

# DM2E: A Linked Data Source of Digitised Manuscripts for the Digital Humanities

**Editor(s):** Christoph Schlieder, Universität Bamberg, Germany

**Solicited review(s):** Günther Görz, Friedrich-Alexander-Universität Erlangen-Nürnberg, Germany; Bernhard Haslhofer, Austrian Institute of Technology, Austria; Eetu Mäkelä, Aalto University, Finland & University of Helsinki, Finland & University of Oxford, UK

Konstantin Baierer<sup>a</sup>, Evelyn Dröge<sup>a</sup>, Kai Eckert<sup>b,\*</sup>, Doron Goldfarb<sup>c</sup>, Julia Iwanowa<sup>a</sup>,  
Christian Morbidoni<sup>d</sup>, Dominique Ritze<sup>e</sup>

<sup>a</sup> *Berlin School of Library and Information Science, HU Berlin, Unter den Linden 6, 10099 Berlin, Germany*  
*E-mail: first.last@ibi.hu-berlin.de*

<sup>b</sup> *Stuttgart Media University, Nobelstr. 10, 70569 Stuttgart, Germany*  
*E-mail: eckert@hdm-stuttgart.de*

<sup>c</sup> *Austrian National Library, Josefsplatz 1, 1015 Vienna, Austria*  
*E-mail: fue@onb.ac.at*

<sup>d</sup> *Dept. of Information Engineering, Università Politecnica delle Marche, Via Brecce Bianche, 60131 Ancona, Italy*  
*E-mail: christian.morbidoni@gmail.com*

<sup>e</sup> *Research Group Data and Web Science, University of Mannheim, B6, 26, 68159 Mannheim, Germany*  
*E-mail: dominique@informatik.uni-mannheim.de*

**Abstract.** The DM2E dataset is a five-star dataset providing metadata and links for direct access to digitized content from various cultural heritage institutions across Europe. The data model is a true specialization of the Europeana Data Model and reflects specific requirements from the domain of manuscripts and old prints, as well as from developers who want to create applications on top of the data. One such application is a scholarly research platform for the Digital Humanities that was created as part of the DM2E project and can be seen as a reference implementation. The Linked Data API was developed with versioning and provenance from the beginning, leading to new theoretical and practical insights.

**Keywords:** Linked Data, Dataset, Cultural Heritage, Digital Humanities, Digital Content, Europeana, EDM, DM2E

## 1. Introduction

The project “Digitised Manuscripts to Europeana” (DM2E)<sup>1</sup> was active from 02/2012 until 01/2015, funded under EU FP7. Its two primary goals were:

1. The transformation of various metadata and content formats describing and representing digital cultural heritage objects (CHOs) in the realm of digitized manuscripts from as many providers

(cf. Section 2) as possible into the Europeana Data Model (EDM) to get it into Europeana, the European digital library.<sup>2</sup>

2. The stable provision of the data as Linked Data and the creation of tools and services to reuse the data in the Digital Humanities. The basis is the possibility to annotate the data, to link the data, and to share the results as new data.

---

\*Corresponding author, e-mail: eckert@hdm-stuttgart.de

<sup>1</sup><http://dm2e.eu/> (10.12.2015)

<sup>2</sup><http://www.europeana.eu> (10.12.2015)

Table 1  
DM2E data sources

Provider	Collection	ML <sup>a</sup>	CL <sup>b</sup>	Type <sup>c</sup>	Count	Vocab <sup>d</sup>	Format	
Berlin Brandenburg Academy of Sciences	Deutsches Textarchiv	de	de	B	1547	1	TEI	
University of Bergen	Wittgenstein Archive Bergen	en	de,en,var.	M	20			
Bulgarian Academy of Sciences	Codex Suprasliensis	en	cu	M	49	2		
Humboldt University Berlin	Polytechnisches Journal	de	de	J	42173	1		
ERC AdG EUROCORR	European Correspondence to Jacob Burckhardt	en	de	L	497	8	METS/ MODS	
University Library JCS Frankfurt am Main	Medieval Manuscript Collection	de	la,de,var.	M	634	1,2,5, 6,7		
	Hebrew Manuscript Collection	de	he,var.	M	378			
	Modern Manuscripts	de	la,it,de,var.	M	279			
	Oriental Manuscripts	de	gez,la,cu,var.	M	29			
	Max Horkheimer Estate	de	var.	A	272			
Georg Eckert Institute for Textbook Research	GEI-Digital	de	de,la	B	3147	1,7		
Brandeis University Library via EAJC <sup>e</sup>	Spanish Civil War Posters	en	es	I	112	4		MARC
Center for Jewish History via EAJC	YIVO Institute for Jewish Research Collection	en	he,yi,ru,var.	B/A	3987			
	Leo Baeck Institute Collection	en	de,en,var.	M/B/J/A	7885			
National Library of Israel	Hebrew & various language Manuscripts	he	he,ar,var.	M	1296	1		
	Hebrew, Yiddish & various language Books	he,en	yi,he,var.	B	7722			
	Archival Material	he	de,he,en,var.	A	2775			
Berlin State Library	Personal Papers of Adelbert von Chamisso	de	de,fr,var.	M/L	4662	1	EAD	
	Personal Papers of Gerhart Hauptmann	de	de,var.	M/L	14295			
	Publisher Archives of Gebauer & Schwetschke	de	de	M/L	43296			
	Western Manuscripts	de	de,la,var.	M	163			
Joint Distribution Committee (JDC) via EAJC	Records of the NYC Office of the JDC, 1914-18	en	en,var.	F	207	1	MAB2	
Austrian National Library	Austrian Books Online	de	de,it,fr,var.	B/J	44425			
	Codices	de	la,de,tr,var.	M	175			
Max Planck Institute for the History of Science	Islamic Scientific Manuscripts Initiative	en	ar	M	763	1,3,4	Custom	
	MPIWG Digital Rare Book Library	en	la,fr,de,var.	M/B	1264			
	The manuscripts of Thomas Harriot	en	en,var.	M	24			
Petőfi Literary Museum	A Tett Magazine	hu	hu	J	183	7	DC	

<sup>a</sup> ML: Metadata Language <sup>b</sup> CL: Content Language

<sup>c</sup> M: Manuscripts / L: Letters / B: Books / I: Images / J: Journal Articles / A: Archival Items / F: Archival File

<sup>d</sup> 1: GND / 2: DBpedia / 3: DDC / 4: LCSH / 5: ZDB / 6: Geonames / 7: VIAF / 8: Freebase

<sup>e</sup> European Association for Jewish Culture

The Linked Data representation of the metadata as described in this paper can be accessed online<sup>3</sup> and is as part of the LOD cloud also registered on Datahub.<sup>4</sup> DM2E is a five-star Linked Data source adhering to the Linked Data Principles [3], i.e., it uses dereferenceable URIs, provides all metadata in RDF using proper content negotiation together with links to other Linked Data sources. The vocabulary achieves four stars of the Five Stars of Linked Data Vocabulary Use [12], cf. Section 3. The data is not only provided as an end in itself, but forms the basis for a scholarly research platform allowing scholars to access the underlying content to annotate it and link it to other sources (Section 4). In order to support the scholars in finding relevant content, the RDF data is enriched – contextualized – as part of the ingestion process (Section 5). A specialty is the provision of full provenance of the data and the support of versioning, as described in Section 6. All the RDF data provided by DM2E can be used without restrictions in accordance with the CC0 public domain dedication.<sup>5</sup>

The rights statements for the described digitized objects are individually assigned by the content providers who have to choose an appropriate statement from the options<sup>6</sup> offered by Europeana and attach this information to each individual item.

## 2. Sources

One major aspect of the DM2E project was publishing metadata about a number of international high profile collections both as Linked Data and through Europeana. Despite its name, DM2E is not restricted to manuscripts but contains also other historical resources like letters, books, images, journal articles, or archival items. Table 1 shows an overview on the content available as Linked Data, broken down by provider and collection name, metadata and content language, type of content, instance count, used reference authorities and metadata source format. The stated counts represent the respective number of instances for which the property

<sup>3</sup><http://data.dm2e.eu/data> (10.12.2015)

<sup>4</sup><http://datahub.io/dataset/dm2e> (10.12.2015)

<sup>5</sup><http://creativecommons.org/publicdomain/zero/1.0/> (10.12.2015)

<sup>6</sup><http://pro.europeana.eu/available-rights-statements> (10.12.2015)

`dm2e:displayLevel` is set to `true` (see Section 3). As can be seen, this dataset is based on the integration of a variety of source metadata formats, reflecting the heterogeneity of the underlying materials, their international character and the flexibility of the DM2E model to effectively represent such diverse content.

Accordingly, there is no common workflow for providers to map their data to the DM2E model. Starting with the tools and the documentation provided by DM2E, providers therefore developed their own metadata transformations mainly based on XSLT, although some chose to directly implement export routines into their collection management systems. Initial consistency checks for mapped data revealed that despite the very detailed specification of the DM2E model some providers showed great creativity in individual interpretations of specific model features, most notably regarding the representation of hierarchies. In order to maintain a homogeneous data representation, specific mapping rules have been established for such cases and distributed amongst the providers in form of a recommendation document. In some cases, transformations created by one provider could be reused or adapted for other providers. This especially proved to be effective for the highly standardized library metadata formats such as MARC, METS/MODS and MAB2. The mapping recommendations and the resulting metadata crosswalks are documented on the DM2E wiki<sup>7</sup>, a more detailed description of the individual transformation workflows is available as project deliverable [4].

### 3. Data Model

The DM2E model is an application profile of EDM, i.e., an application-specific specialization for the representation of manuscripts and similar historical content like old prints, posters, books and old journals [6]. EDM itself is very generic to represent resources provided by museums, libraries, archives and galleries all over Europe. It is based on top-level ontologies like OAI-ORE, Dublin Core and SKOS. Core classes are `edm:ProvidedCHO` for the described cultural heritage object (CHO), `ore:Aggregation` for the metadata record provided for the described CHO and `edm:WebResource` for views of the described CHO, such as images. CHOs can be further quali-

fied by links to contextual resources being instances of `edm:Agent`, `edm:TimeSpan`, `edm:Place`, or `skos:Concept`.

The example of a manuscript by the philosopher Ludwig Wittgenstein shown in Figure 1 illustrates how DM2E data is built-up. Bold resources have been added in the DM2E model, others are part of the underlying EDM. `ore:Aggregation` shows, for example, who has created and mapped the metadata and where the CHO is shown on the Web, while `edm:ProvidedCHO` is about the physical object that is described. The DM2E model allows that CHOs have multiple hierarchical layers. The CHO here has two layers: the manuscript and paragraphs within the manuscript. The type of the CHO is given via `dc:type` in `edm:ProvidedCHO`. Agents are divided into organizations and persons and can be further described. `dm2e:levelOfHierarchy "1"` says that the manuscript is the highest hierarchical level of this object. Hierarchies are very collection-specific; usually the provider knows best which level is most relevant for scholars. Therefore, the property `dm2e:displayLevel` is used to give applications a hint, if this CHO for example should appear in a result list of a search interface or if it should show up in Europeana. Users searching for Wittgenstein usually do not want to see every paragraph of a Wittgenstein manuscript in the search results.

The DM2E model adds mostly subclasses and -properties for the domain of manuscripts to existing elements of the EDM. Main additions have been made for person roles, e.g. `dm2e:composer`, `pro:author` and more properties that specialize `dc:creator`, classes that are used to indicate the CHO's type like `bibo:Letter`, `dm2e:Manuscript` or `fabio:Article` or properties to describe specifics of manuscripts or other document types of the model, e.g. `dm2e:incipit` for the opening words of a manuscript or `dm2e:receivedOn` for a date on which a letter was received. An hierarchical object is described on every level using `edm:ProvidedCHO` and `ore:Aggregation` as the metadata vary between pages or even paragraphs. Additionally, this allows to refer to every CHO as an independent object if desired. The model was created and refined in an iterative, agile process taking into account several mapping workshops and constant feedback by the data providers and application developers. This feedback led to model changes where properties or classes have been added or dropped, or property ranges

<sup>7</sup><http://wiki.dm2e.eu> (10.12.2015)

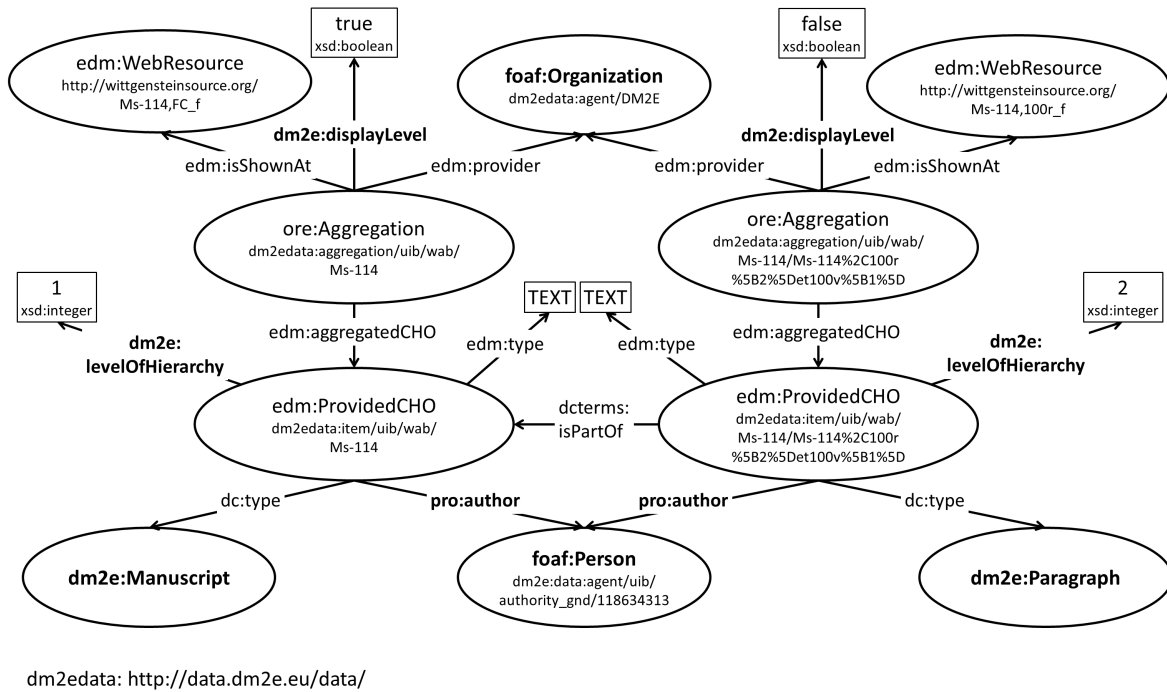


Fig. 1. Example of the DM2E model in use: representation snippet of a Wittgenstein manuscript.

have been adapted. After the intellectual collection of requirements and resulting initial versions of the DM2E model, an evaluation of the mapped data [2] has been conducted to get information about the actual usage of the model. As a result, many properties and classes of the model have been removed despite originally been asked for by providers as they have not been used. In the appendix (Table 3), we provide an overview on the usage of the metadata fields for the class `edm:ProvidedCHO` across the datasets.

Whenever possible, established vocabularies have been reused, precisely: BIBFRAME, BIBO, CIDOC-CRM, FABIO, PRO, rdaGr2, VIVO and VoID. DM2E-specific usage guidelines for each reused element are provided via `dm2e:scopeNote`. On the five stars scale for LOD vocabulary use proposed by [12], the model gets four stars, as it is dereferencable and machine-readable, linked to other vocabularies, has metadata about it but is not (yet) linked to by other vocabularies. As of January 20, 2015, the DM2E model contains 65 additional properties and 23 additional classes. The DM2E dataset currently includes descriptions for 2,670,996 cultural heritage objects, 2,478,765 of which representing single annotatable pages, while 182,259 have `displayLevel` set. Regarding contextual resources, 33,080 objects

of type `skos:Concept` are available, 37,772 are typed as `edm:TimeSpan`, 21,304 as `edm:Agent`, 3,751 as `foaf:Organization`, 104,779 as `foaf:Person` and 27,266 as `edm:Place`. Compared to the EDM data that is available via the Europeana LOD pilot [11], the specialized DM2E data forms a smaller, complementary dataset with RDF statements at a very detailed level. The namespace URI for the DM2E model schema is <http://onto.dm2e.eu/schemas/dm2e/> and [data.dm2e.eu/data/](http://data.dm2e.eu/data/) for instances. The full documentation of the model, including detailed changelogs between model versions, can be accessed via the DM2E wiki.

#### 4. Application

Two applications consuming the DM2E Linked Data have been implemented in the DM2E project to provide the scholarly research platform. The first one is a faceted browser that allows scholars to make sense of the DM2E collections and navigate them along several dimensions, iteratively restricting the search results by language, author, publishing institutions, and

other metadata fields. Facets are derived from a SOLR<sup>8</sup> index populated by running queries against the DM2E SPARQL endpoint. To populate such an index we used the approach and implementation described in [13]. This is based on SPARQL configurations to derive document id, facets and corresponding values from pre-defined graph patterns. It is easy for example to get all datatype properties of a resource and turn them into corresponding facets. In this way facets names (fields in Solr) are not known a priori, so we use dynamic fields: in the following example fields ending with `ss` are multivalued string fields derived from a RDF property. In the prototype implementation we use the `dm2e:displayLevel` property to exclude hierarchical levels that should not show up in search results. As we currently do not provide a public SPARQL endpoint, the RESTful Solr search API is also an important building block to search the DM2E data programmatically. It straightforwardly provides a fast full text search over the main metadata fields, e.g. by using `q=Bartolomeo` as URL query string.<sup>9</sup> Furthermore it can be used to perform non relational queries, such as "all the cultural objects that are written in Italian and issued at a precise year, e.g. 1831", by using the following query string `fq=language_ss:it&fq=issued_ss:18*&wt=json`.<sup>10</sup>

The API response is a JSON array, from which the dereferenceable URL of Linked Data resources are easily retrievable by looking at the "id" slot. Solr does not provide a direct way to get all the facets used in the index, which would facilitate the developer in writing queries. However such a list can be retrieved with a work-around with the following query: `q=*:*&wt=csv&rows=0&facet`.<sup>11</sup>

Based on the Solr API, a end-user faceted browser was developed by customising Ajax-Solr. Each resource shown in the faceted browser can be opened in its provider's own digital library (following the EDM property `edm:isShownAt`) and its Linked Data representation can be reached by clicking on "see RDF data" that points to its dereferenceable URL. The

DM2E faceted search is available online<sup>12</sup> and its usage demonstrated in a short online screencast.<sup>13</sup>

Furthermore, for datasets containing links to annotatable digital objects, users are provided with "Annotate with Pundit" links that direct them to the DM2E semantic annotation environment. The latter constitutes the second web application built on top of the data and is based on Pundit and Feed, two software components developed as part of the DM2E project.

Pundit<sup>14</sup> is an annotation tool that allows users to enrich web pages with semantically structured data. In DM2E substantial improvements were done over the previous versions [9] [10], leading to a completely new user interface and additional annotation functionalities [14]. Annotations in Pundit encode machine readable semantic connections among images, texts and LOD entities in form of RDF triples (using the Open Annotation data model),<sup>15</sup> consumable via SPARQL or a dedicated REST API.<sup>16</sup>

On the other hand, the Feed REST API provides access to Pundit "as a service" by using the URL of a Web page to be annotated as a call parameter. An extension to the Feed API has been developed in DM2E, allowing this call parameter to also be a dereferenceable URL of an RDF description of a digitized object. Feed parses this RDF description to create a customized annotation environment. While the faceted search application only works on resources at display level (e.g. entire books), the annotation environment allows users to go deeper into the hierarchy. The `dcterms:isPartOf` and `edm:isNextInSequence` properties are used to provide basic navigation functionalities such as reaching a specific page of a manuscript or going to the next/previous pages. The actual digital contents that users can annotate are retrieved by following `dm2e:hasAnnotatableVersionAt` links or, alternatively, the `edm:object` links. There can be multiple annotatable objects associated with a resource as, for example, the facsimile image and its HTML transcription. In this case, both contents are shown and made annotatable.<sup>17</sup> In case of the presence of links to

<sup>8</sup><http://lucene.apache.org/solr/> (10.12.2015)

<sup>9</sup><http://141.20.126.236:8080/solr-dm2e/collection1/select?q=Bartolomeo>

<sup>10</sup>[http://141.20.126.236:8080/solr-dm2e/collection1/select?q=\\*:\\*&fq=language\\_ss:it&fq=issued\\_ss:18\\*&wt=json](http://141.20.126.236:8080/solr-dm2e/collection1/select?q=*:*&fq=language_ss:it&fq=issued_ss:18*&wt=json)

<sup>11</sup>[http://141.20.126.236:8080/solr-dm2e/select?q=\\*:\\*&wt=csv&rows=0&facet](http://141.20.126.236:8080/solr-dm2e/select?q=*:*&wt=csv&rows=0&facet)

<sup>12</sup><http://purl.org/net/dm2e/search> (10.12.2015)

<sup>13</sup>[https://youtu.be/\\_rQ\\_7NhewhQ](https://youtu.be/_rQ_7NhewhQ) (10.12.2015)

<sup>14</sup><https://thepund.it/> (10.12.2015)

<sup>15</sup><http://www.openannotation.org/spec/core/> (10.12.2015)

<sup>16</sup>Pundit server API documentation: <http://net7.github.io/pundit2/rest-api.html> (10.12.2015)

<sup>17</sup>Annotatable page example: <http://bit.ly/1FIu6dA> (10.12.2015)

popular LOD datasets such as DBpedia,<sup>18</sup> Feed is able to gather additional metadata (e.g. full names, descriptions, categories) in order to provide additional context to scholars. For this purpose, we established automated contextualization processes, as described in the next section.

By annotating digital objects with Pundit, users in fact create additional RDF knowledge. This could be, for example, links connecting a geographical map of Metz, depicted in a manuscript page, to the city of Metz in DBpedia, or links from a sentence of a manuscript transcription to a DBpedia entity that is mentioned in such a sentence. Such RDF data, in turn, can be indexed to enrich the faceted search interface, thus improving search and discovery. Demonstrative examples of such end-user information enrichments can be seen in an online screencast.<sup>19</sup>

As of September 30, 2014, about 6,600 annotations for about 900 digital objects from the DM2E dataset have been created by scholars using our research platform.

## 5. Contextualization

Linking our datasets to external sources like GND,<sup>20</sup> DBpedia, Geonames,<sup>21</sup> or the Library of Congress Subject Headings<sup>22</sup> enables to easily get information about a resource, either directly by following the link to the external source or by detecting connections between resources based on the same links. While the links to GND often are already present in the original metadata, links to all other sources are generated automatically. To create the links, we use the link discovery framework Silk.<sup>23</sup> Silk generates links based on a linkage rule that is provided by the user. Such a linkage rule specifies the conditions which have to hold to create a link, e.g. the names of two resources need to have a Jaccard measure value above 0.8. All links that are currently in our system are generated with the same configuration which compares the labels using the Jaro Winkler distance and requires a confidence value of 0.9 aiming at a high precision while tolerating spelling

<sup>18</sup><http://wiki.dbpedia.org/> (10.12.2015)

<sup>19</sup><https://youtu.be/tUrdLm43CMA> (10.12.2015)

<sup>20</sup>[http://www.dnb.de/DE/Standardisierung/GND/gnd\\_node.html](http://www.dnb.de/DE/Standardisierung/GND/gnd_node.html) (10.12.2015)

<sup>21</sup><http://www.geonames.org> (10.12.2015)

<sup>22</sup><http://id.loc.gov/authorities/subjects.html>

<sup>23</sup><http://silkframework.org/> (10.12.2015)

Table 2  
Number of links per external source

dbpedia	freebase	geonames	judaica	lcsb	geodata	nytimes	GND
12287	1868	1571	1474	141	5770	570	22698

variations. This might seem like a very simple method, but in most cases we have no further information in the metadata besides a simple string. Our evaluation suggests, however, that even this simple method leads to good results, as many of these strings are not ambiguous. Table 2 shows the number of generated links to each external sources. We link agents, places, and subjects.

Altogether, about 24,000 links have been automatically generated. With a manual analysis of 150 random links from agents to DBpedia and 150 random links from places to Linked Geodata, we evaluate the quality. For agents, 125 correct links have been detected which results in a precision of 0.83. Since DBpedia covers several labels, we can for example correctly link “Jakobä” to the DBpedia agent “Jacqueline Countess of Hainaut.” The incorrect links result either from ambiguous names, e.g. “Heinrich Fischer” who refers to a Swiss rower in DBpedia and not to an author, or from incomplete information, e.g. if only the first name or surname is given.

For places, 128 correct links can be detected, resulting in a precision of 0.85 which is similar to the precision for agents. Since Linked Geodata includes labels in various languages, even places with a German label such as “München” can be linked to “Munich”. The reasons for incorrect links, however, are the same to the ones for agents, e.g. the German city “Heidelberg” is mapped to the city “Heidelberg” located in South Africa due to identical labels.

Across all datasets, about 18% of all agents and 60% of all places are linked on average. With a different linkage rule it is possible to detect more links but with the risk to reduce the precision. Further, the amount of detected links as well as their quality highly depend on the popularity and currency of the resources. Since more than one linkset can be available in our system and the user can track their provenance, more liberal linkage rules can also be applied and the user can be informed about its quality.

## 6. Implementation

At the core of DM2E’s infrastructure is a Jena TDB<sup>24</sup> triplestore, accessible by a Jena Fuseki SPARQL endpoint. After evaluating a few different RDF storage solutions, this combination offered the perfect balance between maintainability, scalability and versatility. In fact, all DM2E’s internal applications and the infrastructure are interfacing with the data exclusively through the SPARQL endpoint, making, in theory, the actual triplestore implementation interchangeable. The RDF data is partitioned into Named Graphs that correspond to individual ingestions (see also Section 6.3 – Dataset Provenance); exporting and importing N-Quad dumps of the full data store as well as specific subsets is straight-forward. Dumps for the DM2E metadata collection<sup>25</sup> and for the contextualized external entities<sup>26</sup> are available to the public.

### 6.1. Data Ingestion

There are two user interfaces that allow data providers or mapping institutions to deliver data to DM2E: A Linked Data-based workflow engine with an HTML5 web interface allows casual users to test their transformations and the ingestion process (Omnom)<sup>27</sup> while a set of command line tools is targeted at power users doing large-scale ingestions and conversions (dm2e-data.sh).<sup>28</sup>

Omnom is centered on the idea that RDF’s flexible graph-based structure combined with the semantic expressivity of ontologies<sup>29</sup> not only allows the definition and execution of intelligent workflows, automating tedious, long-running and error-prone tasks, but solves the problem of tracking data provenance [8]. Combined with the simple Web User Interface,<sup>30</sup> Omnom can be very helpful for the technically-non-too-

savvy to understand the processes of data mapping, data transformation and data ingestion and iteratively improve their own workflow, though Omnom’s approach to use and persist RDF for all data does lead to suboptimal performance when doing full-scale transformations/ingestions.

The command line suite of tools is developed with a server environment in mind and consists of a set of Java tools for DM2E validation, provenance-tracking data ingestion, DM2E-EDM-conversion and EDM validation, as well as shell scripts encapsulating XSLT transformers and RDF serializers and for orchestrating the various operations.

The authoritative source of the DM2E model is the textual/tabular *DM2E Model Specification*, which contains not only the definitions of all properties and classes to be used, but illustrates their usage with examples. The specs are synchronously formalized as an dereferenceable OWL ontology. However, the DM2E model puts restrictions on the usage of properties and classes that cannot be expressed under OWL’s Open World Assumption. These restrictions are targeted towards structural validation of subgraphs of DM2E data rather than inference of new facts. While DM2E is involved in the development of community standards for RDF validation [5], we implemented a custom solution using Java, available on GitHub.<sup>31</sup> While the validation tool is “hard-wired” to DM2E’s model, it is rather meticulous and has proven useful not only for discovering outright model violations (e.g. wrong cardinality of properties or missing conditional statements) but stylistic problems such as unwise characters in URIs and labels or variations in the UTF-8 normalization.

### 6.2. Delivery to Europeana

Being a domain aggregator for Europeana, DM2E has a strong focus on interoperability with the EDM, both on the model and data level. The DM2E model is a specialization of the EDM, i.e., after RDFS inference on the data and removing any statements with properties not contained in EDM, every DM2E-compliant subgraph is an EDM-compliant subgraph. DM2E uses this technique to convert the DM2E data into pure EDM to make the ingestion as easy as possible for the Europeana side, using a synthesis of the two models

<sup>24</sup><http://jena.apache.org> (10.12.2015)

<sup>25</sup>DM2E collection metadata dump: <http://data.dm2e.eu/dm2e-fuseki-direct.2016-01-07.final.nquads.gz> (13.01.2016)

<sup>26</sup>DM2E contextualized links dump: <http://data.dm2e.eu/dm2e-fuseki-direct.2016-01-07.linksets.nquads.gz> (13.01.2016)

<sup>27</sup><http://omnom.dm2e.eu>

<sup>28</sup><https://github.com/DM2E/dm2e-ontologies/blob/master/src/main/bash/dm2e-data.sh> (10.12.2015)

<sup>29</sup><http://onto.dm2e.eu/omnom> (10.12.2015), <http://onto.dm2e.eu/omnom-types> (10.12.2015)

<sup>30</sup><https://github.com/DM2E/dm2e-gui> (10.12.2015)

<sup>31</sup><https://github.com/DM2E/dm2e-ontologies> (10.12.2015)

expressed in OWL.<sup>32</sup> As the last step before delivery to Europeana, the produced EDM representations are validated using a combination of XML Schema and Schematron.<sup>33</sup>

Due to its ubiquitous deployment in the GLAM sector and its proven track record for scalability, DM2E and Europeana agreed on OAI-PMH as the preferred mode of delivery of data for ingestion into Europeana. Using a multi-step process of extracting per-ore:Aggregation-subgraphs from the triplestore, validation against the DM2E model, conversion to EDM, data massaging and validation against the EDM model,<sup>34</sup> an EDM dump of all data in DM2E is created monthly. With OAI-PMH set names corresponding to datasets, these EDM RDF/XML files are then served using the Repox OAI-PMH repository.

### 6.3. Linked Data API

The Linked Data API is implemented using a significantly advanced version of Pubby.<sup>35</sup> The source code for this DM2E-specific version is available via GitHub.<sup>36</sup> An integration of the additional features – which are of general interest – into the main branch of Pubby is planned. The basis for all of them is unleashing the power of SPARQL by allowing arbitrary URI patterns to be mapped to customized SPARQL queries. In the following, we describe how the requirements regarding data access have been accomplished for the DM2E data within Pubby.

*Multiple resource handling.* DM2E implements the OAI-ORE resource map, i.e., whenever the URI of a resource or an aggregation is requested, the client gets redirected to the URI of a resource map. The resource map contains both information about a resource and information about the aggregation – which roughly represents a metadata record in EDM. This implementation also follows practical considerations from the point of view of application developers, as, more often than not, the data about a resource and the data about

the aggregation are used together. So this leads to a substantial reduction of necessary requests to the API.

*Versioning.* All DM2E data is versioned, i.e., the data provided under the URI of a resource map never changes. When updated data is ingested, the API redirects to the new resource map, but the new resource map gets a new URI and contains links to earlier versions of the data in the form of `prov:wasRevisionOf`. This allows the stable identification of triples within the data, a prerequisite for the data to become a trusted subject of scholarly work.

*Dataset provenance.* The full provenance of the DM2E data is provided by linking resource maps and other data pages to superordinate datasets using the VoID vocabulary [1]. The datasets are versioned and all data in a dataset shares the same provenance, following the idea of a common provenance context to support provenance-aware Linked Data applications [7]. The version of a resource map and the provided provenance information then simply corresponds to the version and provenance of the dataset. Versioned datasets are implemented as Named Graphs.

*Statement-level provenance.* To support contextualized resources with statements from various enrichment processes, a special approach has been implemented using statement annotations [7]. Subject URIs are created for all statements and these URIs are linked to the datasets the statements originate from. The statement URIs are identified and described as statements using RDF reification. All reification triples are created on the fly, only where necessary, and can safely be ignored by applications not interested in the provenance of the statements. The HTML representation of the contextualized resources makes use of this information and provides an “Oh Yeah?” button for all – possibly wrong – links to external resources, leading to the provenance information of the statement.<sup>37</sup> To the best of our knowledge this is the first implementation of this button as envisioned by Tim Berners-Lee [3].

## 7. Discussion

Several aspects make the DM2E dataset an interesting and unique source of information. First of all –

<sup>32</sup><https://github.com/DM2E/dm2e-ontologies/blob/master/src/main/resources/edm/edm.owl> (10.12.2015)

<sup>33</sup><https://github.com/DM2E/edm-validation>

<sup>34</sup><https://github.com/DM2E/dm2e-ontologies/blob/master/src/main/bash/dm2e-data.sh> (10.12.2015)

<sup>35</sup><http://wifo5-03.informatik.uni-mannheim.de/pubby/> (10.12.2015)

<sup>36</sup><https://github.com/dm2e/pubby> (10.12.2015)

<sup>37</sup>For example the city Nancy: <http://data.dm2e.eu/data/html/place/onb/abo/Nancy> (10.12.2015)



following the goals of the DM2E project – it contains data from many, carefully selected collections of not only manuscripts, but also old prints, posters, books and old journals with historic value. The data model was developed specifically for this domain where no suitable comprehensive data models existed yet. The DM2E model is also an example of an application profile, an application-specific specialization of the EDM. As such, the data blends well with the huge amount of EDM data available through Europeana. In contrast to many other Linked Datasets, the model and the API have both been tailored to the original data as well as to consuming applications. From a technical point of view, the use of multiple resource representations, versioning and the provision of a full provenance chain have to be mentioned, particularly the proper separation of original, curated metadata from data enrichments generated by automated processes with varying quality. The main short-coming is arguably the lack of a publicly available SPARQL endpoint, mainly due to performance considerations. Fast response times for the scholarly research platform have higher priority. The SOLR-based search and browse interface, however, is provided as convenient entry point to the data and provides a RESTful search API sufficient for most use cases. The data itself also has some shortcomings due to the heterogeneity of the original data. The quality of the metadata ranges from rich descriptions with unambiguous identifiers from authority files for agents, places and subjects to sparse descriptions with few information hidden in free-text fields. It is insofar a dilemma that the contextualization works best for the better data and particularly the poor data is hard to improve. A remedy might be the feedback of data from the annotations provided by the scholars. We plan to investigate this as part of our future work, when more annotations will hopefully be available.

## References

- [1] Keith Alexander, Richard Cyganiak, Michael Hausenblas, and Jun Zhao. Describing linked datasets – on the design and usage of void, the "vocabulary of interlinked datasets". In Christian Bizer, Tom Heath, Tim Berners-Lee, and Kingsley Idehen, editors, *Proceedings of the WWW2009 Workshop on Linked Data on the Web, LDOW 2009, Madrid, Spain, April 20, 2009.*, volume 538 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2009. Available at [http://ceur-ws.org/Vol-538/ldow2009\\_paper20.pdf](http://ceur-ws.org/Vol-538/ldow2009_paper20.pdf).
- [2] Konstantin Baierer, Evelyn Dröge, Vivien Petras, and Violeta Trkulja. Linked data mapping cultures: An evaluation of metadata usage and distribution in a linked data environment. In William E. Moen and Amy Rushing, editors, *Proceedings of the 2014 International Conference on Dublin Core and Metadata Applications, DC 2014, Austin, Texas, USA, October 8-11, 2014*, pages 1–11. Dublin Core Metadata Initiative, 2014. Available at <http://dcpapers.dublincore.org/pubs/article/view/3699>.
- [3] Tim Berners-Lee. Linked Data. <http://www.w3.org/DesignIssues/LinkedData.html>, 2006.
- [4] Kristin Dill, Evelyn Dröge, Øyvind Liland Gjesdal, Doron Goldfarb, Esther Guggenheim, Julia Iwanowa, Marko Knepper, Gerhard Müller, Alois Pichler, Kilian Schmidtner, Klaus Thoden, and Jorge Urzúa. D1.2 – final integration report. [http://dm2e.eu/files/D1.2\\_2.0\\_Final\\_Integration\\_Report\\_140214\\_final.pdf](http://dm2e.eu/files/D1.2_2.0_Final_Integration_Report_140214_final.pdf), 2014.
- [5] Evelyn Dröge, Thomas Bosch, Valentine Charles, Robina Clayphan, Mark Matienzo, Stefanie Rühle, Adrian Pohl, Miika Alonen, Lars Svensson, and Karen Coyle. Report on the current state: use cases and validation requirements [editor's draft]. Deliverable 1, DCMI RDF Application Profiles Task Force, 2014. Available at [http://wiki.dublincore.org/index.php/Deliverable\\_1](http://wiki.dublincore.org/index.php/Deliverable_1).
- [6] Evelyn Dröge, Julia Iwanowa, and Steffen Hennicke. A specialisation of the Europeana Data Model for the representation of manuscripts: The DM2E model. In *Libraries in the Digital Age (LIDA) Proceedings*, volume 13, 2014. Available at <http://ozk.unizd.hr/proceedings/index.php/lida/article/view/117>.
- [7] Kai Eckert. Provenance and annotations for Linked Data. In Muriel Foulonneau and Kai Eckert, editors, *Proceedings of the 2013 International Conference on Dublin Core and Metadata Applications, DC 2013, Lisbon, Portugal, September 2-6, 2013*, pages 9–18. Dublin Core Metadata Initiative, 2013. Available at <http://dcpapers.dublincore.org/pubs/article/view/3669>.
- [8] Kai Eckert, Dominique Ritze, Konstantin Baierer, and Christian Bizer. RESTful open workflows for data provenance and reuse. In Chin-Wan Chung, Andrei Z. Broder, Kyuseok Shim, and Torsten Suel, editors, *23rd International World Wide Web Conference, WWW '14, Seoul, Republic of Korea, April 7-11, 2014, Companion Volume*, pages 259–260. ACM, 2014. doi:10.1145/2567948.2577347.
- [9] Marco Grassi, Christian Morbidoni, Michele Nucci, Simone Fonda, and Giovanni Ledda. Pundit: Semantically structured annotations for web contents and digital libraries. In Annett Mitschick, Fernando Loizides, Livia Predoiu, Andreas Nürnberger, and Seamus Ross, editors, *Proceedings of the 2nd International Workshop on Semantic Digital Archives, Paphos, Cyprus, September 27, 2012*, volume 912 of *CEUR Workshop Proceedings*, pages 49–60. CEUR-WS.org, 2012. Available at <http://ceur-ws.org/Vol-912/paper4.pdf>.
- [10] Marco Grassi, Christian Morbidoni, Michele Nucci, Simone Fonda, and Francesco Piazza. Pundit: Augmenting web contents with semantics. *Literary and Linguistic Computing*, 28(4):640–659, 2013. doi:10.1093/lilc/fqt060.
- [11] Antoine Isaac and Bernhard Haslhofer. Europeana Linked Open Data - data.europeana.eu. *Semantic Web*, 4(3):291–297, 2013. doi:10.3233/SW-120092.
- [12] Krzysztof Janowicz, Pascal Hitzler, Benjamin Adams, Dave Kolas, and Charles Vardeman. Five stars of Linked Data vocabulary use. *Semantic Web*, 5(3):173–176, 2014. doi:10.3233/SW-140135.

- [13] Christian Morbidoni. Linked data and facets to explore text corpora in the humanities: a case study. In Matthew Horridge, Marco Rospocher, and Jacco van Ossenbruggen, editors, *Proceedings of the ISWC 2014 Posters & Demonstrations Track a track within the 13th International Semantic Web Conference, ISWC 2014, Riva del Garda, Italy, October 21, 2014.*, volume 1272 of *CEUR Workshop Proceedings*, pages 413–416. CEUR-WS.org, 2014. Available at [http://ceur-ws.org/Vol-1272/paper\\_125.pdf](http://ceur-ws.org/Vol-1272/paper_125.pdf).
- [14] Christian Morbidoni and Alessio Piccioli. Curating a document collection via crowdsourcing with Pundit 2.0. In Fabien Gandon, Christophe Guéret, Serena Villata, John G. Breslin, Catherine Faron-Zucker, and Antoine Zimmermann, editors, *The Semantic Web: ESWC 2015 Satellite Events - ESWC 2015 Satellite Events Portorož, Slovenia, May 31 - June 4, 2015, Revised Selected Papers*, volume 9341 of *Lecture Notes in Computer Science*, pages 102–106. Springer, 2015. doi:10.1007/978-3-319-25639-9\_20.

