# Interoperability and FAIRness through a novel combination of Web technologies

Mark D Wilkinson [Corresp., 1] , Ruben Verborgh [2] , Luiz Olavo Bonino da Silva Santos [3] , Tim Clark [4] , Morris A Swertz [5] , Fleur D.L. Kelpin [5] , Alasdair J. G. Gray [6] , Erik A. Schultes [7] , Erik M. van Mulligen [8] , Paolo Ciccarese [9] , Mark Thompson [7] , Rajaram Kaliyaperumal [7] , Jerven T. Bolleman [10] , Michel Dumontier [11]

[1] Center for Plant Biotechnology and Genomics - UPM/INIA, Universidad Politécnica de Madrid, Madrid, Spain

[2] Interuniversity Microelectronics Centre (IMEC), Ghent University, Ghent, Belgium

[3] Dutch Techcenter for Life Sciences, Utrecht, The Netherlands

[4] Department of Neurology, Massachusetts General Hospital, Boston, United States of America

[5] Genomics Coordination Center and Department of Genetics, University Medical Center Groningen, Groningen, The Netherlands

[6] Department of Computer Science, School of Mathematical and Computer Sciences, Heriot-Watt University, Edinburgh, United Kingdom

[7] Department of Human Genetics,, Leiden University Medical Center, Leiden, The Netherlands

[8] Department of Medical Informatics, Erasmus University Medical Center, Rotterdam, The Netherlands

[9] Elmer Innovation Lab, Harvard Medical School, Boston, United States of America

[10] Swiss-Prot Group, SIB Swiss Institute of Bioinformatics, Centre Medical Universitaire, Geneva, Switzerland

[11] Stanford Center for Biomedical Informatics Research, Stanford University School of Medicine, Stanford, California, United States of America

Corresponding Author: Mark D Wilkinson
Email address: markw@illuminae.com

Data in the life sciences are extremely diverse and are stored in a broad spectrum of repositories ranging from those designed for particular data types (such as KEGG for pathway data or UniProt for protein data) to those that are general-purpose (such as FigShare, Zenodo, or EUDat). These data have widely different levels of sensitivity and security considerations. For example, clinical observations about genetic mutations in patients are highly sensitive, while observations of species diversity are generally not. The lack of uniformity in data models from one repository to another, and in the richness and availability of metadata descriptions, makes integration and analysis of these data a manual, time-consuming task with no scalability. Here we explore a set of resource-oriented Web design patterns for data discovery, accessibility, transformation, and integration that can be implemented by any general- or special-purpose repository as a means to assist users in finding and reusing their data holdings. We show that by using off-the-shelf technologies, interoperability can be achieved even to the level of an individual spreadsheet cell. We note that the behaviors of this architecture compare favorably to the desiderata defined by the FAIR Data Principles, and can therefore represent an exemplar implementation of those principles. The proposed interoperability design patterns may be used to improve discovery and integration of both new and legacy data, maximizing the utility of all scholarly outputs.

1

2

# Interoperability and FAIRness through a novel combination of Web technologies

5

## Authors:

7

**Mark D. Wilkinson** - Center for Plant Biotechnology and Genomics, UPM-INIA, Madrid, Spain

**Ruben Verborgh** – Ghent University – IMEC, Ghent, Belgium

**Luiz Olavo Bonino da Silva Santos** - Dutch Techcentre for Life Sciences, Utrecht, The Netherlands - Vrije Universiteit Amsterdam, Amsterdam, The Netherlands

**Tim Clark** - Department of Neurology, Massachusetts General Hospital Boston MA and Harvard Medical School, Boston, MA, USA

**Morris A. Swertz** - Genomics Coordination Center and Department of Genetics, University Medical Center Groningen, Groningen, The Netherlands

**Fleur D.L. Kelpin** - Genomics Coordination Center and Department of Genetics, University Medical Center Groningen, Groningen, The Netherlands

**Alasdair J. G. Gray** - Department of Computer Science, School of Mathematical and Computer Sciences, Heriot-Watt University, Edinburgh, UK

**Erik A. Schultes** - Department of Human Genetics, Leiden University Medical Center, Leiden, The Netherlands

**Erik M. van Mulligen** - Department of Medical Informatics, Erasmus University Medical Center Rotterdam, The Netherlands

**Paolo Ciccarese** - Perkin Elmer Innovation Lab, Cambridge MA and Harvard Medical School, Boston MA, USA

**Mark Thompson** - Leiden University Medical Center, Leiden, The Netherlands

**Rajaram Kaliyaperumal** - Leiden University Medical Center, Leiden, The Netherlands

**Jerven T. Bolleman** - Swiss-Prot group, SIB Swiss Institute of Bioinformatics, Centre Medical Universitaire, Geneva, Switzerland

**Michel Dumontier** - Stanford Center for Biomedical Informatics Research, Stanford University, Stanford, California

33

## Corresponding Author:

35

Mark D. Wilkinson

*markw@illuminae.com*, +34 622 784 026

38

39
40
41

42
43

## Abstract

45  Data in the life sciences are extremely diverse and are stored in a broad spectrum of
46  repositories ranging from those designed for particular data types (such as KEGG for
47  pathway data or UniProt for protein data) to those that are general-purpose (such as
48  FigShare, Zenodo, or EUDat). These data have widely different levels of sensitivity and
49  security considerations.  For example, clinical observations about genetic mutations in
50  patients are highly sensitive, while observations of species diversity are generally not. The
51  lack of uniformity in data models from one repository to another, and in the richness and
52  availability of metadata descriptions, makes integration and analysis of these data a manual,
53  time-consuming task with no scalability.  Here we explore a set of resource-oriented Web
54  design patterns for data discovery, accessibility, transformation, and integration that can be
55  implemented by any general- or special-purpose repository as a means to assist users in
56  finding and reusing their data holdings. We show that by using off-the-shelf technologies,
57  interoperability can be achieved even to the level of an individual spreadsheet cell. We note
58  that the behaviors of this architecture compare favorably to the desiderata defined by the
59  FAIR Data Principles, and can therefore represent an exemplar implementation of those
60  principles. The proposed interoperability design patterns may be used to improve discovery
61  and integration of both new and legacy data, maximizing the utility of all scholarly outputs.
62
63

64

## Introduction

66  Carefully-generated data are the foundation for scientific conclusions, new hypotheses,
67  discourse, disagreement and resolution of these disagreements, all of which drive scientific
68  discovery. Data must therefore be considered, and treated, as first-order scientific output,
69  upon which there may be many downstream derivative works, among these, the familiar
70  research article (Starr et al., 2015). But as the volume and complexity of data continue to
71  grow,  a data publication and distribution infrastructure is beginning to emerge that is not *ad*
72  *hoc*, but rather explicitly designed to support discovery, accessibility, (re)coding to
73  standards, integration, machine-guided interpretation, and re-use.
74
75  In this text, we use the word "data" to mean all digital research artefacts, whether they be
76  data (in the traditional sense), research-oriented digital objects such as workflows, or
77  combinations/packages of these (i.e. the concept of a "research object", (Bechhofer et al.,
78  2013)). Effectively, all digital entities in the research data ecosystem will be considered data
79  by this manuscript. Further, we intend "data" to include both data and metadata, and
80  recognize that the distinction between the two is often user-dependent.  Data, of all types,
81  are often published online, where the practice of open data publication is being encouraged
82  by the scholarly community, and increasingly adopted as a requirement of funding agencies
83  (Stein et al., 2015).  Such publications utilize either a special-purpose repository (e.g. model-
84  organism or molecular data repositories) or increasingly commonly will utilize general-
85  purpose repositories such as FigShare, Zenodo, Dataverse, or even institutional
86  repositories.  Special-purpose repositories generally receive dedicated funding to curate and
87  organize data, and have specific query interfaces and APIs to enable exploration of their
88  content.  General-purpose repositories, on the other hand, allow publication of data in
89  arbitrary formats, with little or no curation and often very little structured metadata.  Both of
90  these scenarios pose a problem with respect to interoperability.  While APIs allow
91  mechanized access to the data holdings of a special-purpose repository, each repository has
92  its own API, thus requiring specialized software to be created for each cross-repository
93  query.  Moreover, the ontological basis of the curated annotations are not always
94  transparent (neither to humans nor machines), which thwarts automated integration.
95  General purpose repositories are less likely to have rich APIs, thus often requiring manual
96  discovery and download; however, more importantly, the frequent lack of harmonization of
97  the file types/formats and coding systems in the repository, and lack of curation, results in
98  much of their content being unusable (Roche et al., 2015).
99
100  There are many stakeholders in this endeavour. Scientists themselves, acting as both
101  producers and consumers of these public and private data; public and private research-
102  oriented agencies; journals and professional data publishers both "general purpose" and
103  "special purpose"; research funders who have paid for the underlying research to be
104  conducted; data centres (e.g. the EBI (Cook et al., 2016),  and the SIB (SIB Swiss Institute
105  of Bioinformatics Members, 2016)) who curate and host these data on behalf of the research
106  community; research infrastructures such as BBMRI-ERIC (van Ommen et al., 2015) and

107     ELIXIR (Crosswell & Thornton, 2012), and diverse others. All of these stakeholders have
108     distinct needs with respect to the behaviors of the scholarly data infrastructure. Scientists, for
109     example, need to access research datasets in order to initiate integrative analyses, while
110     funding agencies and review panels may be more interested in the metadata associated with
111     a data deposition - for example, the number of views or downloads, and the selected license.
112     Due to the diversity of stakeholders; the size, nature/format, and distribution of data assets;
113     the need to support freedom-of-choice of all stakeholders; respect for privacy;
114     acknowledgment of data ownership; and recognition of the limited resources available to
115     both data producers and data hosts, we see this endeavour as one of the *Grand Challenges*
116     *of eScience*.
117
118     In January 2014, representatives of a range of stakeholders came together at the request of
119     the Netherlands eScience Center and the Dutch Techcentre for Life Sciences (DTL) at the
120     Lorentz Center in Leiden, the Netherlands, to brainstorm and debate about how to further
121     enhance infrastructures to support a data ecosystem for eScience. From these discussions
122     emerged the notion that the definition and widespread support of a minimal set of
123     community-agreed guiding principles and practices could enable data providers and
124     consumers - machines and humans alike - to more easily find, access, interoperate, and
125     sensibly re-use the vast quantities of information being generated by contemporary data-
126     intensive science. These principles and practices should enable a broad range of integrative
127     and exploratory behaviours, and support a wide range of technology choices and
128     implementations, just as the Internet Protocol (IP) provides a minimal layer that enables the
129     creation of a vast array of data provision, consumption, and visualisation tools on the
130     Internet. The main outcome of the workshop was the definition of the so-called FAIR guiding
131     principles aimed at publishing data in a format that is **Findable**, **Accessible**, **Interoperable**
132     and **Reusable** by both machines and human users. The FAIR Principles underwent a period
133     of public discussion and elaboration, and were recently published (Wilkinson et al., 2016).
134     Briefly, the principles state:
135
136
137         **Findable** - data should be identified using globally unique, resolvable, and persistent
138         identifiers, and should include machine-actionable contextual information that can be
139         indexed to support human and machine discovery of that data.
140
141         **Accessible** - identified data should be accessible, optimally by both humans and
142         machines, using a clearly-defined protocol and, if necessary, with clearly-defined
143         rules for authorization/authentication.
144
145         **Interoperable** - data becomes interoperable when it is machine-actionable, using
146         shared vocabularies and/or ontologies, inside of a syntactically and semantically
147         machine-accessible format.
148
149         **Reusable** - Reusable data will first be compliant with the F, A, and I principles, but
150         further, will be sufficiently well-described with, for example, contextual information, so

151         it can be accurately linked or integrated, like-with-like, with other data sources.

152         Moreover, there should be sufficiently rich provenance information so reused data

153         can be properly cited.

154

155

156 Here we describe a novel interoperability architecture that combines three pre-existing Web

157 technologies and standards to enhance the discovery, integration, and reuse of data in

158 repositories that lack or have incompatible APIs, and/or in formats that normally would not

159 be considered interoperable such as Excel spreadsheets and flat-files. We examine the

160 extent to which the features of this architecture comply with the FAIR Principles, and suggest

161 that this might be considered a "reference implementation" for the FAIR Principles as applied

162 to non-interoperable data formats in any general or special purpose repository.

163

# Methods

164

## Implementation

165

### Overview of technical decisions and their justification

166

167

168 The World Wide Web Consortium's (W3C) Resource Description Framework (RDF) offers

169 the ability to describe entities, their attributes, and their relationships with explicit semantics

170 in a standardized manner compatible with widely used Web application formats such as

171 JSON and XML. The Linked Data Principles (Berners-Lee, 2006) mandate that data items

172 and schema elements are identified by HTTP-resolvable URIs, so the HTTP protocol can be

173 used to obtain the data. Within an RDF description, using shared public ontology terms for

174 metadata annotations supports search and large scale integration. Given all of these

175 features, we opted to use RDF as the basis of this interoperability infrastructure, as it was

176 designed to share data on the Web.

177

178 Beyond this, there was a general feeling that any implementation that required a novel data

179 discovery/sharing "Platform", "Bus", or API, was beyond the minimal design that we had

180 committed to; it would require the invention of a technology that all participants in the data

181 ecosystem would then be required to implement, and this was considered a non-starter.

182 However, there needed to be some form of coalescence around the mechanism for finding

183 and retrieving data. Our initial target-community - that is, the biomedical sciences - have

184 embraced lightweight HTTP interfaces. We propose to continue this direction with an

185 implementation based on REST (Fielding & Taylor, 2002), as several of the FAIR principles

186 map convincingly onto the objectives of the REST architectural style for distributed

187 hypermedia systems, such as having resolvable identifiers for all entities, and a common

188 machine-accessible approach to discovering and retrieving different representations of those

189 entities. The implementation we describe here is largely based on the HTTP GET method,

190 and utilizes rich metadata and hypermedia controls expressed as triples. We use widely-

191 accepted vocabularies not only to describe the data in an interoperable way, but also to

192  describe its nature (e.g. the context of the experiment and how the data was processed) and
193  how to access it. These choices help maximize uptake by our initial target-community,
194  maximize interoperability between resources, and simplify construction of the wide (not pre-
195  defined) range of client behaviors we intend to support.
196
197  Confidential and privacy-sensitive data was also an important consideration, and it was
198  recognized early on that it must be possible, within our implementation, to identify and richly
199  describe data and/or datasets without necessarily allowing direct access to them, or by
200  allowing access through existing regulatory frameworks or security infrastructures. For
201  example, many resources within the International Rare Disease Research Consortium
202  participate in the RD Connect platform (Thompson et al., 2014) which has defined the
203  "disease card" - a metadata object that gives overall information about the individual disease
204  registries, as well as a "disease matrix". The disease matrix provides aggregate data about
205  what disease variants are in the registry, how many individuals represent each disease, and
206  other high-level descriptive data that allows, for example, researchers to determine if they
207  should approach the registry to request full data access.
208
209  Finally, it was important that the data host/provider is not *necessarily* a participant in making
210  their data interoperable - rather, the interoperability solution should be capable of adapting
211  existing data with or without the source provider's participation. This ensures that the
212  interoperability objectives can be pursued for projects with limited resourcing, but more
213  importantly, that those with the needs and the resources, should adopt the responsibility for
214  making their data-of-interest interoperable, even if it is not owned by them. This distributes
215  the problem of migrating data to interoperable formats over the maximum number of
216  stakeholders, and ensures that the most crucial resources - those with the most demand for
217  interoperability - become the earliest targets for migration.
218
219  With these considerations in mind, we were inspired by three existing technologies whose
220  features were used in a novel combination to create an interoperability infrastructure for both
221  data and metadata, that is intended to also addresses the full range of FAIR requirements.
222  Briefly, the selected technologies are:
223
224      1)  The W3C's Linked Data Platform (Speicher, Arwe & Malhotra, 2015). We generated
225          a model for hierarchical dataset containers that is inspired by the concept of a LDP
226          Container, and the LDP's use of the Data Catalogue Vocabulary (DCAT, (Maali,
227          Erickson & Archer, 2014)) for describing datasets, data elements, and distributions of
228          those data elements. We also adopt the DCAT's use of Simple Knowledge
229          Organization System (SKOS, (Miles & Bechhofer, 18 August, 2009)) Concept
230          Schemes as a way to ontologically describe the content of a dataset or data record.
231      2)  The RDF Modelling Language (RML, (Dimou et al.).  RML allows us to describe one
232          or more possible RDF representations for any given dataset, and do so in a manner
233          that is, itself, FAIR: every sub-component of an RML model is Findable, Accessible,
234          Interoperable, and Reusable. Moreover, for many common semi-structured data,
235          there are generic tools that utilize RML models to dynamically drive the

236          transformation of data from these opaque representations into interoperable
237          representations (https://github.com/RMLio/RML-Mapper).
238    3) Triple Pattern Fragments (TPF - (Verborgh et al., 2016)). A TPF interface is a REST
239          Web API to retrieve RDF data from data sources in any native format. A TPF server
240          accepts URLs that represent triple patterns, and returns RDF triples from its data
241          source that match those patterns. These patterns can be used to obtain entire
242          datasets, slices through datasets, or individual data points even down to a single
243          triple (essentially a single cell in a spreadsheet table). Instead of relying on a
244          standardized contract between servers and clients, a TPF interface is self-describing
245          such that automated clients can discover the interface and its data.
246

247 We will now describe in detail how we have applied key features of these technologies, in
248 combination, to provide a novel data discoverability architecture. We will later demonstrate
249 that this combination of technologies also enables both metadata and data-level
250 interoperability even between opaque objects such as flat-files, allowing the data within
251 these objects to be queried in parallel with other data on the Semantic Web.
252

253 **Metadata Interoperability - The "FAIR Accessor" and the Linked Data Platform**
254

255 The Linked Data Platform "*defines a set of rules for HTTP operations on Web resources… to*
256 *provide an architecture for read-write Linked Data on the Web".* All entities and concepts are
257 identified by URLs, with machine-readable metadata describing the function or purpose of
258 each URL and the nature of the resource that will be returned when that URL is resolved.
259

260 Within the LDP specification is the concept of an LDP Container. A basic implementation of
261 LDP containers involves two "kinds" of resources. The first type of resource represents the
262 container - a metadata document that describes the shared features of a collection of
263 resources, and (optionally) the membership of that collection. This is analogous to, for
264 example, a metadata document describing a data repository, where the repository itself has
265 features (ownership, curation policy, etc.) that are independent from the individual data
266 records within that repository (i.e. the members of the collection). The second type of
267 resource describes a member of the contained collection and (optionally) provide ways to
268 access the record itself.
269

270 Our implementation utilizes this container concept described by the LDP, however, it does
271 not require a full implementation of LDP, as we only need read functionality, while LDP
272 defines a read/write interface. In addition, other requirements of LDP would have added
273 complexity without notable benefit. Our implementation, which we refer to as the "FAIR
274 Accessor", has two resource types, with the following features:
275

276 **Container resource:** This is a composite research object (of any kind - repository,
277 repository-record, database, dataset, data-slice, workflow, etc.). Its representation could
278 include scope or knowledge-domain covered, authorship/ownership of the object, latest
279 update, version number, curation policy, and so forth. This metadata may or may not include

280   URLs representing MetaRecord resources (described below) that comprise the individual
281   elements within the composite object. Notably, the Container URL provides a resolvable
282   identifier independent from the identifier of the dataset being described; in fact, the dataset
283   may not have an identifier, as would be the case, for example, where the container
284   represents a dynamically-generated data-slice. In addition, Containers may be published by
285   anyone - that is, the publisher of a Container may be independent from the publisher of the
286   research object it is describing. This enables one of the objectives of our interoperability
287   layer implementation - that anyone can publish metadata about any research object, thus
288   making those objects more FAIR.
289
290   **MetaRecord resource**: This is a specific element within a collection (data point, record,
291   study, service, etc.). Its representation should include information regarding licensing and
292   accessibility, access protocols, rich citation information, and other descriptive metadata. It
293   also includes a reference to the container(s) of which it is a member (the Container URL).
294   Finally, the MetaRecord may include further URLs that provide direct access to the data
295   itself, with an explicit reference to the associated data format by its MIME type (e.g.
296   text/html, application/json, application/vnd.ms-excel, text/csv, etc.).  As with Container
297   resources, MetaRecords may be published by anyone, and independently of the original
298   data publisher.
299
300   In summary, the FAIR Accessor shares commonalities with the Linked Data Platform, but
301   additionally recommends the inclusion of rich contextual metadata, based on the FAIR
302   Principles, that facilitate discovery and interoperability of repository and record-level
303   information. The FAIR Accessor is read-only, utilizing only HTTP GET together with widely-
304   used semantic frameworks to guide both human and machine exploration. Importantly, the
305   lack of a novel API means that the information is accessible to generic Web-crawling agents,
306   and may also be processed if that agent "understands" the vocabularies used.
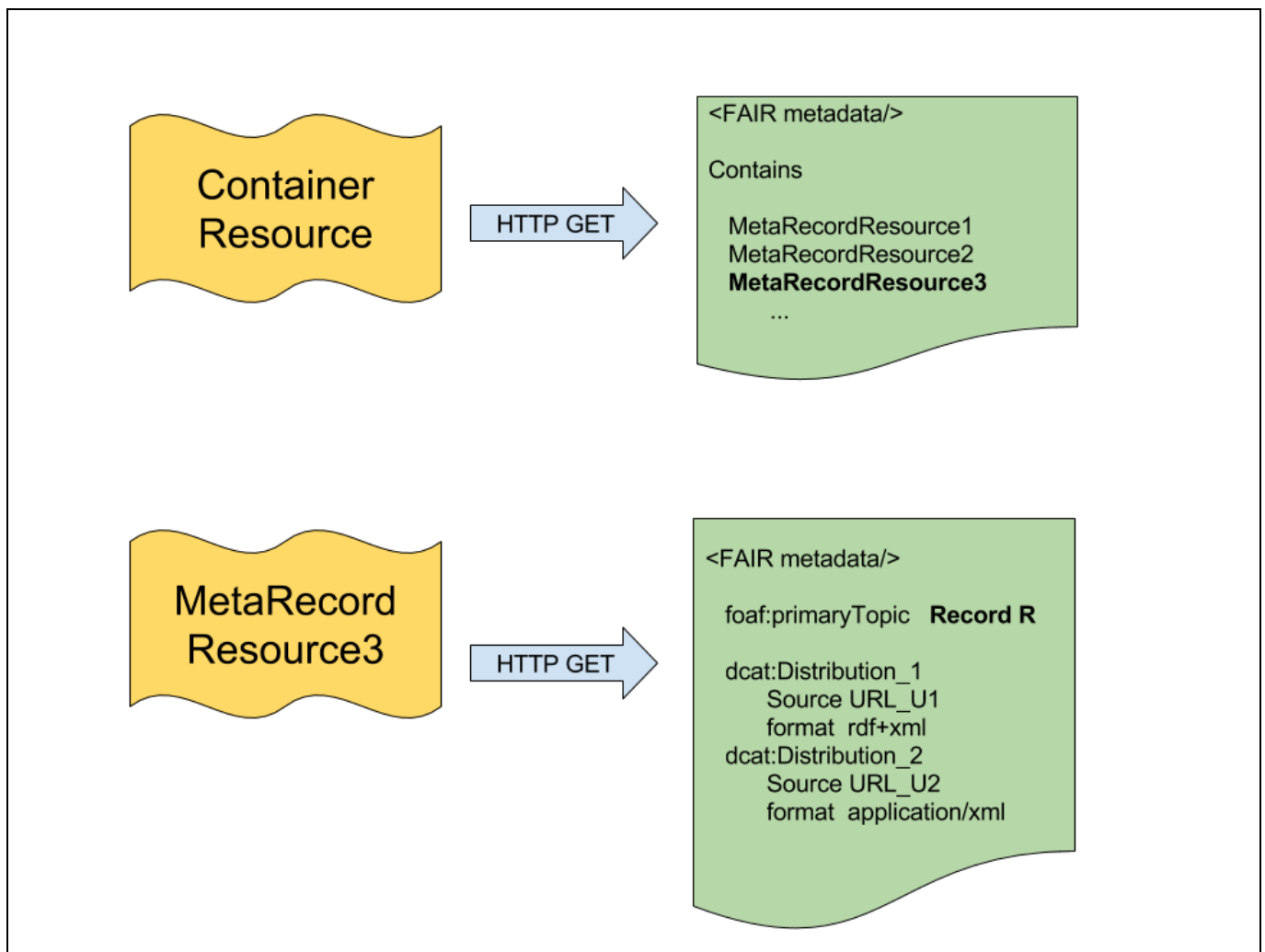307

Figure 1 The two layers of the FAIR Accessor. Inspired by the LDP Container, there are two resources in the FAIR Accessor. The first resource is a Container, which provides metadata, following FAIR Principles, about a composite research object, and optionally a list of URLs representing MetaRecords that describe individual components within the collection. The MetaRecord resources resolve to documents containing metadata about an individual data component and, optionally, a set of links structured as DCAT Distributions that lead to various representations of that data.

308
309
310    At the metadata level, therefore, this portion of the interoperability architecture provides a
311    high degree of FAIRness; however, it does not significantly enhance the FAIRness and
312    interoperability of the data itself, which was a key goal for this project. We will now describe
313    the application of two recently-published Web technologies - Triple Pattern Fragments and
314    RML - to the problem of data-level interoperability. We will show that these two technologies
315    are capable of transforming non-FAIR data into FAIR data, and will demonstrate how they
316    can be integrated into the FAIR Accessor to provide a machine-traversable path for
317    incremental drill-down from high-level repository metadata all the way through to individual
318    data points within a record.
319

320  **Data Interoperability: Compatible data discovery through RML-based FAIR**
321  **Profiles**
322
323  In our approach to data-level interoperability, we first identified a number of desiderata that
324  the solution should exhibit:
325
326  ● View-harmonization over dissimilar datatypes, allowing discovery of *potentially*
327      *integrable* data within non-integrable formats.
328  ● Support for a multitude of source data formats (XML, Excel, CSV, JSON, etc.)
329  ● "Cell-level" discovery and interoperability (referring to a "cell" in a spreadsheet)
330  ● Modular, such that a user can make interoperable only the data component they
331      require
332  ● Reusable, avoiding "one-solution-per-record" and minimizing effort/waste
333  ● Must use standard technologies, and reuse existing vocabularies
334  ● Should not require the participation of the data host (for public data)
335
336  The approach we selected was based on the premise that data, in any format, could be
337  metamodeled as a first step towards interoperability; i.e., the salient data-types and
338  relationships within an opaque data "blob" could be described in a machine-readable
339  manner. The metamodels of two data sources could then be compared to determine if their
340  contained data was, in principle, integrable.
341
342  We referred to these metamodels as "FAIR Profiles", and we further noted that there could
343  be multiple FAIR Profiles for any given data, where the Profiles might differ in structure, or
344  ontological/semantic framework. For example, a data record containing blood pressure
345  information might describe this data facet using the SNOMED vocabulary in one Profile, and
346  the ICD10 vocabulary in another Profile. We acknowledge that these meta-modelling
347  concepts are not novel, and have been suggested by a variety of other projects such as
348  DCAT (called a "DCAT Profile", though never implemented) and Dublin Core (the DC
349  Application Profile (http://dublincore.org/documents/profile-guidelines/), and have been
350  extensively described by the ISO 11179 standard ("metadata registries": http://metadata-
351  standards.org/11179/).
352
353  Our investigation into relevant existing technologies and implementations revealed a
354  relatively new, unofficial specification for a generic mapping language called "RDF Mapping
355  Language" (RML (Dimou et al.)). RML is an extension of R2RML (Das, Sundara & Cyganiak,
356  27 September, 2012), a W3C Recommendation for mapping relational databases to RDF,
357  and is described as "*a uniform mapping formalization for data in different format, which*
358  *[enables] reuse and exchange between tools and applied data*" (Dimou et al.). An RML map
359  describes the triple structure of an RDF representation, the semantic types, and the
360  constituent URI structures, that would result from a transformation of non-RDF data into RDF
361  data. RML maps are modular RDF documents where each component is a template,
362  identified by a URI, that describes the schema for a single-resource-centric graph (i.e. a
363  graph with all triples that share the same subject). The "object" position in each of these

364 triple templates may be mapped to a literal, or may be mapped as the value defined by
365 another RML module. These modules therefore assemble into a complete map of an RDF
366 representation of a data source. Finally, RML maps can also be used as templates to guide
367 the data transformation itself, using file-format-specific (but content-agnostic) software such
368 as RML Mapper (http://github.com/RMLio/RMLMapper). RML therefore fulfils each of the
369 desiderata for FAIR Profiles, and we have selected this technology as the candidate for their
370 implementation.
371
372 FAIR Profiles enable view harmonization and facilitate search/discovery of compatible but
373 structurally non-integrable data, possibly in distinct repositories. The Profiles of one data
374 resource can be compared to the Profiles of another data resource to identify commonalities
375 at the semantic level (even if the underlying data is semantically opaque) - a key step toward
376 Interoperability. FAIR Profiles created *ab initio* to fully describe a data resource, therefore,
377 have utility independent of any *actuated* transformation of the underlying data. We believe,
378 however, that it is unlikely that repository owners, or third parties, will undertake the effort of
379 creating FAIR Profiles for this purpose. We believe there is an alternative, needs-directed,
380 community-oriented approach to creating FAIR Profiles that distributes the burden of
381 designing these profiles over a broader number of researchers, and in particular, transfers
382 most of the the burden onto those who need the resources (many) rather than those who
383 own the resources (few).
384
385 Data transformation is a near-daily task for bioinformaticians, however once complete, this
386 effort is largely wasted. It would be more efficient, economical, and collaborative, to capture
387 and reuse the expert knowledge behind those transformations in a FAIR manner. Such
388 knowledge capture must not require any coordinated effort - individual researchers transform
389 data in different ways, at different times, depending on their needs - and preferably should
390 integrate into the researcher's existing work-habits. At present, these transformations - often
391 accomplished by small one-off scripts - are not published, cannot be discovered, and cannot
392 be described.
393
394 We propose that we could use the concept of a FAIR Profile to capture the cognitive effort
395 that is invested in these numerous small data transformations, where the Profile explicitly
396 expresses each individual researcher's perspective of the implicit meaning of the data. To
397 accomplish this, we propose that individuals who create transformation scripts, not only
398 publish those scripts in a publicly-accessible location, but additionally publish small RML
399 models describing the output of that transformation. (additional incentives for doing so will be
400 described in the discussion section). Specifically, we propose the following approach:
401
402   ● Data transformers independently publish one or more RML maps, where each map is
403     constrained to describe a single triple pattern that their transformation generates from
404     the underlying data source. We call these single-triple RML maps "Triple
405     Descriptors", and the structure of a Triple Descriptor is shown in Figure 2.
406   ● We further propose (but do not demonstrate here) that these Triple Descriptors may
407     later be aggregated to generate a FAIR Profile containing all triple patterns

408     associated with a given data source, from all providers. This would give us the view-
409     harmonization we desire for search and discovery, without the centralized effort. It
410     would, in fact, provide a more comprehensive view-harmonization because it would
411     also likely be redundant, containing different interpretations/representations of the
412     same data element based on the perspectives of different researchers.
413
414     RML is fully tolerant to both redundancy, and distribution of its model subcomponents. Any
415     given data point may be mapped to any number of RML models, and RML models utilize
416     URIs to identify every model component, allowing individual components to be located
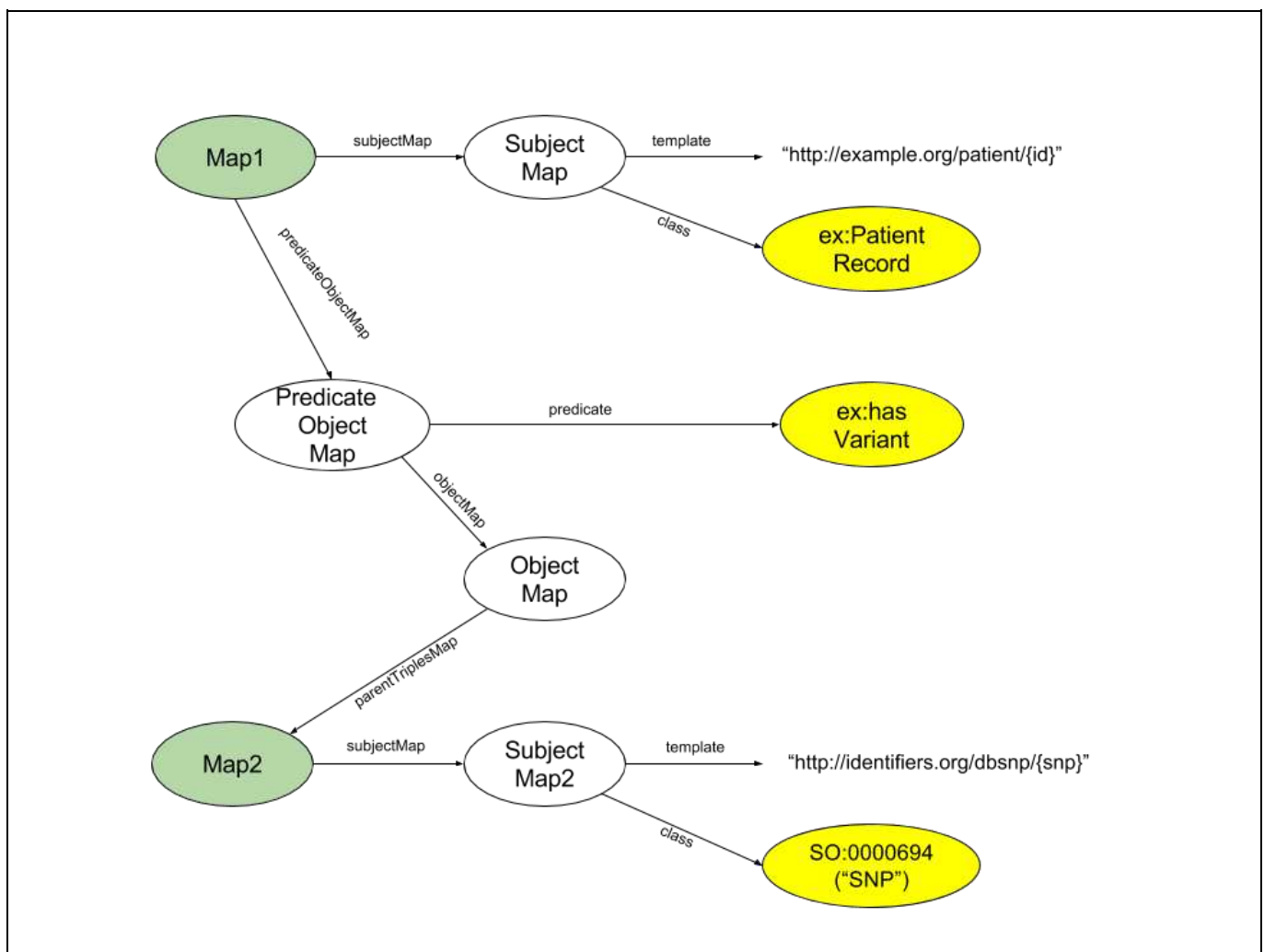417     anywhere on the Web.
418
419



Figure 2: Diagram of the structure of an exemplar Triple Descriptor representing a hypothetical record of a SNP in a patient's genome. In this descriptor, the Subject will have the URL structure *http://example.org/patient/{id}*, and the subject is of type PatientRecord. The predicate is hasVariant, and the object will have URL structure http://identifiers.org/dbsnp/{snp} with the rdf:type from the sequence ontology "0000694" (which is the concept of a "SNP"). The two nodes shaded green are of

the same ontological type, showing the iterative nature of RML, and how individual RML Triple
Descriptors will be concatenated into full FAIR Profiles. The three nodes shaded yellow are the nodes
that define the subject type, predicate and object type of the triple being described.

420
421
422   FAIR Profiles, therefore, are RML models - authored *ab initio*, and/or aggregated from Triple
423   Descriptors - that describe one or more whole or partial RDF representations for a given
424   data source. Triple Descriptors, and sometimes entire FAIR Profiles, may be re-used to
425   describe other sources if each source shares the aspects described within the model. In this
426   way, it is possible to identify, with considerable precision (i.e. potentially at the level of a
427   spreadsheet column or individual cell) potentially integrable data from two distinct sources,
428   based on the two sources sharing one or more Triple Descriptors in their FAIR Profile.
429
430

431   **Data Interoperability: Data transformation with FAIR Projectors and Triple**
432   **Pattern Fragments**
433
434   The ability to identify potentially integrable data within opaque file formats is, itself, a notable
435   achievement compared to the *status quo*. Nevertheless, beyond just discovery of relevant
436   data, our interoperability layer aims to support and facilitate cross-resource data integration
437   and query answering. This requires that the data is not only semantically described, but is
438   also semantically and syntactically transformed into a common structure.
439
440   Above, we presented a mechanism to describe structure and semantics - Triple Descriptors
441   in RML - what remains lacking, however, is a way to execute data transformations that
442   provide output consistent with a given Triple Descriptor.  Although in the previous section we
443   proposed that those who undertake data transformations should publish their transformation
444   script, together with its associated Triple Descriptors, this does not address a critical barrier
445   to interoperability - opaque, non-machine-readable interfaces and API proliferation
446   (Verborgh & Dumontier, 2016).  We propose, therefore, that what is required is a universally-
447   applicable way of retrieving data from any transformation script (or any data source), without
448   inventing a new API. We now describe our suggestion for how to achieve this behavior, and
449   we refer to such transformation tools as "FAIR Projectors".
450
451   Triple Pattern Fragments (TPF) defines a REST interface through which clients can request
452   triples based on a triple pattern [S,P,O] where any component of that pattern is either a
453   constant or a variable. In response, a TPF server returns pages with all triples from its data
454   source that match the incoming pattern.   We use the TPF interface for FAIR Projectors, and
455   therefore all Projectors share a common URL pattern, defined by the Triple Pattern
456   Fragments specification (Verborgh et al., 2016).  In addition, we require that the semantics of
457   the output triple patterns are defined by (one or more) Triple Descriptors, thus allowing a
458   client to select the appropriate FAIR Projector for its needs.
459

460   A FAIR Projector, therefore, is a Web resource that is associated with *both* a particular data
461   source, and particular Triple Descriptor(s). Calling HTTP GET on the URL of the FAIR
462   Projector produces RDF triples from the data source that match the format defined by that
463   Projector's Triple Descriptor.  The originating data source behind a Projector may be a
464   database, a data transformation script, an analytical web service, another FAIR Projector, or
465   any other data-source.
466

467   **Linking the Components: FAIR Projectors and the FAIR Accessor**
468
469   At this point, we have a means for obtaining triples with a specific structure - TPF Servers -
470   and we have a means of describing the structure and semantics of those triples - Triple
471   Descriptors. Together these two elements define a FAIR Projector. However, we still lack a
472   formal mechanism for linking these two components, such that the discovery of a Triple
473   Descriptor with the desired semantics, also provides its associated TPF Server (Projector)
474   URL.
475
476   We propose that this association can be easily accomplished, without defining any novel API
477   or standard, if the output of a FAIR Projector is considered a type of DCAT Distribution. In
478   this way, it may be included as another distribution component of the MetaRecord metadata
479   from a FAIR Accessor, where the URL of the Projector, and its Triple Descriptor, are
480   metadata elements of that Distribution. This is diagrammed in Figure 3, where Distribution_3
481   and Distribution_4 include Triple Pattern Fragment URLs served by a FAIR Projector, and
482   also include the Triple Descriptor RML model that describes the structure and semantics of
483   the data that will be produced by calling that URL.
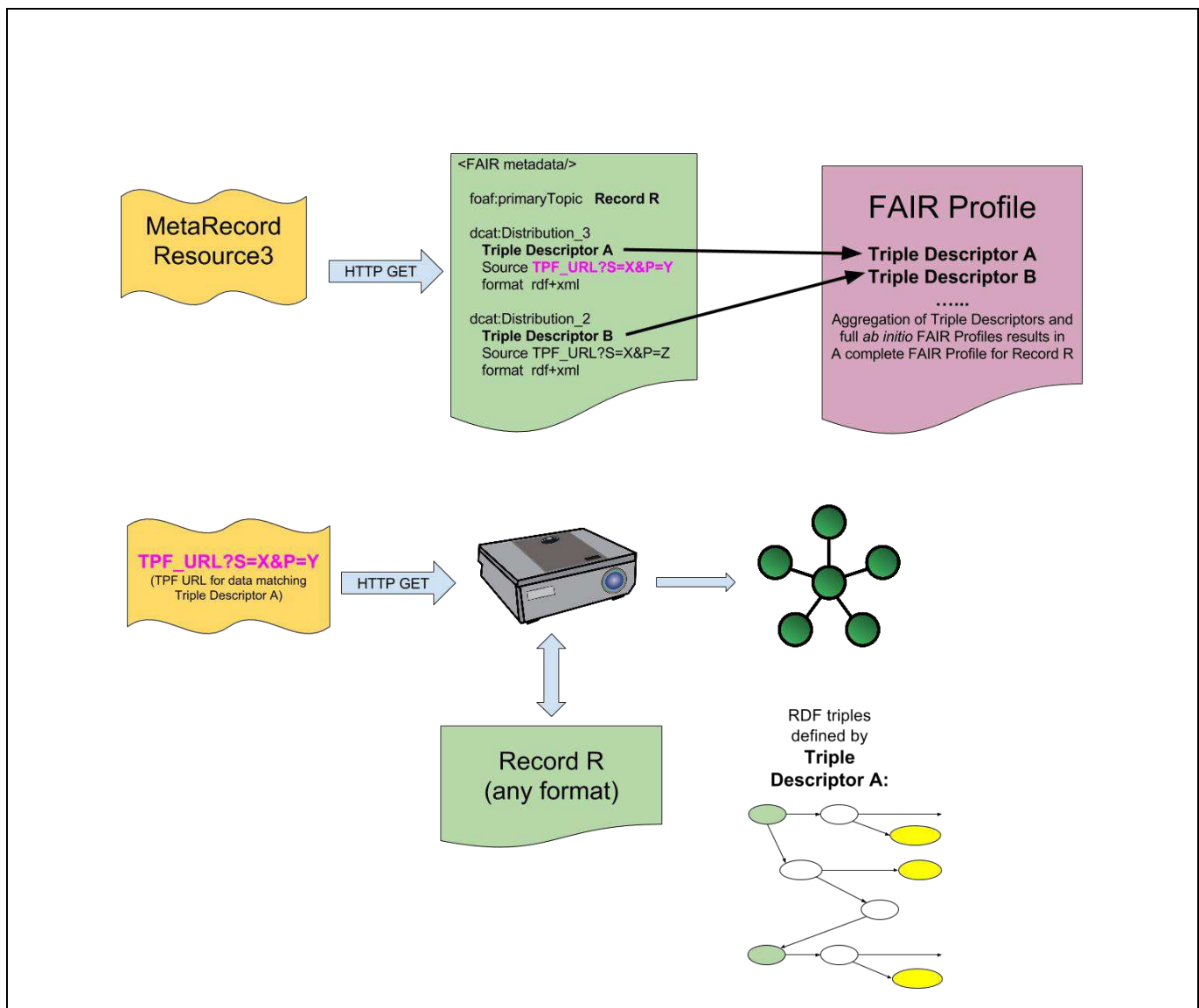484
485
486
487

Figure 3. Integration of FAIR Projectors into the FAIR Accessor. Resolving the MetaRecord resource returns a metadata document containing multiple DCAT Distributions for a given record, as in Figure 1. When a FAIR Projector is available, additional DCAT Distributions are included in this metadata document. These Distributions contain a URL (purple text) representing a Projector, and a Triple Descriptor that describes, in RML, the structure and semantics of the Triple(s) that will be obtained from that Projector resource if it is resolved. These Triple Descriptors may be aggregated into FAIR Profiles, based on the Record that they are associated with (Record R, in the figure) to give a full mapping of all available representations of the data present in Record R.

488
489

## Results

491

492    To demonstrate the interoperability layer, we will explore an example involving UniProt. In
493    this example, we create a FAIR Accessor for a dataset that consists of a specific "slice" of

494 the Protein records within the UniProt database - that is, the set of proteins in *Aspergillus*
495 *nidulans* (taxon 16245) that are annotated as being involved in RNA Metabolism (GO
496 0006396). We first demonstrate the functionality of the two layers of the FAIR Accessor. We
497 then demonstrate a FAIR Projector, and show how its metadata integrates into the FAIR
498 Accessor. In this example, the Projector modifies the ontological framework of the UniProt
499 data such that the ontological terms used by UniProt are replaced by the terms specified in
500 EDAM. We will demonstrate that this transformation is specified, in a machine-readable way,
501 by the FAIR Triple Descriptor that accompanies each Projector's metadata.
502

503 **The two-step FAIR Accessor**
504
505 The example FAIR Accessor serves the results of the following query against the UniProt
506 SPARQL endpoint:
507
```
508        PREFIX up:<http://purl.uniprot.org/core/>
509        PREFIX taxon:<http://purl.uniprot.org/taxonomy/>
510        PREFIX rdf:<http://www.w3.org/1999/02/22-rdf-syntax-ns#>
511        PREFIX rdfs:<http://www.w3.org/2000/01/rdf-schema#>
512        SELECT distinct ?id
513
514        WHERE
515        {
516               ?protein a up:Protein .
517               ?protein up:organism ?organism .
518               ?organism rdfs:subClassOf taxon:162425 .
519               ?protein up:classifiedWith ?go .
520               ?go rdfs:subClassOf* <http://purl.obolibrary.org/obo/GO_0006396> .
521
522               bind(replace(str(?protein), "http://purl.uniprot.org/uniprot/", "",
523        "i") as ?id)
524        }
```
525
526 Accessor output is retrieved from the Container Resource URL:
527
```
528        http://linkeddata.systems/Accessors/UniProtAccessor
```
529
530 The result of calling GET on the Container Resource URL is visualized in Figure 4, where
531 Tabulator (Tim Berners-lee et al., 2006) is used to render the output as HTML for human-
532 readability.
533
534

| UniProt Slice FAIR Accessor - Aspergillus RNA Processing proteins | creator | wilkinsonlab.info/ |
| | language | eng |
| | license | cc by nd4.0 |
| | title | UniProt Slice FAIR Accessor - Aspergillus RNA Processing proteins |
| | authored By | 0000 0002 9699 485X |
| | entities | 412 |
| | term has Principal Investigator | Dr. Mark Wilkinson |
| | type | Dataset |
| | | Basic Container |
| | | Collection |
| | contact Point | Wilkinson.rdf |
| | description | Takes a SPARQL query of the UniProt database specific to proteins and their GO annotations related to RNA Procssing proteins in Aspergillus and makes it a FAIR Accessor source. The precise query is: |

```
PREFIX up:<http://purl.uniprot.org/core/>
PREFIX taxon:<http://purl.uniprot.org/taxonomy/>
PREFIX rdf:<http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX rdfs:<http://www.w3.org/2000/01/rdf-schema#>
SELECT distinct ?id

WHERE
{
    ?protein a up:Protein .
    ?protein up:organism ?organism .
    ?organism rdfs:subClassOf taxon:162425 .
    ?protein up:classifiedWith ?go .
    ?go rdfs:subClassOf* <http://purl.obolibrary.org/obo/GO_0006396> .

    bind(replace(str(?protein), "http://purl.uniprot.org/uniprot/", "", "i") as ?id)
}
```

| | identifier | Uni Prot Accessor |
| | keyword | Aspergillus nidulans |
| | | Aspergillus |
| | | Proteins |
| | | RNA Processing |
| | landing Page | uniprot.org/ |
| | language | en |
| | publisher | wilkinsonlab.info/ |
| | theme | RNA Processing conceptscheme.rdf |
| | contains | C8UZX9 |
| | | C8UZY5 |
| | | C8V0B4 |
| | | C8V0M2 |
| | | C8V0U7 |

Figure 4. A representative portion of the output from resolving the Container Resource of the FAIR Accessor, rendered into HTML by the Tabulator Firefox plugin. The three columns show the label of the Subject node of all RDF Triples (left), the label of the URI in the predicate position of each Triple (middle), and the value of the Object position (right), where blue text indicates that the value is a Resource, and black text indicates that the value is a literal.

535
536
537    Of particular note are the following metadata elements:
538
539

| http://purl.org/dc/elements/1.1/license | https://creativecommons.org/licenses/by-nd/4.0/ |

| http://purl.org/pav/authoredBy | http://orcid.org/0000-0002-9699-485X |
|---|---|
| http://rdfs.org/ns/void#entities | 411 |
| a | http://purl.org/dc/dcmitype/Dataset<br><br>http://www.w3.org/ns/ldp#BasicContainer<br><br>http://www.w3.org/ns/prov#Collection |
| http://www.w3.org/ns/dcat#contactPoint | http://biordf.org/DataFairPort/MiscRDF/Wilkinson.rdf |
| http://www.w3.org/ns/dcat#keyword | "Aspergillus nidulans", "Aspergillus", "Proteins", "RNA Processing"; |
| http://www.w3.org/ns/dcat#theme | http://linkeddata.systems/ConceptSchemes/RNA_Processing_conceptscheme.rdf |
| http://www.w3.org/ns/ldp#contains | http://linkeddata.systems/cgi-bin/Accessors/ UniProtAccessor/C8UZX9<br><br>http://linkeddata.systems/cgi-bin/Accessors/ UniProtAccessor/C8UZY5<br><br>… |

540
541  ● License information is provided as an HTML + RDFa document, following one of the
542     primary standard license forms published by Creative Commons. This allows the
543     license to be unambiguously interpreted by both machines and people prior to
544     accessing any data elements, an important feature that will be discussed later.
545  ● Authorship is provided by name, using the Academic Research Project Funding
546     Ontology (ARPFO), but is also unambiguously provided by a link to the author's
547     ORCID, using the Provenance Authoring and Versioning (PAV) ontology.
548  ● The repository descriptor is typed as being a Dublin Core Dataset, a Linked Data
549     Platform container, and a Provenance Collection, allowing it to be interpreted by a
550     variety of client agents, and conforming to several best-practices, such as the
551     Healthcare and Life Science Dataset Description guidelines (Dumontier et al., 2016))
552  ● Contact information is provided in a machine-readable manner via the FoaF record of
553     the author, and the DCAT ontology "contactPoint" property.
554  ● Human readable keywords, using DCAT, are mirrored and/or enhanced by a
555     machine-readable RDF document which is the value of the DCAT "theme" property.
556     This RDF document follows the structure determined by the SKOS ontology, and lists
557     the ontological terms that describe the repository for machine-processing.
558  ● Finally, individual records within the dataset are represented as the value of the
559     Linked Data Platform "contains" property, and provided as a possibly paginated list of
560     URLs (a discussion of machine-actionable pagination will not be included here).
561     These URLs are the MetaRecord Resource URLs shown in Figure 1.

562
563
564    Following the flow in Figure 1, the next step in the FAIR Accessor is to resolve a
565    MetaRecord Resource URL. For clarity, we will first show the metadata document that is
566    returned if there are no FAIR Projectors for that dataset. In the subsequent section, we will
567    show how FAIR Projectors enhance this basic metadata with additional features.
568
569    Calling HTTP GET on a MetaRecord Resource URL returns a document with the structure
570    shown in Figure 5.
571
572

| UniProt Protein C8UZX9 | bibliographic Citation | The UniProt Consortium (2015). UniProt: a hub for protein information. Nucleic Acids Res. 43: D204-D212 |
| | creator | UniProt Consortium |
| | language | eng |
| | license | cc by nd3.0 |
| | title | UniProt Protein C8UZX9 |
| | in Dataset | Uni Prot Accessor/ |
| | contact Point | contact |
| | description | KRR1 small subunit processome componentKRR-R motif-containing protein 1 |
| | distribution | C8UZX9.rdf |
| | | C8UZX9.html |
| | identifier | C8UZX9 |
| | keyword | Annotation |
| | | Aspergillus nidulans |
| | | Aspergillus |
| | | Functinal Annotation |
| | | GO |
| | | Gene Ontology |
| | | Proteins |
| | | RNA Processing |
| | landing Page | uniprot.org |
| | language | en |
| | publisher | uniprot.org |
| | page | sparql |
| | | uniprot.org/ |
| | primary topic | C8UZX9 |
| C8UZX9 | ... | |
| C8UZX9.rdf | format | application/rdf+xml |
| | type | Dataset |
| | | Dataset |
| | | Distribution |
| | download URL | C8UZX9.rdf |
| C8UZX9.html | format | text/html |
| | type | Dataset |
| | | Distribution |
| | download URL | C8UZX9.html |

Figure 5. A representative portion of the output from resolving the MetaRecord Resource of the FAIR Accessor for record C8UZX9, rendered into HTML by the Tabulator Firefox plugin. The columns have the same meaning as in Figure 4.

573
574
575    Many properties in this metadata document are similar to those at the higher level of the
576    FAIR Accessor, however, the primary topic of this document is the original UniProt record.

577 Therefore, the values of these facets now reflect the authorship and contact information for
578 the record itself. We do, however, recognize that MetaRecords are themselves scholarly
579 works and should be properly cited. The MetaRecord includes the "in dataset" predicate,
580 which referrs back to the first level of the FAIR Accessor, thus this provides one avenue for
581 capturing the provenance information for the MetaRecord. If additional provenance detail is
582 required, we propose (but no not describe furrther here) that this information could be
583 contained in a separate named graph, in a manner akin to that used by
584 NanoPublications(Kuhn et al., 2016).
585
586 The important distinctive property in this document is the "distribution" property, from the
587 DCAT ontology. For clarity, an abbreviated document in Turtle format is shown in Figure 6,
588 containing only the "distribution" elements and their values.
589
590

```
@prefix dcat: <http://www.w3.org/ns/dcat#>.
@prefix Uni: <./>.
@prefix n0: <http://purl.org/dc/elements/1.1/>.
@prefix void: <http://rdfs.org/ns/void#>.

Uni:C8UZX9
    dcat:distribution
        <http://www.uniprot.org/uniprot/C8UZX9.rdf>,
        <http://www.uniprot.org/uniprot/C8UZX9.html> .

<http://www.uniprot.org/uniprot/C8UZX9.rdf>
    n0:format
        "application/rdf+xml";
    a     n0:Dataset, void:Dataset, dcat:Distribution;
    dcat:downloadURL
        <http://www.uniprot.org/uniprot/C8UZX9.rdf>.

<http://www.uniprot.org/uniprot/C8UZX9.html>
    n0:format
        "text/html";
    a     n0:Dataset, dcat:Distribution;
    dcat:downloadURL
        <http://www.uniprot.org/uniprot/C8UZX9.html>.
```

Figure 6. Turtle representation of the subset of triples from the MetaRecord metadata pertaining to the two DCAT Distributions. Each distribution specifies an available representation (media type), and a URL from which that representation can be downloaded.

591
592
593 There are two DCAT Distributions in this document. The first is described as being in format
594 "application/rdf+xml", with its associated download URL. The second is described as being
595 in format "text/html", again with the correct URL for that representation. Both are typed as

596  Distributions from the DCAT ontology. These distributions are published by UniProt
597  themselves, and the UniProt URLs are used. The additional metadata in the FAIR Accessor
598  explicitly describes the keywords that relate to that record (both machine and human-
599  readable), access policy, license, and format, allowing machines to more accurately
600  determine the utility of this record prior to retrieving it.
601
602  Several things are important to note before moving to a discussion of FAIR Projectors. First,
603  the two levels of the FAIR Accessor are not interdependent. The Container layer can
604  describe relevant information about the scope and nature of a repository, but might not
605  provide any further links to MetaRecords. Similarly, whether or not to provide a distribution
606  within a MetaRecord is entirely at the discretion of the data owner. For sensitive data, an
607  owner may chose to simply provide (even limited) metadata, but not provide any direct link to
608  the data itself, and this is perfectly conformant with the FAIR guidelines. Further, when
609  publishing a single data record, it is not obligatory to publish the Container level of the FAIR
610  Accessor; one could simply provide the MetaRecord document describing that data file,
611  together with an optional link to that file as a Distribution.
612

## The FAIR Projector

614
615  FAIR Projectors can be used for many purposes, including (but not limited to) transformation
616  of a data source from non-Linked Data to Linked Data, transformation of a Linked Data
617  source into a different Linked Data structure or ontological framework, load-
618  management/query-management, or as a means to explicitly describe the ontological
619  structure of an underlying data source in a searchable manner. In this demonstration, the
620  FAIR Projector transforms the semantics of the native RDF provided by UniProt into a
621  different ontological framework (EDAM).
622
623  The address of this FAIR Projector's TPF interface is:
624
625  `http://linkeddata.systems:3001/fragments`
626
627  The TPF API requires  a subject and/or predicate and/or object node to be specified as
628  parameters; a request for the all-variable pattern will (currently) return nothing. How can a
629  software agent know what parameters are valid, and what will be returned from such a call?
630
631  In this interoperability infrastructure, we propose that Projectors should be considered as
632  DCAT Distributions, and thus TPF URLs, with appropriate parameters added and bound, are
633  included in the distribution section of the MetaRecord metadata. An example is shown in
634  Figure 6, again rendered using Tabulator.
635
636

| UniProt Protein C8UZX9 | bibliographic Citation | The UniProt Consortium (2015). UniProt: a hub for protein information. Nucleic Acids Res. 43: D204-D212 |
| | creator | UniProt Consortium |
| | language | eng |
| | license | cc by nd3.0 |
| | title | UniProt Protein C8UZX9 |
| | in Dataset | Uni Prot Accessor/ |
| | contact Point | contact |
| | description | KRR1 small subunit processome componentKRR-R motif-containing protein 1 |
| | distribution | fragments?subject=http%3A%2F%2Fidentifiers%2Eorg%2Funiprot%2FC8UZX9&predicate=http%3A%2F%2Fpurl%2Euniprot%2Eorg%2Fcore%2Fclassified With |
| | | fragments?subject=http%3A%2F%2Fidentifiers%2Eorg%2Funiprot%2FC8UZX9&predicate=http%3A%2F%2Fpurl%2Euniprot%2Eorg%2Fcore%2Forganism |
| | | C8UZX9.rdf |
| | | C8UZX9.html |
| | identifier | C8UZX9 |
| | keyword | Annotation |
| | | Aspergillus nidulans |
| | | Aspergillus |
| | | Functinal Annotation |
| | | GO |
| | | Gene Ontology |
| | | Proteins |
| | | RNA Processing |
| | landing Page | uniprot.org |
| | language | en |
| | publisher | uniprot.org |
| | page | sparql |
| | | uniprot.org/ |
| | primary topic | C8UZX9 |

Figure 7. A portion of the output from resolving the MetaRecord Resource of the FAIR Accessor for record C8UZX9, rendered into HTML by the Tabulator Firefox plugin. The columns have the same meaning as in Figure 4. Comparing the structure of this document to that in Figure 5 shows that there are now four values for the "distribution" predicate. An RDF and HTML representation, as in Figure 5, and two additional distributions with URLs conforming to the TPF design pattern (highlighted).

637
638
639 Note that, in addition to the two distributions C8V1J1.html and C8V1J1.rdf that were seen in
640 Figure 5, there are now two additional distributions that include both a subject and predicate
641 parameter in their URLs. These are the URLs for two FAIR Projections of that data.
642
643 Again, looking at an abbreviated and simplified Turtle document for clarity (Figure 8) we can
644 see the metadata structure of one of these two new distributions.
645
646

```
Uni:C8UZX9
    dcat:distribution
<http://linkeddata.systems:3001/fragments?subject=http%3A%2F%2Fidentifiers%2Eorg%2Funiprot%2
FC8UZX9&predicate=http%3A%2F%2Fpurl%2Euniprot%2Eorg%2Fcore%2FclassifiedWith> .

<http://linkeddata.systems:3001/fragments?subject=http%3A%2F%2Fidentifiers%2Eorg%2Funiprot%2
FC8UZX9&predicate=http%3A%2F%2Fpurl%2Euniprot%2Eorg%2Fcore%2FclassifiedWith>
```

```
    n0:format
        "application/rdf+xml", "application/x-turtle", "text/html";
    a    FAI:Projector, n0:Dataset, void:Dataset, dcat:Distribution;
    dcat:downloadURL
<http://linkeddata.systems:3001/fragments?subject=http%3A%2F%2Fidentifiers%2Eorg%2Funiprot%2
FC8UZX9&predicate=http%3A%2F%2Fpurl%2Euniprot%2Eorg%2Fcore%2FclassifiedWith>.


loc:Source3C0D4EAA-8497-11E6-99DD-D5545D07C3DD
    rml:hasMapping
        loc:Mappings3C0D4EAA-8497-11E6-99DD-D5545D07C3DD;
    rml:referenceFormulation
        ql:TriplePatternFragments;
    rml:source
<http://linkeddata.systems:3001/fragments?subject=http%3A%2F%2Fidentifiers%2Eorg%2Funiprot%2
FC8UZX9&predicate=http%3A%2F%2Fpurl%2Euniprot%2Eorg%2Fcore%2FclassifiedWith> .


loc:Mappings3C0D4EAA-8497-11E6-99DD-D5545D07C3DD
    rml:logicalSource
        loc:Source3C0D4EAA-8497-11E6-99DD-D5545D07C3DD;
    rr:predicateObjectMap
        loc:POMap3C0D4EAA-8497-11E6-99DD-D5545D07C3DD;
    rr:subjectMap
        loc:SubjectMap3C0D4EAA-8497-11E6-99DD-D5545D07C3DD.
```

Figure 8. Turtle representation of the subset of triples from the MetaRecord metadata pertaining to one of the FAIR Projector DCAT Distributions of the MetaRecord shown in Figure 6. The text is color-coded to assist in visual exploration of the RDF.  The DCAT Distribution blocks of the two Projector distributions (black bold) have multiple media-type representations (red), and are connected to an RML logicalSource (purple), which itself is linked by the hasMapping predicate to a Mappings (blue) block of RML that semantically describes the subject, predicate, and object (green and orange) of the Triple Descriptor for that Projector.  The full RML model is shown separately in Figure 9.

647
648 Following the Triple Pattern Fragments behavior, requesting the downloadURL with HTTP
649 GET will trigger the Projector to generate all triples where the subject is UniProt record
650 C8UZX9, and the predicate is "classifiedWith" from the UniProt Core ontology.  Those triples
651 will match the semantics and structure defined in the Mappings (blue) block. The
652 interpretation of the Dublin Core "format" predicate in this context is noteworthy, as its value
653 is only loosely defined by Dublin Core. A Projector is a RESTful resource that will respond to
654 HTTP content-negotiation to select the representation of the requested resource. The values
655 of the "format" predicate in this example should be interpreted as a list of the possible
656 formats available, for example, to be used as valid values for the HTTP Accept Header when
657 calling that resource. In this case, there are three available representations - Turtle, HTML,
658 and RDF/XML.
659
660 The schematic structure of the Mapping RML is visualized in Figure 2, with the actual output
661 from the Accessor shown in Figure 9, color-coded to assist visual exploration. The RML
662 describes a Triple where the subject will be of type `edam:data_0896` ("Protein record"), the

663    predicate will be "`classifiedWith`" from the UniProt Core ontology, and the object will be
664    of type `edam:data_1176` ("GO Concept ID").
665
666

```
loc:Mappings3C0D4EAA-8497-11E6-99DD-D5545D07C3DD
    rml:logicalSource
        loc:Source3C0D4EAA-8497-11E6-99DD-D5545D07C3DD;
    rr:predicateObjectMap
        loc:POMap3C0D4EAA-8497-11E6-99DD-D5545D07C3DD;
    rr:subjectMap
        loc:SubjectMap3C0D4EAA-8497-11E6-99DD-D5545D07C3DD.


loc:SubjectMap3C0D4EAA-8497-11E6-99DD-D5545D07C3DD
    rr:class ed:data_0896; rr:template "http://identifiers.org/uniprot/{ID}".


loc:POMap3C0D4EAA-8497-11E6-99DD-D5545D07C3DD
    rr:objectMap
        loc:ObjectMap3C0D4EAA-8497-11E6-99DD-D5545D07C3DD;
    rr:predicate
        core:classifiedWith.


loc:ObjectMap3C0D4EAA-8497-11E6-99DD-D5545D07C3DD
    rr:parentTriplesMap loc:SubjectMap23C0D4EAA-8497-11E6-99DD-D5545D07C3DD.

loc:SubjectMap23C0D4EAA-8497-11E6-99DD-D5545D07C3DD
    rr:class ed:data_1176; rr:template "http://purl.obolibrary.org/obo/{GO}".
```

Figure 9. Turtle representation of a Triple Descriptor within the MetaRecord metadata shown in Figures 6 and 7. The text is color-coded to assist in visual exploration of the RDF. The RDF structure shown here is represented schematically in Figure 2. The black bold text shows the locations where the semantic type of the projected subject, predicate, and object are stored.

667
668    The triples that are returned by calling HTTP GET on that Projector URL are:
669
670        @prefix uni: <http://identifiers.org/uniprot/>.
671        @prefix obo: <http://purl.obolibrary.org/obo/>.
672        uni:C8UZX9 core:classifiedWith obo:GO_0000447, obo:GO_0000462
673        .
674
675    This is accompanied by a block of hypermedia controls (not shown) using the Hydra
676    vocabulary (Lanthaler & Gütl; Das, Sundara & Cyganiak, 27 September, 2012) that provide
677    machine-readable instructions for how to navigate the remainder of that dataset - for
678    example, how to get the entire row, or the entire column for the current data-point.
679

680  Though the subject and object are not explicitly typed in the output from this call to the
681  Projector, further exploration of the Projector's output, via those TPF's hypermedia controls,
682  would reveal that the Subject and Object are in fact typed according to the EDAM ontology
683  (Ison et al., 2013), as declared in the RML meta-descriptor.  Thus, this FAIR Projector
684  transformed data from UniProt Core semantic types, to the equivalent data, now represented
685  within the EDAM semantic framework, as shown in Figure 10.  Also note that the URI
686  structure for the UniProt entity has been changed from the UniProt URI scheme to the more
687  interoperable Identifiers.org scheme.
688
689  This example was chosen because, in UniProt and the Gene Ontology consortium's
690  representation, Gene Ontology terms do not have a richer classification than "owl:Class".
691  With respect to interoperability, this is problematic, as the lack of rich semantic typing
692  prevents them from being used for automated discovery of resources that could potentially
693  consume them, or use them for integrative, cross-domain queries.
694
695

| In UniProt | `http://purl.uniprot.org/uniprot/C8UZX9`<br>    `a`<br>        `http://purl.uniprot.org/core/Protein ;`<br><br>    `http://purl.uniprot.org/core/classifiedWith`<br>        `http://purl.obolibrary.org/obo/GO_0000462 .`<br><br>`http://purl.obolibrary.org/obo/GO_0000462`<br>    `a`<br>        `http://www.w3.org/2002/07/owl#Class` |
| --- | --- |
| After Projection | **`http://identifiers.org/uniprot/`**`C8UZX9`<br>    `a`<br>        **`http://edamontology.org/data_0896 ;`**<br><br>    `http://purl.uniprot.org/core/classifiedWith`<br>        `http://purl.obolibrary.org/obo/GO_0000462 .`<br><br>`http://purl.obolibrary.org/obo/GO_0000462`<br>    `a`<br>        **`http://edamontology.org/data_1176`** |

Figure 10:  Data before and after FAIR Projection.  Bolded segments show how the URI structure and the semantics of the data were modified, according to the mapping defined in the Triple Descriptor (data_0896 = "Protein report" and data_1176 = "GO Concept ID").  URI structure transformations may be useful for integrative queries against datasets that utilize the Identifiers.org URI scheme such as OpenLifeData (González et al., 2014). Semantic transformations allow integrative queries across datasets that utilize diverse and redundant ontologies for describing their data, and in this example, may also be used to add semantics where there were none before.

696

# Discussion

Interoperability is hard. It was immediately evident that, of the four FAIR principles, Interoperability was going to be the most challenging. Here we have designed a novel infrastructure with the primary objective of interoperability for both metadata and data, but with an eye to all four of the FAIR Principles. We wished to provide discoverable and interoperable access to a wide range of underlying data sources - even those in computationally opaque formats - as well as supporting wide array of both academic and commercial end-user applications above these data sources. In addition, we imposed constraints on our selection of technologies; in particular, that the implementation should re-use existing technologies as much as possible, and should support multiple and unpredictable end-uses. Moreover, it was accepted from the outset that the tradeoff between simplicity and power was one that could not be avoided. While other interoperability projects such as caBIO (Covitz et al., 2003) and TAPIR (De Giovanni et al., 2010) created rich APIs or query languages, enabling extremely powerful cross-resource data exploration and integration, this was done at the expense of broad-scale uptake and/or with the explicit and unavoidable participation of the individual providers. Thus, with the goal of maximizing global uptake and adoption of this interoperability infrastructure, and democratizing the cost of implementation over the entire stakeholder community - both users and providers - we opted for lightweight, weakly integrative, REST-based solutions, that nevertheless lend themselves to significant degrees of mechanization in both discovery and integration.

We now look more closely at how this interoperability infrastructure meets the expectations within the FAIR Principles.

**FAIR facet(s) addressed by the Container Resource:**
- **Findable** - The container has a distinct globally unique and resolvable identifier, allowing it to be discovered and explicitly, unambiguously cited. This is important because, in many cases, the dataset being described does not natively possess an identifier, as in our example above where the dataset represented the results of a query. In addition, the container's metadata describes the research object, allowing humans and machines to evaluate the potential utility of that object for their task.
- **Accessible** - the Container URL resolves to a metadata record using standard HTTP GET. In addition to describing the nature of the research object, the metadata record should include information regarding licensing, access restrictions, and/or the access protocol for the research object. Importantly, the container metadata exists independently of the research object it describes, where FAIR Accessibility requires metadata to be persistently available even if the data itself is not.
- **Interoperable** - The metadata is provided in RDF - a globally-applicable syntax for data and knowledge sharing. In addition, the metadata uses shared, widely-adopted public ontologies and vocabularies to facilitate interoperability at the metadata level.

741   ● **Reusable** - the metadata includes citation information related to the
742      authorship of the container and/or its contents, and license information related
743      to the reuse of the data, by whom, and for what purpose.

745   **Other features of the Container Resource**
746      ● **Privacy protection** - The container metadata provides access to a rich
747         description of the content of a resource, without exposing any data within that
748         resource. While a provider may choose to include MetaRecord URLs within
749         this container, they are not required to do so if, for example, the data is highly
750         sensitive, or no longer easily accessible; however, the contact information
751         provided within the container allows potential users of that data to inquire as
752         to the possibility of gaining access in some other way.  As such, this container
753         facilitates a high degree of FAIRness, while still providing a high degree of
754         privacy protection.

756   **FAIR Facet(s) Addressed by the MetaRecord:**
757      ● **Findable** - The MetaRecord URL is a globally-unique and resolvable identifier
758         for a data entity, regardless of whether or not it natively possesses an
759         identifier. The metadata it resolves to allows both humans and machines to
760         interrogate the nature of a data element before deciding to access it.
761      ● **Accessible** - the metadata provided by accessing the MetaRecord URL
762         describes the accessibility protocol and license information for that record,
763         and describes all available formats.
764      ● **Interoperable** - as with the Container metadata, the use of shared ontologies
765         and RDF ensures that the metadata is interoperable.
766      ● **Reusable** - the MetaRecord metadata should carry record-level citation
767         information to ensure proper attribution if the data is used.  We further
768         propose, but do not demonstrate, that authorship of the MetaRecord itself
769         could be carried in a second named-graph, in a manner similar to that
770         adopted by the NanoPublication community.

772   **Other features of the MetaRecord**
773      ● **Privacy protection** - the MetaRecord provides for rich descriptive information
774         about a specific member of a collection, where the granularity of that
775         description is entirely under the control of the data owner. As such, the
776         MetaRecord can provide a high degree of FAIRness at the level of an
777         individual record, without necessarily exposing any identifiable information. In
778         addition, the provider may choose to stop at this level of FAIRness, and not
779         include further URLs giving access to the data itself.
780      ● **Symmetry of traversal** - Since we predict that clients will, in the future, query
781         over indexes of FAIR metadata searching for dataset or records of interest, it
782         is not possible to predict the position at which a client or their agent will enter
783         your FAIR Accessor. While the container metadata provides links to individual
784         MetaRecords, the MetaRecord similarly provides a reference back "upwards"

785            to its container. Thus a client can access repository-level metadata (e.g.
786            curation policy, ownership, linking policy) for any given data element it
787            discovers. This became particularly relevant as a result of the European Court
788            of Justice decision
789            (http://curia.europa.eu/jcms/upload/docs/application/pdf/2016-
790            09/cp160092en.pdf) that puts the burden of proof on those who create
791            hyperlinks to ensure the document they link to is not, itself, in violation of
792            copyright.
793       ●   **High granularity of access control** - individual elements of a collection may
794            have distinct access constraints or licenses. For example, individual patients
795            within a study may have provided different consent. MetaRecords allow each
796            element within a collection to possess, and publish, its own access policy,
797            access protocol, license, and/or usage-constraints, thus providing fine-
798            grained control of the access/use of individual elements within a repository.
799
800
801
802        **FAIR Facet(s) Addressed by the Triple Descriptors and FAIR Projectors:**
803       ●   **Findable** - Triple Descriptors, in isolation or when aggregated into FAIR
804            Profiles, provide one or more semantic interpretations of data elements. By
805            indexing these descriptors, it would become possible to search over datasets
806            for those that contain data-types of interest. Moreover, FAIR Projectors, as a
807            result of the TPF URI structure, create a unique URL for every data-point
808            within a record. This has striking consequences with respect to scholarly
809            communication. For example, it becomes possible to unambiguously refer-to,
810            and therefore "discuss" and/or annotate, individual spreadsheet cells from any
811            data repository.
812       ●   **Accessible** - Using the TPF design patterns, all data retrieval is
813            accomplished in exactly the same way - via HTTP GET. The response
814            includes machine-readable instructions that guide further exploration of the
815            data without the need to define an API. FAIR Projectors also give the data
816            owner high granularity access control; rather than publishing their entire
817            dataset, they can select to publish only certain components of that dataset,
818            and/or can put different access controls on different data elements, for
819            example, down to the level of an individual spreadsheet cell.
820       ●   **Interoperable** - FAIR Projectors provide a standardized way to export any
821            type of underlying data in a machine-readable structure, using widely used,
822            public shared vocabularies. Data linkages that were initially implicit in the
823            datastore, identifiers for example, become explicit when converted into URIs,
824            resulting in qualified linkages between formerly opaque data deposits.
825       ●   **Reusable** - All data points now possess unique identifiers, which allows them
826            to be explicitly connected to their citation and license information (i.e. the
827            MetaRecord). In this way, every data point, even when encountered in
828            isolation, provides a path to trace-back to its reusability metadata.

829
830    **Other features of FAIR Projection**
831        ● **Native formats are preserved** - As in many research domains,
832           bioinformatics has created a large number of data/file formats. Many of
833           these, especially those that hold "big data", are specially formatted flat-files
834           that focus on size-efficient representation of data, at the expense of general
835           machine-accessibility. The analytical tooling that exists in this domain,
836           therefore, is capable of consuming these various formats. While the FAIR
837           Data community has never advocated for wholesale Interoperable
838           representations of these kinds of data - which would be inefficient, wasteful,
839           and lacking in utility given that no tooling exists to consume such
840           representations - the FAIR Projector provides a middle-ground. Projection
841           allows software to query the core content of a file in a repository prior to
842           downloading it; for example, to determine if it contains data about an entity or
843           identifier of interest. FAIR Projectors, therefore, enable efficient efficient
844           discovery of data of-interest, without requiring wasteful transformation of all
845           data content into a FAIR format.
846        ● **Semantic conversion of existing Triplestores** - It is customary to re-cast
847           the semantic types of entities within triplestores using customized SPARQL
848           BIND or CONSTRUCT clauses. FAIR Projectors provide a standardized,
849           SPARQL-free, and discoverable way to accomplish the same task. This
850           further harmonizes data, and simplifies interoperability.
851        ● **Standardized interface to (some) Web Services** - Many Web Services in
852           the biomedical domain have a single input parameter, generally representing
853           an identifier for some biochemical entity. FAIR Projectors can easily replace
854           these myriad Web Services with a common TPF interface, thus dramatically
855           enhancing discoverability, machine-readability, and interoperability between
856           these currently widely disparate services.
857
858    **Incentives - why will this happen?**
859    Looking forward, there is every indication that FAIRness will be a requirement of funding
860    agencies and/or journals. As such, infrastructures such as the one described in this
861    exemplar will almost certainly become a natural part of scholarly data publishing in the
862    future. We indicated earlier, however, that we also believe that the creation of FAIR layers
863    over pre-existing data will become a natural part of the daily data transformation activities of
864    the global bioinformatics community. We suggest that this will happen because, though we
865    utilize RML in this demonstration only for its modelling properties, there exists tooling for a
866    wide variety of common file formats such as CSV and Excel that allow RML models to drive
867    the data transformation itself. As such, we predict that those who need to transform data will
868    begin to create, publish, and use RML models together with these generic transformation
869    tools to enact their data transformations, rather than continuing to write one-off scripts. This
870    may be incentivized even more by creating repositories of RML models that can be reused
871    by those needing to do data transformations. Though the infrastructure for capturing these
872    user-driven transformation events and formalizing them into FAIR Projectors does not yet

873 exist, it does not appear on its surface to be a complex problem. Thus, we expect that such
874 infrastructure should appear soon after FAIRness becomes a scholarly publishing
875 requirement.
876
877 Indeed, several communities of data providers are currently planning to use this, or related
878 FAIR implementations, to assist their communities to find, access, and reuse their valuable
879 data holdings.  For example, the Biobanking and Rare disease communities will be given
880 end-user tools that utilize/generate such FAIR infrastructures to:  guide discovery by
881 researchers; help both biobankers and researchers to re-code their data to standard
882 ontologies building on the SORTA system (Pang et al., 2015); assist to extend the
883 MOLGENIS/BiobankConnect system (Pang et al., 2016); add FAIR interfaces to the BBMRI
884 and RD-connect national and European biobank data and sample catalogues.  There are
885 also a core group of FAIR infrastructure authors who are creating large-scale indexing and
886 discovery systems that will facilitate the automated identification and retrieval of relevant
887 information, from any repository, in response to end-user queries, portending a day when
888 currently unused - "lost" - data deposits once again provide return-on-investment through
889 their discovery and reuse.
890

891 # Conclusions
892
893 There is a growing movement of governing bodies and funding organizations towards a
894 requirement for open data publishing, following the FAIR Principles. It is, therefore, useful to
895 have an exemplar "reference implementation" that demonstrates the kinds of behaviours that
896 are expected from FAIR resources.
897
898 Of the four FAIR Principles, Interoperability is arguably the most difficult FAIR facet to
899 achieve, and has been the topic of decades of informatics research. Several new standards
900 and frameworks have appeared in recent months that addressed various aspects of the
901 Interoperability problem. Here, we apply these in a novel combination, and show that the
902 result is capable of providing interoperability between formerly incompatible data formats
903 published anywhere on the Web. In addition, we note that the other three aspects of FAIR -
904 Findability, Accessibility, and Reusability - are easily addressed by the resulting
905 infrastructure. The outcome, therefore, provides machine-discoverable access to richly
906 described data resources in any format, in any repository, with the possibility of
907 interoperability of the contained data down to the level of an individual "cell". No new
908 standards or APIs were required; rather, we rely on RESTful behaviours, with all entities
909 being resolvable resources that allow hypermedia-driven "drill-down" from the level of a
910 repository descriptor, all the way to an individual data point in the record.
911
912 Such an interoperability layer may be created and published by anyone, for any data source,
913 without necessitating an interaction with the data owner. Moreover, the majority of the
914 interoperability layer we describe may be achieved through dynamically generated files from
915 software, or even (for the Accessor portion) through static, manually-edited files deposited in

916 any public repository. As such, knowledge of how to build or deploy Web infrastructure is not
917 required to achieve a large portion of these FAIR behaviors.
918
919 The trade-off between power and simplicity was considered acceptable, as a means to
920 hopefully encourage wide adoption. The modularity of the solution was also important
921 because, in a manner akin to crowdsourcing, we anticipate that the implementation will
922 spread through the community on a needs-driven basis, with the most critical resource
923 components being targeted early - the result of individual researchers requiring interoperable
924 access to datasets/subsets of interest to them. The interoperability design patterns
925 presented here provide a structured way for these individuals to contribute and share their
926 individual effort - effort they would have invested anyway - in a collaborative manner, piece-
927 by-piece building a much larger interoperable and FAIR data infrastructure to benefit the
928 global community.

929

## Acknowledgements

## References

951 Bechhofer S., Buchan I., De Roure D., Missier P., Ainsworth J., Bhagat J., Couch P.,

952      Cruickshank D., Delderfield M., Dunlop I., Gamble M., Michaelides D., Owen S.,

953      Newman D., Sufi S., Goble C. 2013. Why linked data is not enough for scientists. *Future*

954     *generations computer systems: FGCS* 29:599–611.

955    Berners-Lee T. 2006.Linked Data. *Available at*

956     *https://www.w3.org/DesignIssues/LinkedData.html* (accessed September 27, 2016).

957    Cook CE., Bergman MT., Finn RD., Cochrane G., Birney E., Apweiler R. 2016. The

958     European Bioinformatics Institute in 2016: Data growth and integration. *Nucleic acids*

959     *research* 44:D20–6.

960    Covitz PA., Hartel F., Schaefer C., De Coronado S., Fragoso G., Sahni H., Gustafson S.,

961     Buetow KH. 2003. caCORE: a common infrastructure for cancer informatics.

962     *Bioinformatics*  19:2404–2412.

963    Crosswell LC., Thornton JM. 2012. ELIXIR: a distributed infrastructure for European

964     biological data. *Trends in biotechnology* 30:241–242.

965    Das S., Sundara S., Cyganiak R. 27 September, 2012. *R2RML: RDB to RDF Mapping*

966     *Language*. W3C Recommendation.

967    De Giovanni R., Copp C., Döring M., Hobern D. 2010. *TAPIR - TDWG Access Protocol for*

968     *Information Retrieval*. TSWG Standards .

969    Dimou A., Vander Sande M., Colpaert P., Verborgh R., Mannens E., de Walle R. RML: A

970     Generic Language for Integrated RDF Mappings of Heterogeneous Data. In:

971     *Proceedings of the 7th Workshop on Linked Data on the Web*.

972    Dimou A., Sande MV., Slepicka J., Szekely P., Mannens E., Knoblock C., Walle RV de.

973     Mapping Hierarchical Sources into RDF Using the RML Mapping Language. In: *2014*

974     *IEEE International Conference on Semantic Computing*. IEEE, 151–158.

975    Dumontier M., Gray AJG., Marshall MS., Alexiev V., Ansell P., Bader G., Baran J., Bolleman

976     JT., Callahan A., Cruz-Toledo J., Gaudet P., Gombocz EA., Gonzalez-Beltran AN.,

977     Groth P., Haendel M., Ito M., Jupp S., Juty N., Katayama T., Kobayashi N.,

978     Krishnaswami K., Laibe C., Le Novère N., Lin S., Malone J., Miller M., Mungall CJ.,

979     Rietveld L., Wimalaratne SM., Yamaguchi A. 2016. The health care and life sciences

980    community profile for dataset descriptions. *PeerJ* 4:e2331.

981  Fielding RT., Taylor RN. 2002. Principled design of the modern Web architecture. *ACM*

982    *Transactions on Internet Technology* 2:115–150.

983  González AR., Callahan A., Cruz-Toledo J., Garcia A., Egaña Aranguren M., Dumontier M.,

984    Wilkinson MD. 2014. Automatically exposing OpenLifeData via SADI semantic Web

985    Services. *Journal of biomedical semantics* 5:46.

986  Ison J., Kalas M., Jonassen I., Bolser D., Uludag M., McWilliam H., Malone J., Lopez R.,

987    Pettifer S., Rice P. 2013. EDAM: an ontology of bioinformatics operations, types of data

988    and identifiers, topics and formats. *Bioinformatics* 29:1325–1332.

989  Kuhn T., Tobias K., Christine C., Michael K., Núria Q-R., Ruben V., George G., Ngomo A-

990    CN., Raffaele V., Michel D. 2016. Decentralized provenance-aware publishing with

991    nanopublications. *PeerJ Computer Science* 2:e78.

992  Lanthaler M., Gütl C. Hydra: A Vocabulary for Hypermedia-Driven Web APIs. In:

993    *Proceedings of the 6th Workshop on Linked Data on the Web (LDOW2013)*.

994  Maali F., Erickson J., Archer P. 2014. *Data Catalog Vocabulary (DCAT)*. W3C

995    Recommendation .

996  Miles A., Bechhofer S. 18 August, 2009. *SKOS Simple Knowledge Organization System*

997    *Reference*. W3C Recommendation .

998  van Ommen G-JB., Törnwall O., Bréchot C., Dagher G., Galli J., Hveem K., Landegren U.,

999    Luchinat C., Metspalu A., Nilsson C., Solesvik OV., Perola M., Litton J-E., Zatloukal K.

1000   2015. BBMRI-ERIC as a resource for pharmaceutical and life science industries: the

1001   development of biobank-based Expert Centres. *European journal of human genetics:*

1002   *EJHG* 23:893–900.

1003  Pang C., Sollie A., Sijtsma A., Hendriksen D., Charbon B., de Haan M., de Boer T., Kelpin

1004   F., Jetten J., van der Velde JK., Smidt N., Sijmons R., Hillege H., Swertz MA. 2015.

1005   SORTA: a system for ontology-based re-coding and technical annotation of biomedical

1006        phenotype data. *Database: the journal of biological databases and curation* 2015. DOI:

1007            10.1093/database/bav089.

1008    Pang C., van Enckevort D., de Haan M., Kelpin F., Jetten J., Hendriksen D., de Boer T.,

1009            Charbon B., Winder E., van der Velde KJ., Doiron D., Fortier I., Hillege H., Swertz MA.

1010            2016. MOLGENIS/connect: a system for semi-automatic integration of heterogeneous

1011            phenotype data with applications in biobanks. *Bioinformatics*  32:2176–2183.

1012    Roche DG., Kruuk LEB., Lanfear R., Binning SA. 2015. Public Data Archiving in Ecology and

1013            Evolution: How Well Are We Doing? *PLoS biology* 13:e1002295.

1014    SIB Swiss Institute of Bioinformatics Members. 2016. The SIB Swiss Institute of

1015            Bioinformatics' resources: focus on curated databases. *Nucleic acids research* 44:D27–

1016            37.

1017    Speicher S., Arwe J., Malhotra A. 2015. *Linked data platform 1.0*. W3C Recommendation.

1018    Starr J., Castro E., Crosas M., Dumontier M., Downs RR., Duerr R., Haak LL., Haendel M.,

1019            Herman I., Hodson S., Hourclé J., Kratz JE., Lin J., Nielsen LH., Nurnberger A., Proell

1020            S., Rauber A., Sacchi S., Smith A., Taylor M., Clark T. 2015. Achieving human and

1021            machine accessibility of cited data in scholarly publications. *PeerJ. Computer science* 1.

1022            DOI: 10.7717/peerj-cs.1.

1023    Stein LD., Knoppers BM., Campbell P., Getz G., Korbel JO. 2015. Data analysis: Create a

1024            cloud commons. *Nature* 523:149–151.

1025    Thompson R., Johnston L., Taruscio D., Monaco L., Béroud C., Gut IG., Hansson MG., 't

1026            Hoen P-BA., Patrinos GP., Dawkins H., Ensini M., Zatloukal K., Koubi D., Heslop E.,

1027            Paschall JE., Posada M., Robinson PN., Bushby K., Lochmüller H. 2014. RD-Connect:

1028            an integrated platform connecting databases, registries, biobanks and clinical

1029            bioinformatics for rare disease research. *Journal of general internal medicine* 29 Suppl

1030            3:S780–7.

1031    Tim Berners-lee T., Chen Y., Chilton L., Connolly D., Dhanaraj R., Hollenbach J., Lerer A.,

1032   Sheets D. 2006. Tabulator: Exploring and analyzing linked data on the semantic web.

1033   In: *Proceedings of the 3rd International Semantic Web User Interaction Workshop*.

1034 Verborgh R., Ruben V., Sande MV., Olaf H., Van Herwegen J., De Vocht L., De Meester B.,

1035   Gerald H., Pieter C. 2016. Triple Pattern Fragments: A low-cost knowledge graph

1036   interface for the Web. *Web Semantics: Science, Services and Agents on the World*

1037   *Wide Web* 37-38:184–206.

1038 Verborgh R., Dumontier M. 2016. A Web API ecosystem through feature-based reuse.

1039 Wilkinson MD., Dumontier M., Aalbersberg IJJ., Appleton G., Axton M., Baak A., Blomberg

1040   N., Boiten J-W., da Silva Santos LB., Bourne PE., Bouwman J., Brookes AJ., Clark T.,

1041   Crosas M., Dillo I., Dumon O., Edmunds S., Evelo CT., Finkers R., Gonzalez-Beltran A.,

1042   Gray AJG., Groth P., Goble C., Grethe JS., Heringa J., 't Hoen PAC., Hooft R., Kuhn T.,

1043   Kok R., Kok J., Lusher SJ., Martone ME., Mons A., Packer AL., Persson B., Rocca-

1044   Serra P., Roos M., van Schaik R., Sansone S-A., Schultes E., Sengstag T., Slater T.,

1045   Strawn G., Swertz MA., Thompson M., van der Lei J., van Mulligen E., Velterop J.,

1046   Waagmeester A., Wittenburg P., Wolstencroft K., Zhao J., Mons B. 2016. The FAIR

1047   Guiding Principles for scientific data management and stewardship. *Scientific data*

1048   3:160018.

1049