ],
779  ],
780  ],
781  ],
782  ],
783  ],
784  ],
785  ],
786  ],
787  ],
788  ],
789  ],
790  ],
791  ],
792  ],
793  ],
794  ],
795  ],
796  ],
797  ],
798  ],
799  ],
800  ],
801  ],
802  ],
803  ],
804  ],
805  ],
806  ],
807  ],
808  ],
809  ],
810  ],
811  ],
812  ],
813  ],
814  ],
815  ],
816  ],
817  ],
818  ],
819  ],
820  ],
821  ],
822  ],
823  ],
824  ],
825  ],
826  ],
827  ],
828  ],
829  ],
830  ],
831  ],
832  ],
833  ],
834  ],
835  ],
836  ],
837  ],
838  ],
839  ],
840  ],
841  ],
842  ],
843  ],
844  ],
845  ],
846  ],
847  ],
848  ],
849  ],
850  ],
851  ],
852  ],
853  ],
854  ],
855  ],
856  ],
857  ],
858  ],
859  ],
860  ],
861  ],
862  ],
863  ],
864  ],
865  ],
866  ],
867  ],
868  ],
869  ],
870  ],
871  ],
872  ],
873  ],
874  ],
875  ],
876  ],
877  ],
878  ],
879  ],
880  ],
881  ],
882  ],
883  ],
884  ],
885  ],
886  ],
887  ],
888  ],
889  ],
890  ],
891  ],
892  ],
893  ],
894  ],
895  ],
896  ],
897  ],
898  ],
899  ],
900  ],
901  ],
902  ],
903  ],
904  ],
905  ],
906  ],
907  ],
908  ],
909  ],
910  ],
911  ],
912  ],
913  ],
914  ],
915  ],
916  ],
917  ],
918  ],
919  ],
920  ],
921  ],
922  ],
923  ],
924  ],
925  ],
926  ],
927  ],
928  ],
929  ],
930  ],
931  ],
932  ],
933  ],
934  ],
935  ],
936  ],
937  ],
938  ],
939  ],
940  ],
941  ],
942  ],
943  ],
944  ],
945  ],
946  ],
947  ],
948  ],
949  ],
950  ],
951  ],
952  ],
953  ],
954  ],
955  ],
956  ],
957  ],
958  ],
959  ],
960  ],
961  ],
962  ],
963  ],
964  ],
965  ],
966  ],
967  ],
968  ],
969  ],
970  ],
971  ],
972  ],
973  ],
974  ],
975  ],
976  ],
977  ],
978  ],
979  ],
980  ],
981  ],
982  ],
983  ],
984  ],
985  ],
986  ],
987  ],
988  ],
989  ],
990  ],
991  ],
992  ],
993  ],
994  ],
995  ],
996  ],
997  ],
998  ],
999  ],
1000 ]

```

Figure 4: A programmatically accessible data structure based on a contextual data model.

¹¹ <https://www.fda.gov/drugs/electronic-submissions-cder/electronic-common-technical-document-ectd>

¹² <https://www.ohdsi.org/data-standardization/the-common-data-model/>

¹³ <https://gdc.cancer.gov/developers/gdc-data-model>

¹⁴ <https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1004743>

Data and Solution Governance

Simply put, consistent governance enables the implementation and maintenance of data and information practices and solutions such as those described above. By thoughtfully linking models and dictionaries to metadata, data, and business processes, an organization lays a solid foundation through FAIR data practices to meet business and regulatory requirements.

Rancho BioSciences provides its clients with FAIR data governance capabilities, providing data modeling services, ontology management, and solution development capabilities supported by documented procedures and policies. Our scientists, business analysts, ontologists, and solution developers work with you to understand your FAIR data requirements and partner with you to provide the best solutions implementing across your private solution space, with vendor and open-source solutions, and through the development and implementation of custom-fit FAIR solutions by facilitating data accessibility and availability, data usability, data integrity, and data security.

