

# X3ML mapping framework for information integration in cultural heritage and beyond

Yannis Marketakis<sup>1</sup> · Nikos Minadakis<sup>1</sup> · Haridimos Kondylakis<sup>1</sup> ·  
Konstantina Konsolaki<sup>1</sup> · Georgios Samaritakis<sup>1</sup> · Maria Theodoridou<sup>1</sup> ·  
Giorgos Flouris<sup>1</sup> · Martin Doerr<sup>1</sup>

Received: 28 December 2015 / Revised: 20 May 2016 / Accepted: 25 May 2016 / Published online: 6 June 2016  
© Springer-Verlag Berlin Heidelberg 2016

**Abstract** The aggregation of heterogeneous data from different institutions in cultural heritage and e-science has the potential to create rich data resources useful for a range of different purposes, from research to education and public interests. In this paper, we present the X3ML framework, a framework for information integration that handles effectively and efficiently the steps involved in schema mapping, uniform resource identifier (URI) definition and generation, data transformation, provision and aggregation. The framework is based on the *X3ML mapping definition language* for describing both schema mappings and URI generation policies and has a lot of advantages when compared with other relevant frameworks. We describe the architecture of the framework as well as details on the various available components. Usability aspects are discussed and performance metrics are demonstrated. The high impact of our work is

verified via the increasing number of international projects that adopt and use this framework.

**Keywords** Data mappings · Schema matching · Data aggregation · URI generation · Information integration

## 1 Introduction

Managing heterogeneous data is a challenge for cultural heritage institutions, such as archives, libraries and museums, but equally for research institutes of descriptive sciences such as earth sciences [1], biodiversity [2,3], clinical studies and e-health [4,5]. These institutions host and maintain various collections with heterogeneous material, usually stored in relational databases and often described by different metadata schemata. For exploiting this wealth of data, handling these metadata as a unified set is vital in many applications, including information retrieval, data integration [6], data migration [7] and creating rich domain-specific applications. In this direction, complex query and integration mechanisms have to be designed and implemented to enable uniform access to heterogeneous and autonomous data sources [8].

To allow data transformation and aggregation, it is required to produce mappings, to relate equivalent concepts or relationships from the source schemata to the aggregation schema, i.e., the target schema, in a way that facts described in terms of the source schema can automatically be translated into descriptions in terms of the target schema (also known as the “enterprise model” [9]). This is the mapping definition process and the output of this task is the mapping, i.e., a collection of mapping rules.

In this paper, we describe the *X3ML mapping framework* designed to support information integration for resource

---

✉ Yannis Marketakis  
marketak@ics.forth.gr

Nikos Minadakis  
minadakn@ics.forth.gr

Haridimos Kondylakis  
kondylak@ics.forth.gr

Konstantina Konsolaki  
konsolak@ics.forth.gr

Georgios Samaritakis  
samarita@ics.forth.gr

Maria Theodoridou  
maria@ics.forth.gr

Giorgos Flouris  
fgeo@ics.forth.gr

Martin Doerr  
martin@ics.forth.gr

<sup>1</sup> Institute of Computer Science, FORTH-ICS, Heraklion, Greece

discovery. In such a scenario, alternative knowledge violating cardinality constraints is equally relevant for resource discovery search. Thus, the current version of the framework does not take into account cardinality constraints or other rules in the target side. It guarantees only strict inheritance and subsumption and does not enforce any other constraints.

The framework is able to support the data aggregation process by providing mechanisms of data transformation and URI generation. Mappings are specified using the *X3ML mapping definition language*, a declarative, human readable language that supports the cognitive process of a mapping. Unlike XSLT that is intended to be comprehensible only by IT technicians, the *X3ML mapping definition language* can be understood by non-technical people as well. Thus, a domain expert is capable of both verifying the semantics and reading and validating the schema matching. This model carefully distinguishes between mapping activities carried out by the domain experts, who know and provide the data, and from the IT technicians, who actually implement data translation and integration solutions. *X3ML mapping framework* serves as an interface between them.

Usually, schema matching is used to describe the process of identifying that two different concepts are semantically related. This allows the definition of the appropriate mappings that are being used as input for the transformation process. However, a common problem is that the IT experts do not fully understand the semantics of the schema matching and the domain experts do not understand how to use the technical solutions. For this reason, in our approach the schema matching and the URI generation processes are separated. The schema matching can be fully performed by the domain expert and the URI generation by the IT expert, therefore solving the bottleneck that requires the IT expert to fully understand the mapping. Furthermore, this approach keeps the schema mappings between different systems harmonized since their definitions do not change, in contrast to the URIs that may change between different institutions and are independent of the semantics. Our approach completely separates the definition of the schema matching from the actual execution. This is important because different processes might have different life cycles; in particular, the schema matching definition has a different life cycle compared to the URI generation process. The former is subject to more sparse changes compared to the latter.

In this paper, we extend our previous work in the area [10] in many ways. First, we present more information and details on the background, the related work and the *X3ML mapping framework*. Although in our past work we focused on only one of the components (particularly, the *X3ML Engine* component), in this paper we present in detail the framework as a whole and describe all its relevant components; just indicatively the key components of the framework

are: (a) *Mapping Memory Manager*, (b) *3M Editor*, and (c) *X3ML Engine*; however, there are also additional components that are described in detail in Sect. 4. In addition, in this paper we demonstrate the feasibility of our solution out of various project experiences. More specifically, we describe how the proposed framework has been exploited for supporting the mapping and transformation of several archives and databases from various (even heterogeneous) domains to CIDOC CRM [11] and its extensions [12]. Finally, we discuss the usability aspects of some of these components and demonstrate our current advances in the area.

In the sequel, we use the following terminology: (a) we refer to the *X3ML mapping framework* as the X3ML framework and (b) to the *X3ML mapping definition language* simply as X3ML.

The remainder of this paper is organized as follows: Sect. 2 discusses the related work, whereas Sect. 3 presents the background for understanding the context of the X3ML. Section 4 presents the overall architecture of the X3ML framework providing more details on the individual components. Section 5 presents experiences on using the X3ML framework demonstrating the feasibility of our solution and the great advantages gained. Finally, Sect. 6 concludes this paper and discusses the future directions of our work.

## 2 Related work

In the latest years, an active field of research is focused on mapping relational databases (RDB) to RDF, since the majority of data currently published on the Web are still stored in relational databases with local schemata and local identifiers. Bridging the conceptual gap between the relational model and RDF is the key to make the data available as linked data, materializing the vision of the Semantic Web [13].

One approach toward this direction is the Direct Mapping [14, 15] which maps automatically relational tables to classes, and attributes of tables to properties using an RDF vocabulary. The URIs of the instances and the classes are automatically generated based on the RDB schema and data. One implementation that exploits Direct Mapping is SquirrelRDF [16]. This approach is based on mapping discovery, and URI generation is tied to the implementation of the system. This limits the ways to generate and use URIs, making therefore difficult the support of complex structures and information integration.

Besides Direct Mapping, there are also XML-based solutions such as D2R MAP [17] which is a declarative language for describing mappings between relational databases and OWL/RDFS ontologies. D2R MAP is extended by eD2R [18] to map databases that are not in first normal form and by D2RQ [19], which is formally defined by an RDFS schema.

The mappings are based on SQL queries that extract records from the RDB and transformation functions that can be applied to the extracted values.

Another XML-based approach, focusing on expressivity, is R2O [20] which is able to cope with complex mapping cases where one model is richer than the other. For that sake, Virtuoso RDF Views [21] are created, exposing RDBs, using a declarative meta schema language for defining the mapping of SQL data to RDF vocabularies. Triplify [22], on the other hand, maps HTTP-URI requests onto RDB queries and translates the resulting relations into RDF statements.

R2RML [23] is a mapping language proposed by W3C to standardize RDB to RDF mappings. There exist already several implementations [24] and R2RML is lately extended by RML [25] to simultaneously support different mapping sources into RDF. Similar to the previous solutions is the Relational.OWL [26] which is a representation format for mapping relational models to OWL full ontologies.

Besides approaches that try to map relational data, there are other approaches which try to map CSV files to RDF. For example, XLWrap's mapping language [27] provides conversions from CSV and spreadsheets to the RDF data model. Mapping Master's M2 [28] converts data from spreadsheets into OWL statements and Vertere [29] is a conversion tool based on a templating mechanism. Each line results in one or more RDF resources, while each column value can result in one or more triples about this resource. TaRQL [30] is used also for the same purpose. Other tools that provide mappings from XML to RDF lead to mappings in the syntactic level rather than on the semantic level and fail to provide solutions applicable to broader domains. There are tools in this category that are based on XSLT (Krextor [31], AstroGrid-D [32]), on XPATH (Tripliser [33]) and XQUERY (XSPARQL [34]) and tools that are based on algorithms (GRDDL [35]) that provide links between data and RDF. Finally, there are existing tools that provide mappings from several formats to RDF. Tools in this category include Datalift [36], the DataTank [37], OpenRefine [38], RDFizers [39], and Virtuoso Sponger [40]. A fundamental problem when using highly expressive languages such as XSLT is the fact that even the programmer himself has difficulties in understanding the mapping logic. The only way to verify the mapping is testing its output behavior.

Other approaches exploit mapping technologies to publish their data as linked data. For example, the Smithsonian American Art Museum<sup>1</sup> used KARMA [41] to publish their data as linked data, a tool trying to automate the mapping process and allowing the users to adjust the generated mappings. However, there is still no clear distinction on the work of the domain and the IT experts, which perplexes the whole workflow. KARMA uses R2RML model inheriting the issue

of tight coupling between the schema matching and the URI generation.

One work with similar goals to our approach is the SIP Creator [42], created for Europeana<sup>2</sup> in 2009 to bring together more than 150 different sources. Whereas it only dealt with flat formats at the beginning, it was later expanded to handle hierarchical formats as output, however only in XML. Motivated by the goal of transparency, the mapping file format created by the SIP Creator was seen to come closer to the structure that was required for generic mapping, since it appeared in an easy-to-comprehend, human readable XML format with the source and target paths clearly identified. However, it combined interactive schema matching resolving not formally structured elements (they call it "syntax normalization"), on the basis of the Groovy [43] programming language. Groovy however, as a programming language, cannot be used or understood by a domain expert.

Most of the described approaches tightly couple the URI generation and the schema matching processes and lack general conditions where mapping is dependent on particular data values, or data structures, limiting the ways that the URIs can be constructed and making it impossible to select other forms of URIs to be generated. The latter may frequently be required to adapt transformed data to different targets, such as Linked Open Data (LOD), or to look-up the identifiers used in the target system. Furthermore, in the described works, there is no distinction between IT and domain experts, so the IT experts are assumed to be responsible for the entire process. This is reflected in the complexity and lack of user friendliness of the employed mapping languages, even if some easy parts of the mapping are done graphically. Furthermore, IT experts cannot easily understand the domain semantics they are dealing with and URI generation is still based on a sort of unique key generation from data, as in relational database management systems (RDBMS).

All these different approaches prove that there is no standard model to support mapping of data sources other than relational, the technologies used are too complex to be used by the domain experts and the whole workflow is not well-defined. Compared to these approaches our work:

- (a) uses a simple model for defining the mappings in a way that is comprehensible and readable by the domain experts,
- (b) is generic because the mapping definitions are not tied to the implementation of the data transformation engine,
- (c) supports incremental changes of source and target schema,
- (d) supports customized URI generation policies, and
- (e) promotes the collaborative work of experts with different roles on the mapping process.

<sup>1</sup> <http://americanart.si.edu/>.

<sup>2</sup> <http://www.europeana.eu/>.

**Table 1** A categorization of the related works

	Model/language/specification	Software platform/ tool
Relational-based	Direct mapping [14] R2RML [23] RML [25] Relational.OWL [26]	Squirrel RDF [16] RDB2RDF [24] Virtuoso RDF Views [21] Triplify [22] D2RQ [19]
XML-based	D2R Map [17] eD2R [18] R2O [20] GRDDL [35]	Krextor [31] AstroGrid-D [32] Tripliser [33] XSPARQL [34]
CSV-based	XLWrap [27]	Master's M2 [28] Vertere [29] TaRQL [30]
Multiple		DataLift (DB, CSV, XML, GML, etc.) [36] The DataTank (CSV, XLS, XML, JSON, etc.) [37] OpenRefine (CSV, XML, JSON, RDF, etc.) [38] RDFizer (many) [39] Virtuoso Sparger (many) [40] KARMA (DB, CSV, JSON, XML, XLS) [41]

Table 1 shows the related approaches, categorized with respect to the different source types.

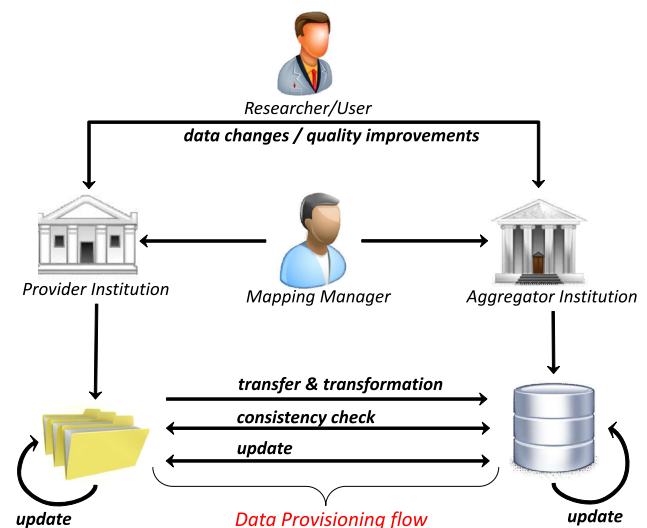
### 3 Background

Our work has been based on two main pillars: (a) the Synergy Reference Model that describes the best practices for the data provisioning and aggregation processes and (b) the *X3ML mapping definition language* (or simply X3ML), for describing the mappings per se. Below we will describe them in detail.

#### 3.1 Synergy Reference Model

The Synergy Reference Model (for short SRM) which is an initiative of the CIDOC CRM Special Interest Group,<sup>3</sup> is a reference model for a better practice of data provisioning and aggregation processes, primarily in the cultural heritage sector, but also for e-science. It is based on experience and evaluation of national and international information integration projects. It defines a consistent set of business processes, user roles, generic software components, and open interfaces that form a harmonious whole. Currently, a draft version of the model is available online,<sup>4</sup> still being evolved and enriched. The goal of SRM is to: (a) describe the provision of data between providers and aggregators including associated data mapping components, (b) address the lack of functionality in current models (i.e., OAIS [44]), (c) incorporate the necessary knowledge and input needed from providers to create quality sustainable aggregations, and (d) define a modular architecture that can be developed and optimized by different developers with minimal interdependencies and without hindering integrated UI development for the different user roles involved.

SRM aims at identifying, supporting or managing the processes needed to be executed or maintained between a

**Fig. 1** The data provisioning process

provider (the source) and an aggregator (the target) institution. It supports the management of data between source and target models and the delivery of transformed data at defined times, including updates. This includes a mapping definition, i.e., specification of the parameters for the data transformation process, such that complete sets of data records can automatically be transformed. A high-level view of the data provisioning process is shown in Fig. 1.

In more details, the main steps of the data provisioning workflow are:

- *Schema matching* source and target schema experts (a.k.a the domain experts) define a schema matching which is documented in a schema matching definition file. This file should be human and machine readable and it is the ultimate communication means on the semantic correctness of the mapping.
- *Instance generation specification* in this step, the URI generation and datatype conversion policies are defined for each instance of a target schema class referred to in the

<sup>3</sup> [http://www.cidoc-crm.org/who\\_we\\_are.html](http://www.cidoc-crm.org/who_we_are.html).

<sup>4</sup> [http://www.cidoc-crm.org/docs/SRM\\_v1.4.pdf](http://www.cidoc-crm.org/docs/SRM_v1.4.pdf).

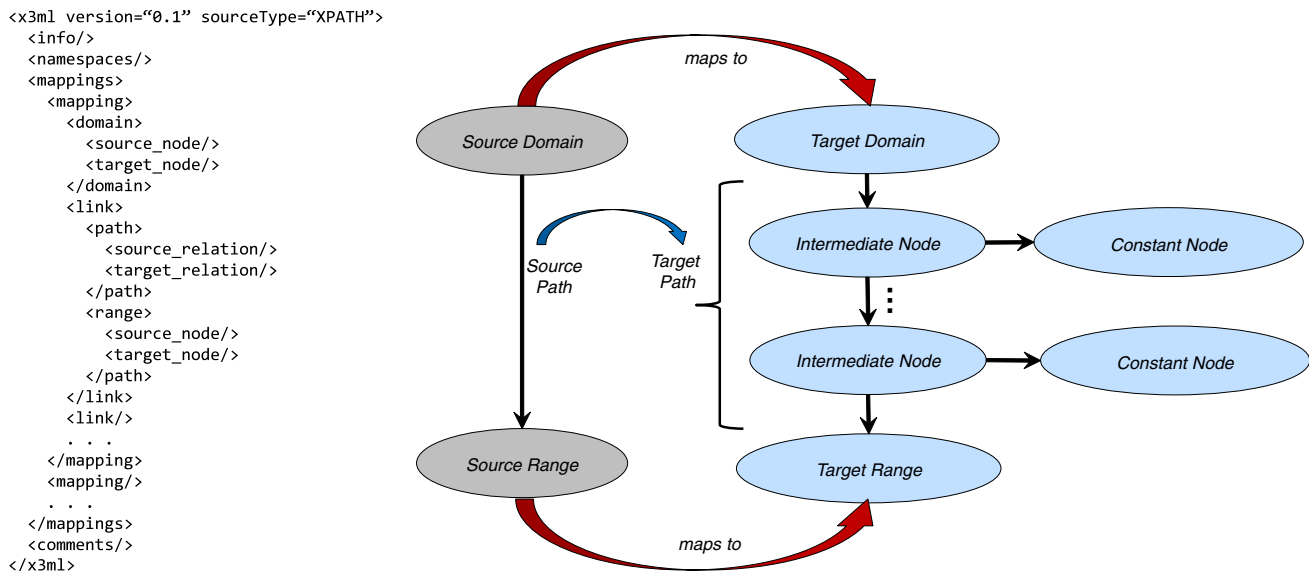


Fig. 2 The structure and the XML representation of an X3ML mapping

matching. In this step, only IT experts are involved and domain experts have no interest or knowledge about it.

- *Terminology mapping* the terminology mappings between source and target data/terms are defined. Providers may use anything from intuitive lists of uncontrolled terms up to highly structured third party thesauri.
- *Transformation* once the mapping definition has been finalized (and all syntax errors are resolved), the data needs to be transformed, producing a set of valid target records. The transformation process itself may run completely automatically. In the case where any issues arise, the aggregator can resolve them on a temporary or permanent basis, but it is also possible that these records are sent back to the provider for further analysis and resolution.
- *Ingestion* once records are transformed, an automated translation for source terms using a terminology map follows. The transformed records will then be ingested into the target system.
- *Change detection* after the ingestion of the records, all changes that may affect the consistency of the provider and aggregator data are monitored. SRM foresees a series of distinct update processes in all partner sites and is the only framework at the moment which takes the maintenance into account.

### 3.2 X3ML mapping definition language

X3ML is an XML-based language which describes schema mappings in such a way that they can be collaboratively created and discussed by experts. X3ML was designed on the basis of work that started in FORTH in 2006 [45] and emphasizes on establishing a standardized mapping description which lends itself to collaboration and the building of a

mapping memory to accumulate knowledge and experience. It was adapted primarily to be more aligned with the DRY principle [46] (avoiding repetition) and to be more explicit in its contract with the URI generation process. X3ML separates schema mapping from the concern of generating proper URIs, so that different expertise can be applied to these two very different responsibilities.

Schema matching is performed by domain experts who need to be concerned only with the correct interpretation of the source schema. The structure of X3ML is quite easy to understand consisting of: (a) a *header* that contains basic information (e.g., title, description, contact persons, the source and target schemata, sample records etc.) and (b) a series of *mappings*, each containing a domain (the main entity that is being mapped), and a number of links which consist of a path and a range. Each link describes the relation (path) of the domain entity to the corresponding range entity.

The basic mapping scheme and the XML representation of an X3ML mapping are shown in Fig. 2. Each *domain-path-range* of the source schema is mapped individually to the target schema and can be seen as a self-explanatory, context-independent proposition. As evident from the figure, X3ML mappings allow the definition of rich structures by adding intermediate nodes, even if the source relations are much simpler. An X3ML structure consists of:

- the mapping between the source domain and the target domain,
- the mapping between the source range and the target range,
- the proper source path,
- the proper target path,
- the mapping between source path and target path.



Below, we will describe the main concepts of the X3ML mapping definition language.

**Info and comment** Since X3ML is intended to bridge the gap between human authors and machines, it has to allow textual comments to be placed in the mapping specification. This is intended for describing alternatives to be discussed between experts and for harmonizing parallel mappings between providers. For this purpose, the info and the comment elements have been defined in the mapping specification. These blocks contain relevant information for humans to understand the specific mappings and can also contain useful provenance information recorded automatically by the tools used to manage the X3ML files, such as the date of creation or the author of the mapping file.

**Mapping** Each mapping element consists of a domain element and a number of links. It is quite common to have a single domain mapping and some range mappings, so by using this particular format for mappings, the single domain mapping does not have to be declared again. This is an ergonomic choice good for tree-dominated source schemata or sets of relational tables, which helps user orientation. It further provides an intuitive default local scope to define that the same instance of a class in a mapping rule appears as domain value in multiple target target propositions.

**Domain** The domain element is used to specify the mappings between a source (`source_node`) entity (table, class, non-leaf element) that can be regarded as domain of a source proposition and an equivalent target (`target_node`) entity. The `source_node` provides information on how to navigate to the source record and in case of XML it is an XPATH expression. The `target_node` defines an entity element that will lead to the generation of resource URIs or datatype values for the output graph. It may also contain `if` conditions (described below) upon which the mapping depends.

**Link** Inside the link element there is a path element. It allows mapping a source relation from the above defined source domain to a target relation to the above defined target domain. The path element must be followed by a range element, which is used to map the source and the target entities that are the equivalent range of the respective paths. The target relation might contain `if` conditions as well. A source/target range pair may reappear as a subsequent domain in an X3ML mapping.

**Conditional** The conditional expressions in *X3ML mapping definition language* can check for existence, equality and narrowness of values. They are expressed in the form of `if` statements and can be combined into Boolean expressions.

**Intermediate node** Sometimes, a path in the source schema needs to be further analyzed to a sequence of paths in the output with respect to the target schema. For this reason, the

user can define the generation of an intermediate node (or intermediate entity).

**Additional** Regularly constant properties and entities are needed to be added to a target entity, either from background knowledge or to characterize the meaning of a classification by the source schema rather than by data. For instance, a database about museum objects may not mention at all the museum as current keeper. A table “coin” may be mapped to “physical object”, but each instance of “coin” must have the type “coin”. For that purpose, an additional element can be used, containing the entity which will be attached to the target entity, the relationship describing the link, and the respective constant values.

**Variables** Sometimes it is necessary to generate an instance in X3ML only once in the scope of a given domain entity and then re-use it in a number of links of this domain. This is most frequently the case for intermediate target nodes. For example, a description of a museum object may reuse the same production event for mapping its “creator” link and its “date” link. In these cases, an entity can be assigned to a variable.

**Join operator** Sometimes, it is required to combine values from different tables in the source and produce new values in the output. This is the definition of the relational join operation. X3ML contains a specific operator for support (n-ary) join operation between different tables. The operator that is used is ‘==’ and is being used inside a link element. More specifically, it is being expressed inside the path element and expresses the equality of the value of the left-hand side (its table is the one defined in the domain of the corresponding mapping) with the value of the right-hand side (its table is the one defined on the range of the corresponding link).

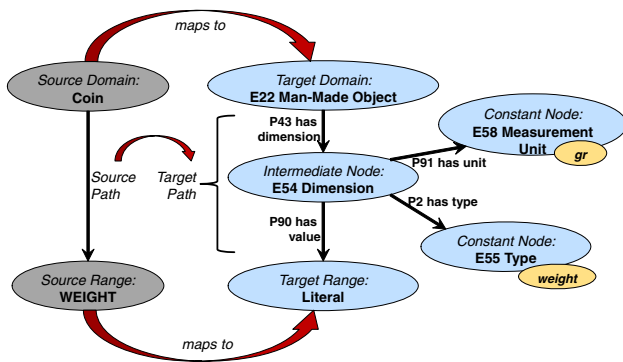
**Instance generation policy** The definition of the URI generation policy follows the schema matching and is performed usually by an IT expert who must ensure that the generated URIs match certain criteria such as consistency and uniqueness. A set of predefined URI generators (UUIDs, literals) and templates are available but any URI-generating function can be implemented and incorporated in the system. In the X3ML definition, the target domain and range contain the functions that generate URIs or literals.

Figure 3 shows how a simple relational database entry that specifies the weight of a coin is mapped and expressed with respect to the CIDOC CRM schema [11]. The mapping of this example can be found online.<sup>5</sup>

## 4 The X3ML mapping framework

The X3ML framework consists of a set of software components that assist the data provisioning process for infor-

<sup>5</sup> <http://www.ics.forth.gr/isl/OEAWcoins-Published>.



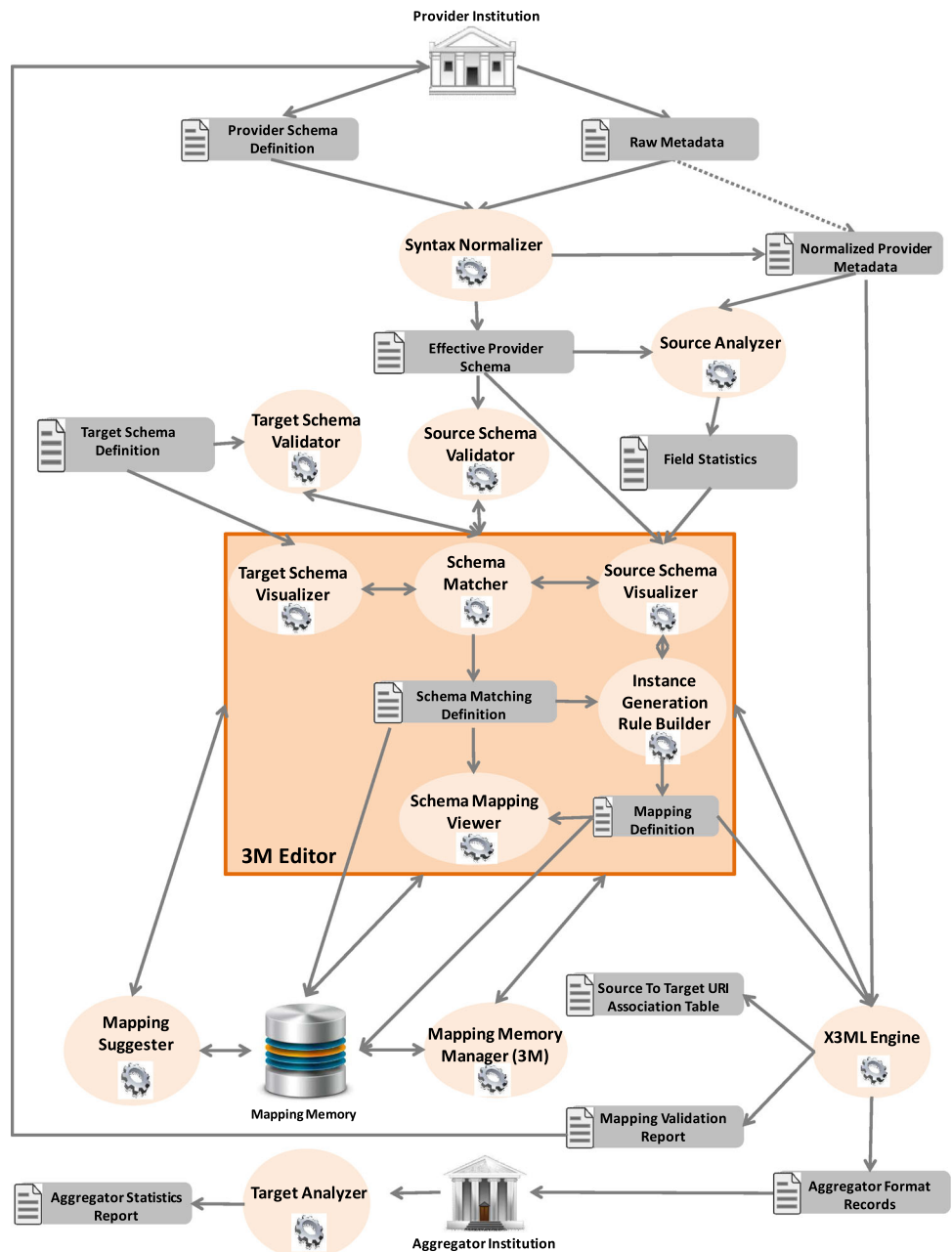
**Fig. 3** Mapping relational data to CIDOC CRM

information integration. A diagram presenting these components and the flow of data among them is shown in Fig. 4.

The process starts with the provider institution and is completed at the aggregator institution, when all records are transformed to the target format and are ingested to the target system.

Starting from the provider institution, the *syntax normalizer* can be used to normalize the provider's records. It exploits local syntax rules and produces a new provider schema definition, called *effective provider schema*. Normalization is quite often needed in date fields or in fields

**Fig. 4** The X3ML flow network



that contain concatenated information. For example, in the source schema the field:

```
<Acquisition>
  bequeathed; 1936-07-07; Young, Arthur W.
</Acquisition>
```

contains information about the actor, the date and the type of an acquisition and needs to be normalized before it is actually mapped to the target schema. In this example, the *effective provider schema* will be:

```
<Acquisition>
  <AcqType>bequeathed</AcqType>
  <AcqDate>1936-07-07</AcqDate>
  <AcqActor>Young, Arthur W.</AcqActor>
</Acquisition>
```

The next step of the provisioning process is the definition of the mappings. The X3ML framework provides various components that assist the experts to complete this time-consuming and error-prone task. One of these is the *Mapping Memory Manager* (for short *3M*). 3M is a managing system suitable for handling the mapping files. It offers a variety of actions that help both provider and target schema experts manage their files and also communicates with another component, called *3M Editor*. 3M Editor is an application suite that helps the experts complete the mapping definition process. The first step of the mapping definition is the schema matching. Provider schema experts with target schema experts exploit the *Schema Matcher* component, to define a schema matching which is documented in a *Schema matching definition*. The *Source Schema Validator* and the *Target Schema Validator* components assist the experts in selecting the valid paths with respect to the corresponding schemata, preventing them from making mistakes, while the *Source Schema Visualizer* and the *Target Schema Visualizer* assist users navigating through all source and target elements. The matching process is also supported by the *Mapping Suggester*, which makes use of “mapping memories” of similar cases as they are collected and cached from the user community.

The next step is the specification of the instance generation rules that define the URI generation policy for each target class. This task is accomplished by the *Instance Generator Rule Builder*, which complements the Schema Matching Definition with the instance generation policies, producing the *Mapping Definition* (which is also called X3ML file). Both the Schema Matching Definition and the Mapping Definition may be viewed with the *Schema Mapping Viewer*. The

files are stored to the *Mapping Memory*, which is an XML database.

Afterward, the mapping definition and the normalized input from the provider are given as input to the *X3ML Engine* component which is responsible for transforming the records to the aggregator format, so that they can be ingested in the aggregator institution.

The components of the X3ML framework have been designed with respect to the following design principles:

- *Collaborative mapping memory* The X3ML mapping descriptions must lend themselves to being stored and handled by collaborative tools, as well as potentially written by hand. This was the motivation for choosing a simple syntax in XML, and one which does not depend on implicit knowledge.
- *Facilitating instance matching* This involves extracting semantic information with the intent of generating correct instance URIs.
- *Transparency* The most important feature of X3ML framework is its generic application to mapping creation and execution and hopefully its longevity. The cleaner the core design of the framework and X3ML specification, and the clearer its documentation, the more readily will it get traction and become the basis for future mappings.
- *Reuse of standards and technologies* The best way to build a new software module is to carefully choose its dependencies and keeping them as small as possible. Building on top of proven technologies is the quickest way to a dependable result.

In the following subsections, we will describe these components in more detail.

#### 4.1 Syntax normalizer

The Syntax normalizer component is responsible for converting all data structures that are necessary for the transformation into a standard form. The reason for converting them is because data transformation components can deal with a limited set of standard data structures.

In many cases, the local syntax rules of the provider can be so complicated, or even non-deterministic, that it is often more effective to use a set of custom filtering rules that resolve one structural feature at a time and verify it with the source schema experts. For instance, if italics are used to tag a particular kind of field, it is better to convert first all italics to XML tags, so that the semantics of the field will be properly preserved and transformed.

To this end, the Syntax Normalizer component contains the necessary filters for converting the input data structures properly. Currently, three different syntax normalization fil-



ters have been implemented. The functionality and the design of the filters were based on actual requirements coming from the users of the X3ML framework.

The first filter normalizes the syntax of an XML file by rearranging its structure to fix syntax issues, so that it becomes compliant with the concepts of a semantic model. The second filter assists the users inspecting the values of specific elements in an XML file and update specific values while preserving the original ones. The third filter's main functionality is to split the values of an element, to create sub-elements and assign these values to them.

The following example shows an element *events* containing various dates, actors, and places:

```
<events>
  <date>2014-11-23</date>
  <date>2015-01-15</date>
  <actor>James</actor>
  <actor>Bob</actor>
  <place>Greece</place>
  <place>Italy</place>
</events>
```

In the above example, clearly the sequence of the elements is important; practically, there are two different events with the corresponding date, actor and place. Therefore, the Syntax Normalizer component re-constructs this block and produces the normalized one, which is shown below:

```
<events>
  <event>
    <date>2014-11-23</date>
    <actor>James</actor>
    <place>Greece</place>
  </event>
  <event>
    <date>2015-01-15</date>
    <actor>Bob</actor>
    <place>Italy</place>
  </event>
</events>
```

## 4.2 Mapping memory manager (3M)

3M is a tool for managing mapping definition files. It provides a number of administrative actions that assist the experts to manage their mapping definition files. It is available online<sup>6</sup> and is free to use. Users are requested to register online to obtain a username and password. Using these credentials,

users can login and see all the available mappings. Although all the mappings are visible to all users, each user is allowed to edit only his own mappings and the mappings he has access rights to do so. Users can also create a new mapping file, by giving a title and selecting one or more of the proposed target schemata. The newly created mapping file can be further edited with the 3M Editor (described in the subsequent section). Since the number of mappings can become quite big, several actions (such as searching, filtering and sorting) are provided. Moreover, users are able to export their mappings for off-line use, import mapping files, make versions of a specific mapping file, delete it, make a copy and also give rights to other users to edit their own file.

## 4.3 3M editor

The 3M Editor component is a Web application suite containing several software sub-components and exploits several external services. It is available online<sup>7</sup> and its main functionality is to assist users during the mapping definition process, using a human-friendly user interface and a set of sub-components that either suggest or validate the user input.

The user interface of the 3M editor is tabular and resembles modern Web browsers. The rationale for this design is to allow users to quickly familiarize with the editor and navigate instantly through the offered options. The main task of the editor is to support the creation of a complete X3ML file and check how the actual source data are mapped to the defined target output.

To create a complete X3ML file, the user has to fill in two different categories of resources:

- *Informational resources* they contain some general information about the mapping and its creators, the source, the target schema and the corresponding namespaces. Additionally, they might contain examples about the source or target records or even a generator policy file (we discuss them in detail Sect. 4.7).
- *Mapping resources* they contain information about the mappings themselves. Users are facilitated with the definition of the mappings using a set of supporting sub-components. A core sub-component is the Schema Matcher, which supports the definition of mappings. To define them, users are working with a tabular format of the mappings. We use this design since we found out that most users were accustomed to maintaining mappings of their own in spreadsheets. Thus, each mapping is represented as a table. Figure 5 shows an actual example. The header of the table represents the domain of the mapping, and the rows represent the links. Since each link contains two elements, one path and one range, the rows

<sup>6</sup> <http://www.ics.forth.gr/isl/3M/>.

<sup>7</sup> Through the 3M component.



(described in Sect. 4.2) when creating the mapping file. In addition, the user can also upload files that are going to be used as source or target schemata, example source or target records and a generator policy file. The uploaded files are then parsed, analyzed and used by the corresponding components.

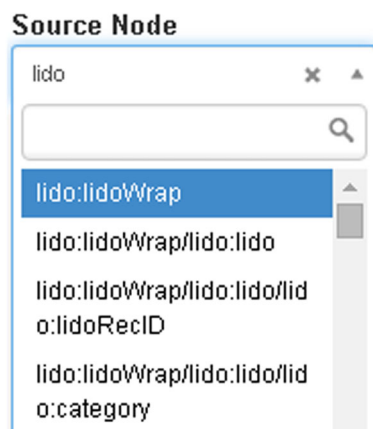
The component performs its functionality on an interactive manner; at any given moment, the user can switch from view mode (that uses the Schema Mapping Viewer component, described in Sect. 4.3.5) to edit mode, for inspecting how the mapping process is progressing. Furthermore, if the users are familiar with the *X3ML language* syntax, there is even a raw XML mode.

#### 4.3.2 Source schema visualizer

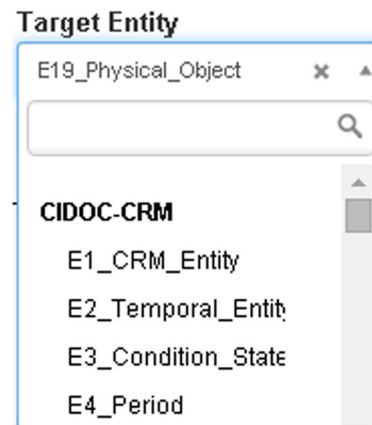
The Source Schema Visualizer component is responsible for assisting users in selecting the appropriate source path for the definition of the mapping. This component replaces the input fields with select boxes that contain values from the source records, so that users can select the values instead of writing them explicitly. The component analyzes the structure of the source schema and “proposes” values to the user, to avoid mistakes in writing the XPATH expressions. It is evident that the source schema file or the source records should be uploaded beforehand. The component can be activated or deactivated through the 3M Editor component. An indicative screenshot of the component is shown at Fig. 7. The source schema that is being exploited in this example is LIDO [47].

#### 4.3.3 Target schema visualizer

The Target Schema Visualizer is a software component for assisting users in selecting the appropriate target paths. When there are no target schema files available, users have to fill in the target paths using text input fields. If at least one target schema file is available, users can exploit this component



**Fig. 7** Selecting values from the source schema (LIDO) using Source Schema Visualizer



**Fig. 8** Selecting values from the target schema (CIDOC CRM) using the Target Schema Visualizer

for selecting the appropriate concepts (i.e., classes) from the target schema for their mappings.

The visualizer is activated after a target schema has been uploaded and validated. For the latter, the Target Schema Validator sub-component is exploited (described in more detail in Sect. 4.4). Figure 8 shows an indicative screenshot of the Target Schema Visualizer. The available options that are shown are classes of the target schema; in this case the target schema is CIDOC CRM.

#### 4.3.4 Instance generation rule builder

After the creation of the *Schema Matching Definition*, the specification for the instances generation is defined. This is accomplished using the Instance Generation Rule Builder software component. The user interface is similar to the Schema Matcher component (described in Sect. 4.3.1). However, users can only edit details about the instance and value generators.

The component uses templates for specifying how a URI or a label will be created and they can be exploited throughout the mapping. The generator templates are defined separately and linked to the actual mapping, to be used by the generators for producing the actual identifier or label. Each target entity must have only one instance generator and any number of label generators. Generators use their names as identifiers; therefore the same generator can be exploited for many different target entities. There are some generators already built in the X3ML framework; however, more user-defined generators can be defined. A more detailed discussion about generators can be found in the sequel (in Sect. 4.7).

Figure 9 shows an indicative example; the upper part defines how the URIs for instances of the class *E41\_Appellation* will be assigned, while the lower part defines how the corresponding labels will be created (e.g., *rdfs:label*).

**Fig. 9** Instance generation rule builder

#### 4.3.5 Schema mapping viewer

Schema Mapping Viewer is a sub-component used to inspect X3ML files. There are multiple interaction modes:

- *View mode* This mode shows the mappings using a human-understandable representation. More specifically, it uses a tabular format where each mapping is modeled as a different mapping, the header of the table represents the domain of the mapping, and the rows represent the links. This is the default mode for showing the mappings and is also used to present the publicly available mappings with respect to CIDOC CRM.<sup>8</sup>
- *Raw XML mode* X3ML is presented as raw XML. It may be the entire X3ML file or specific subtrees. This mode is intended to be used by users who are familiar with the syntax of the X3ML.
- *Graph view mode* X3ML is presented as a graph (shown in Fig. 6). This mode is intended for users non-experienced with the X3ML, or for quickly getting an overview of the defined mappings.

#### 4.4 Source schema and target schema validators

The Source Schema Validator is an external tool used to validate and traverse the source schemata and produce all possible XPATH expressions. The identified expressions are given as input to the Source Schema Visualizer for “proposing” them to the user. The supported source schemata may be in XSD or XML schema format.

The Target Schema Validator is a tool used to analyze the target schemata and to produce valid options according to previous user selections. The target schemata may be in RDFS, RDF, OWL, XSD or XML format. It consists of a set of engines for validating the schema based on their format.

Currently, the following Target Schema Validator engines are implemented:

- Using XML structure: It works with RDFS or RDF schema files. The engine is based on queries evaluated over the RDFS schemata stored in an XML database. Schemata have to be well-formed XML files containing certain tags (`rdfs:Class`, `rdf:Property`, `rdfs:domain` etc.). The results are sorted and grouped by schema file.
- Using semantic reasoner:<sup>9</sup> It works with RDFS, RDF or OWL schema files. The engine is based on Jena<sup>10</sup> and exploits also inference on the schema.
- Using XML analysis: It works with XSD or XML schema files. The engine produces XPATH expressions as a result.

#### 4.5 Source and target analyzers

The Source and Target Analyzers are two components under development that aim to provide information regarding the coverage of the source and target schemata, respectively.

The Source Analyzer contains three different metrics related to the source schema. The first metric displays the percentage coverage of parent elements, while the second counts only the leaves of the source schema that exist in the mapping. Finally, the third metric calculates the total coverage of the source schema.

The Target Analyzer allows users to explore the target schema, offering an effective schema analysis. The metrics are divided into two categories:

- Direct
- Ancestors/Descendants.

At the category of Direct metrics, each node is a separate element. Three different rates are provided, about classes, properties and resources. An element is considered to be covered only if it is mentioned in the mapping directly. On the other hand, the category of Ancestors/Descendants provides the same rates (classes, properties and resources), but it is differentiated on how the covered elements are calculated. An element is considered to be covered, if at least one ancestor or the element itself is referred in the mapping.

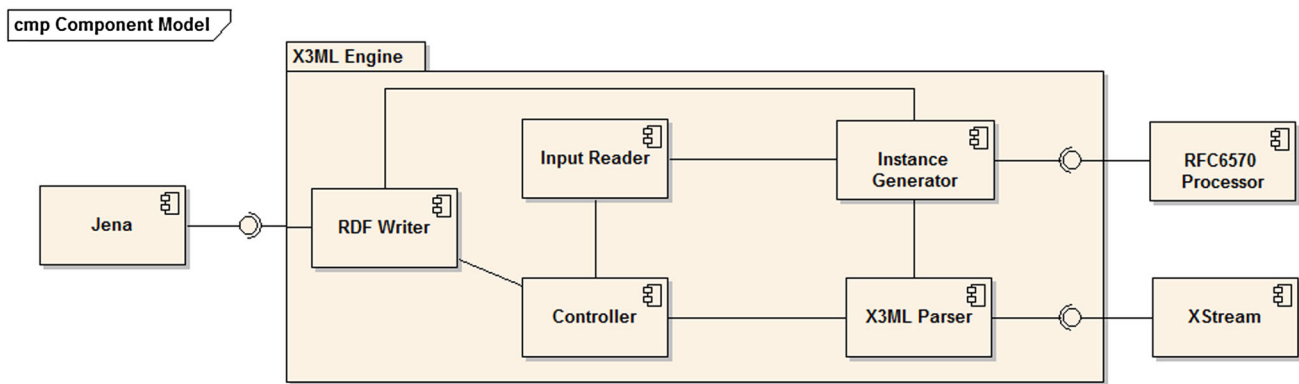
#### 4.6 Mapping suggester

The Mapping Suggester [48] is a software component under development used to suggest mappings to the user. These suggestions make use of “mapping memories” of similar

<sup>8</sup> <http://www.ics.forth.gr/isl/3M-PublishedMappings/>.

<sup>9</sup> <https://github.com/isl/Reasoner>.

<sup>10</sup> <https://jena.apache.org/>.



**Fig. 10** The main components of the X3ML engine

cases collected from the user community and are recalculated with each new mapping decision. The user can either accept or reject the suggestion. When the user creates a new mapping file, the Mapping Suggester runs a schema matching with the source schema they provide and the existing mapped source schemata in the mapping memory. The correspondences/crosswalks found during the schema matching are used by the mapping suggester to suggest mappings to the user.

#### 4.7 The X3ML engine

The X3ML Engine realizes the transformation of the source records to the target format. The engine takes as input the source data (currently in the form of an XML document), the description of the mappings in the X3ML mapping definition file and the URI generation policy file and is responsible for transforming the source document into a valid RDF document which corresponds to the input XML file, with respect to the given mappings and policy.

Figure 10 depicts the main subcomponents of the X3ML engine. The *Input Reader* component is responsible for reading the input data (currently, we support XML documents; however, as we describe in Sect. 6 more formats will be supported in the future using the proper extensions). The *X3ML Parser* component is responsible for reading and manipulating the X3ML mapping definitions. The component *RDF Writer* outputs the transformed data into the RDF format. The *Instance Generator* component produces the URIs and the labels based on the descriptions that exist in the mappings and finally the *Controller* component coordinates the entire process.

One of the most important tasks of the X3ML engine is the generation of values. Values can be either literal values or identifiers for resources (i.e., URIs). The generation of values is handled by the *InstanceGenerator* component. In general, this component supports the generation of: (a) instances and (b) labels. For each entity, there must exist one

*instance\_generator* and any number of subsequent *label\_generator* blocks in the corresponding mapping definition file.

The following block shows a template for defining an *instance\_generator*.

```

<instance_generator name='[gen-name]''>
  <arg name='[arg-name]'' type='[arg-type]''>
    [arg-value]
  </arg>
  ...
</instance_generator>

```

This type of generator is responsible for constructing an identifier for a resource. Each *instance\_generator* should be assigned a name and associated with a list of arguments. The arguments are exploited to provide the text segments that are required for constructing an identifier. Each argument has a name, which should be unique in the context of the generator, and a type and value pair. The type determines how the value of the argument will be used; there are three different options for the type:

1. *constant* In this case, the value of the argument should be used as it is defined. This type is used when we want to assign constant values to the generated identifiers.
2. *xpath* In this case, the value of the argument is an XPATH expression, that should be evaluated with the given input. The result of the XPATH expression will be used for the generation of the identifier. This type is used when it is required to exploit data from the input file in the generated identifiers.
3. *position* In this case, the value of the argument is ignored and the index position of the corresponding source node, within its context, is used.

The X3ML engine provides the default implementation for producing URIs and UUIDs.



The following block shows the configuration of a `label_generator`.

```
<label_generator name='' [gen-name] ''>
  <arg name='' [arg-name] '' type='' [arg-type] ''>
    [arg-value]
  </arg>
  <arg name='' language '' type='' constant ''>
    [language-code]
  </arg>
  ...
</label_generator>
```

Each `label_generator` has a name and a list of arguments (that are similar in spirit with the arguments that are defined for the `instance_generator`). One extra argument that is exploited is the argument for defining the language tag of the generated value. If it is empty, then it is implied that the generated value will not have it (i.e., in the case of number values). The X3ML engine provides default implementations for producing literal values in the form of `rdfs:label` and `skos:prefLabel` and constant values.

The *InstanceGenerator* component is configured through an XML file (which is given as input in the X3ML engine). When URIs are to be generated on the basis of source record content, it is wise to leverage existing standards and reuse the associated implementations. For the template-based URI generation, the RFC 6570 [49] standard is available and is used. Whenever the required URIs or labels cannot be generated by the default generators, the simple templates, or the URI templates, it is always possible to insert a special generator in the form of a class implementing the *InstanceGenerator* component interfaces.

The X3ML engine keeps an association table for the produced “values”. This practically means that whenever a new “value” is created (either a URI, UUID, or literal) the exact XPATH from the source file is associated with the generated value. The contents of this association table can also be exported in XML format and can be exploited from the user as a quick overview of the generated values for particular resources of the source.

#### 4.8 Implementation details and software licenses

The X3ML framework has been implemented in JAVA respecting the principles that were described in Sect. 3.2. The Web application components of the framework have been designed and implemented following a three-tier architecture [50]. More specifically, it comprises the following layers:

- *Data tier* Data are stored in an XML database. For the implementation of the Mapping Memory, we used eXist DB.<sup>11</sup>
- *Logic tier* Server code is implemented using J2EE. Whenever it was required, we reused existing and widely adopted libraries. An indicative list contains: XStream<sup>12</sup> and jOOX<sup>13</sup> for parsing XML-based documents, and Handy URI Templates<sup>14</sup> to support the generation of valid URIs.
- *Presentation tier*: Client code is implemented using Web standards such as HTML5,<sup>15</sup> CSS3,<sup>16</sup> Javascript and well-known libraries and frameworks such as jQuery,<sup>17</sup> AngularJS,<sup>18</sup> and Bootstrap.<sup>19</sup>

All the components of the X3ML framework have been developed as open source components. However, since they have been developed in the context of different EU projects, they are published under different open source licenses. More specifically, the 3M and 3M Editor components have been released under the European Union Public License,<sup>20</sup> whereas the X3ML engine has been released under the Apache 2 license.<sup>21</sup>

Furthermore, the source code of the X3ML framework has been deposited in a free and open repository. The source code and the complete documentation of the components are available from the following links:

- <https://github.com/isl/Mapping-Memory-Manager>
- <https://github.com/isl/3MEditor>
- <https://github.com/isl/x3ml>

## 5 Exploitation and evaluation

The X3ML framework has been validated, evaluated and widely used by several projects. Below, we first describe our experiences from using the X3ML framework in several EU projects (in Sect. 5.1) and then discuss about the evaluation results (in Sect. 5.2).

<sup>11</sup> <http://www.exist-db.org/exist/apps/homepage/index.html>.

<sup>12</sup> <http://x-stream.github.io/>.

<sup>13</sup> <https://github.com/jOOQ/jOOX>.

<sup>14</sup> <https://github.com/damnhandy/Handy-URI-Templates>.

<sup>15</sup> <http://www.w3.org/TR/html5/>.

<sup>16</sup> [http://www.w3schools.com/css/css3\\_intro.asp](http://www.w3schools.com/css/css3_intro.asp).

<sup>17</sup> <https://jquery.com/>.

<sup>18</sup> <https://angularjs.org/>.

<sup>19</sup> <http://getbootstrap.com/>.

<sup>20</sup> [https://joinup.ec.europa.eu/community/eupl/og\\_page/eupl](https://joinup.ec.europa.eu/community/eupl/og_page/eupl).

<sup>21</sup> <http://www.apache.org/licenses/LICENSE-2.0>.

### 5.1 X3ML exploitation in EU projects

The ResearchSpace project<sup>22</sup>—a collaboration of the British Museum, the Rijksmuseum, Oxford e-Research Centre, Yale Center for British Art (YCBA) and others—is developing a collaborative environment for humanities and cultural heritage research. The project has been using the X3ML framework for the mapping and transformation of their data using CIDOC-CRM as the target schema. A significant part of the British Museum and Rijksmuseum collection data have been transformed to CRM and integrated in ResearchSpace, showing the great value of the X3ML engine for transforming big amounts of real, heterogeneous cultural data.

Furthermore, the ResearchSpace project organized the CIDOC CRM Mapping workshop for humanities scholars and cultural heritage professionals, first at Yale University (USA) in August 2015, and then at Oxford University (UK) in October 2015. The workshop aimed to help professionals develop skills and understand data mapping techniques using CIDOC CRM. The participation exceeded the available positions (30 participants in Yale and 20 in Oxford), showing great interest in our approach and the potential of high impact of our work. Just indicatively the workshops attended well-known institutions in the area of cultural heritage domain, like GETTY museum,<sup>23</sup> Canadian Heritage,<sup>24</sup> Yale,<sup>25</sup> Frick<sup>26</sup> and others. The feedback we received from the two workshops was very positive<sup>27</sup> and it is very promising to see that, after attending the workshops, several institutions started mapping their data, an attempt continuing to date.

Another project using and exploiting our framework is the ARIADNE project [51]. The project initiated several mapping activities using the X3ML framework, to convert existing schemata of archaeological data to CIDOC-CRM and its extension suite. Using the X3ML framework, ARIADNE intended to promote an open sharing of data in the archeology sector, supporting effectively and efficiently all involved tasks. One of the tasks carried out in the context of ARIADNE project was to develop an integrated scenario where several Roman coin databases are being mapped and transformed to the common CIDOC CRM schema through X3ML and finally integrated in a semantic repository. So far, five archives have been integrated:

1. a set of 72 numismatic records from the dFMRÖ archive[52], an online MySQL database of the Numismatic Research Group of the Austrian Academy of Sciences,
2. a set of 1670 numismatic records coming from the Cambridge Fitzwilliam Museum archive,<sup>28</sup>
3. a set of 630 records coming from the Archaeological Superintendence of Rome<sup>29</sup> database,
4. a set of 517 coins of the Pergamon project from the Arachne object database<sup>30</sup> of the German Archaeological Institute,
5. the collections of MuseiDItalia, the digital library integrated in CulturaItalia with 25000 records already in CIDOC-CRM form.

The details of this integration use case are described in detail in [53]. This use case played an important role during the development of the X3ML framework both as a testbed and as a requirements provider. Furthermore, Vast-Lab<sup>31</sup> have been using the framework for mapping the Italian archaeological documentation system to CIDOC CRM [54].

The PARTHENOS project<sup>32</sup> aims at strengthening the cohesion of research in the broad sector of Linguistic Studies, Humanities, Cultural Heritage, History, Archaeology and related fields through a thematic cluster of European Research Infrastructures, integrating initiatives, e-infrastructures and other world-class infrastructures, and building bridges between different, though tightly interrelated, fields. PARTHENOS is committed to CRM-based integration and decided to use X3ML framework to support mappings, proving again its wide applicability and value. Towards this direction, X3ML framework has been integrated in the gCube infrastructure [55], and is publicly available.<sup>33</sup>

In addition, X3ML has been used for mappings of the Europeana Data Model[56] and LIDO[47] to CIDOC CRM, whereas VRE4EIC [57] will use the X3ML framework to support mappings to CERIF [58]. Moreover, the ITN-DCH project [59] has also initiated mapping activities using X3ML such as the mapping of the MayaArch3D database.<sup>34</sup>

X3ML has also been exploited in the biodiversity domain, in the context of the Lifewatch Greece project [60]. More specifically, it has been used for describing the mappings of data from the biodiversity domain that were derived from various relational databases, as well as from CSV files, and

<sup>22</sup> <http://www.researchspace.org/>.

<sup>23</sup> <http://www.getty.edu/museum/>.

<sup>24</sup> <http://www.pch.gc.ca/eng/1266037002102>.

<sup>25</sup> <http://www.yale.edu/>.

<sup>26</sup> <http://www.frick.org/>.

<sup>27</sup> Some feedback from the USA workshop as published in the call of the UK workshop can be found at <http://www.researchspace.org/home/project-updates/cidoccrmmappingworkshopoxforduniversity>.

<sup>28</sup> <http://www.fitzmuseum.cam.ac.uk>.

<sup>29</sup> <http://archeoroma.beniculturali.it/en>.

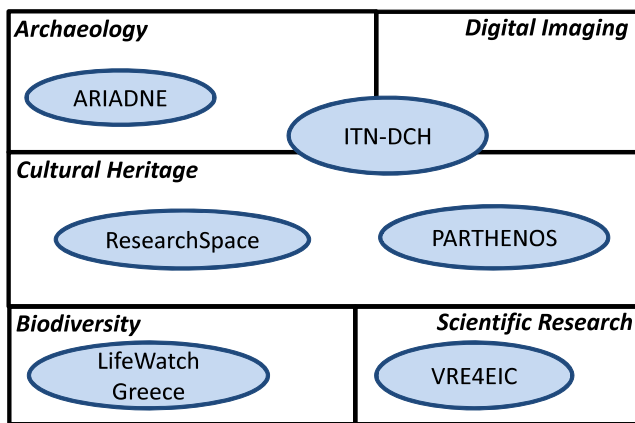
<sup>30</sup> <http://arachne.uni-koeln.de/drupal/>.

<sup>31</sup> <http://vast-lab.org>.

<sup>32</sup> <http://www.parthenos-project.eu/>.

<sup>33</sup> <http://mapping-d-parthenos.d4science.org/3M/>.

<sup>34</sup> <http://www.mayaarch3d.org>.



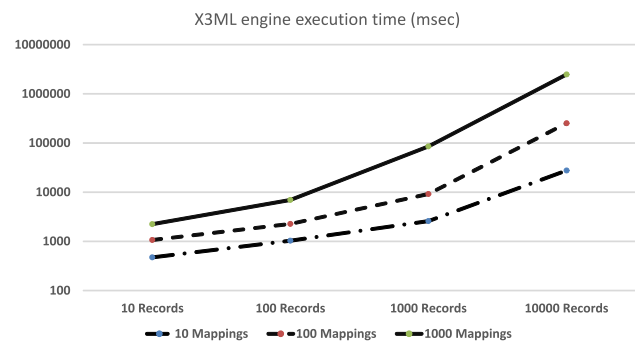
**Fig. 11** The current exploitation of the X3ML framework from various projects

were transformed into instances of the CIDOC CRM and its extensions for the biodiversity domain [2,3]. The aggregated records were then used for supporting the discovery of datasets and the scientific research processes related to the biodiversity domain.

All these diverse cases show an ever increasing interest in the X3ML framework and prove the effectiveness of our approach in real-world cases. Figure 11 illustrates the different domains (rectangles) and the various projects (circles) where X3ML framework has been exploited.

## 5.2 Performance evaluation

Apart from the evaluation over real institutions (serving real needs), we wanted to evaluate how efficient the framework is as the input grows. Therefore, we decided to evaluate the performance of the X3ML engine, since it is the component that requires the most computational resources. For the evaluation,<sup>35</sup> we used an XML input and a X3ML mapping example coming from the ARIADNE project as a base to produce synthetic data that was provided as input to the X3ML engine. Three X3ML mapping files were created containing 10, 100 and 1000 mappings and 4 XML input files containing 10, 100, 1000 and 10,000 records, respectively. Figure 12 depicts the evaluation results. We can observe that the overall time depends on both the number of mappings and the size of the input. More specifically, the time required for data transformation is approximately 1 s when the size of the input is low (10 records), even if the mappings are many (from 10 to 1000). As the size of the input increases, however, the overall time that is required increases as well. Note that the total number of output records is the total number of input records multiplied by the number of mappings (i.e., 10 input



**Fig. 12** X3ML engine evaluation results

records with 10 mappings will produce 100 output records). Concluding, we can see that the execution time is affected equally by the number of the mappings and the records, and it is related with the number of the links that are created during the transformation process.

## 6 Conclusions and future work

### 6.1 Discussion

This paper presented a novel framework for the management of the core processes needed to create, maintain and manage mapping relationships between different data sources. We described the X3ML framework and its building blocks for describing both schema mappings and URI generation policies, as well as tools for managing, editing, visualizing and executing those mappings. Our framework supports the mapping definition and transformation process and the generation of URIs and values, and is characterized by usability, scalability and ease of learning.

We demonstrated some of our experiences on using the aforementioned framework and discussed the evaluation of some of the related components. Finally, we demonstrated some of our experiences on using the aforementioned framework for mapping and transforming data from heterogeneous sources to CIDOC CRM and its extensions in various projects. To the best of our knowledge, the X3ML framework is the most complete mapping framework existing today with many tools supporting all different tasks involved.

### 6.2 Future work

There are many directions that we are currently exploring or plan to work on in the immediate future. First of all, as already discussed, most of the components are under continuous development and testing. More specifically, we are working on making the X3ML framework more generic and scalable.

<sup>35</sup> The experiments were carried out on a PC with an Intel i7 processor, 8GB RAM, running Windows 7 32 bit.

Moreover, an important future effort will be the consideration of alternative types of sources (at least for the X3ML Engine component). As already described in Sect. 4.7, the current version of the X3ML Engine requires the source records to be expressed as an XML document, but our plan is to extend it to support other types of input, such as, in particular, RDF input. This requires several modifications in the design and implementation of the engine. More importantly, the basic construct that we use for reading the source data will be an RDF model (e.g., Jena, Sesame), so instead of XPATH queries the SPARQL [61] language will be used for querying. Furthermore, the Instance Generator component should be able to carry the URIs from the source records to the target records if needed.

One apparent advantage of this approach is that the framework will support input and output of the same format. This sparked the light to investigate another direction: that of *invertible X3ML mappings*. In an invertible X3ML mapping, one can identify in a unique manner (and consequently regenerate) the data in the source dataset that led to the creation of each piece of data in the target dataset. Based on this idea, we can formalize the notion of *invertibility* by trying to identify how X3ML maps the source data to the target data.

Invertibility requires regenerating both the data and the URIs. For the latter, the association table (that was mentioned earlier in Sect. 4.7) can be used to support this functionality. The former is more difficult and consists in determining the triple(s) in the source dataset that contributed in the generation a given (set of) triple(s) in the target dataset. Note that it is not always possible to determine that in some cases, a triple in the target dataset can be (potentially) generated by two (or more) different sets of triples in the source dataset; therefore, the source cannot be uniquely determined. Defining necessary and sufficient conditions for invertibility, and identifying the inverse X3ML mapping (when one exists), is an ongoing work; a preliminary discussion appears in [10].

Solving the basic problem of invertibility will also allow addressing more complex problems that are interesting from the curator's point of view, such as:

1. Given a set of target records that have been derived from a set of source records using the appropriate X3ML mappings, and some updates that are applied over the target records, how can one propagate the necessary changes back to the source records?
2. Consider an integrated target dataset that has been aggregated from the application of multiple sets of X3ML mappings upon multiple source records. Suppose that the aggregation of information from the different records allows the inference of additional information, not present in (or inferrable by) any of the source records. How can this inferred information be propagated back

to the source records, and how can one determine the sources that will get each new information?

Finally as the number of mappings grow, it is becoming important to categorize them for assisting users finding easily the desired mappings within the 3M component. Apart from the search and filter functionalities that are already available, we plan to support also a categorization of mappings with respect to different thematic categories.

**Acknowledgements** This work was partially supported by the following projects: *ARIADNE* (FP7 Research Infrastructures, 2013–2017), *PARTHENOS* (H2020 Research Infrastructures, 2015–2019), *Blue-BRIDGE* (H2020 Research Infrastructures, 2015–2018), and *VRE4EIC* (H2020 Research Infrastructures, 2015–2018). The authors would also like to thank Nikos Anyfantis for working with the Source and Target Analyzer components and Korina Doerr for designing the user interfaces of the X3ML framework.

## References

1. Marketakis, Y., Tzitzikas, Y., Tona, C., Argenti, M., Marelli, F., Albani, M., Guarino, R., Polsinelli, B., Bitto, R.: On harmonizing earth science policies, semantics, metadata and ontologies. In: Ensuring the Long-Term Preservation and Value Adding to Scientific and Technical Data (PV'2013) (2013)
2. Tzitzikas, Y., Allocca, C., Bekiari, C., Marketakis, Y., Fafalios, P., Doerr, M., Minadakis, N., Patkos, T., Candela, L.: Integrating heterogeneous and distributed information about marine species through a top level ontology. In: Metadata and Semantics Research, pp. 289–301. Springer, New York (2013)
3. Tzitzikas, Y., Allocca, C., Bekiari, C., Marketakis, Y., Fafalios, P., Doerr, M., Minadakis, N., Patkos, T., Candella, L.: Unifying heterogeneous and distributed information about marine species through the top level ontology MarineTLO. Program Electron. Library Inf. Syst. **50**(1), 16 (2015)
4. Kondylakis, H., Plexousakis, D., Hrgovcic, V., Woitsch, R., Premm, M., Schüle, M.: Agents, models and semantic integration in support of personal ehealth knowledge spaces. In: Proceedings of Web Information Systems Engineering—WISE 2014—15th International Conference, Thessaloniki, Greece, October 12–14, 2014, Part I, pp. 496–511 (2014)
5. Kondylakis, H., Spanakis, E.G., Sfakianakis, S., Sakkalis, V., Tsiknakis, M., Marias, K., Zhao, X., Yu, H., Dong, F.: Digital patient: personalized and translational data management through the myhealthavatar EU project. In: 37th Annual International Conference of the IEEE Engineering in Medicine and Biology Society, EMBC 2015, Milan, Italy, August 25–29, 2015, pp. 1397–1400 (2015)
6. Kondylakis, H., Flouris, G., Plexousakis, D.: Ontology and schema evolution in data integration: review and assessment. In: On the Move to Meaningful Internet Systems: OTM 2009, Confederated International Conferences, CoopIS, DOA, IS, and ODBASE 2009, Vilamoura, Portugal, November 1–6, 2009, Proceedings, Part II, pp. 932–947 (2009)
7. Kondylakis, H., Plexousakis, D.: Ontology evolution: assisting query migration. In: Proceedings of Conceptual Modeling—31st International Conference ER 2012, Florence, Italy, October 15–18, 2012, pp. 331–344 (2012)
8. Kondylakis, H., Plexousakis, D.: Exelixis: evolving ontology-based data integration system. In: Proceedings of the ACM



- SIGMOD International Conference on Management of Data, SIGMOD 2011, Athens, Greece, June 12–16, 2011, pp. 1283–1286 (2011)
9. Calvanese, D., De Giacomo, G., Lenzerini, M., Nardi, D., Rosati, R.: Description logic framework for information integration. In: KR, pp. 2–13 (1998)
  10. Minadakis, N., Marketakis, Y., Kondylakis, H., Flouris, G., Theodoridou, M., Doerr, M., de Jong, G.: X3ML framework: an effective suite for supporting data mappings. In: Workshop for Extending, Mapping and Focusing the CRM—co-located with TPD'2015 (2015)
  11. Doerr, M.: The CIDOC conceptual reference module: an ontological approach to semantic interoperability of metadata. *AI Mag.* **24**(3), 75 (2003)
  12. Doerr, M.: CIDOC-CRM Family of Models. <http://www.ics.forth.gr/isl/CRMext>
  13. Berners-Lee, T., Hendler, J., Lassila, O., et al.: The semantic web. *Sci. Am.* **284**(5), 28–37 (2001)
  14. Berners-Lee, T.: Relational databases on the semantic web (2013)
  15. Hert, M., Reif, G., Gall, H.C.: A comparison of RDB-to-RDF mapping languages. In: Proceedings of the 7th International Conference on Semantic Systems, pp. 25–32. ACM (2011)
  16. SquirrelRDF. <http://jena.sourceforge.net/SquirrelRDF/>. Accessed Dec 2015
  17. Bizer, C.: D2R MAP—A Database to RDF Mapping Language. WWW (Posters) (2003)
  18. Barrasa, J., Corcho, O., Gómez-Pérez, A.: Fund finder: a case study of database-to-ontology mapping. In: Semantic Integration Workshop, p 9. Citeseer (2003)
  19. Bizer, C., Seaborne, A.: D2RQ—treating non-RDF databases as virtual RDF graphs. In: Proceedings of the 3rd International Semantic Web Conference (ISWC2004), vol. 2004. Citeseer, Hiroshima (2004)
  20. Barrasa Rodríguez, J., Corcho, O., Gómez-Pérez, A.: R2O, an extensible and semantically based database-to-ontology mapping language (2004)
  21. Openlink Software: Mapping Relational Data to RDF with Virtuoso's RDF Views. <http://virtuoso.openlinksw.com/whitepapers/relational%20rdf%20views%20mapping.html>. Accessed Dec 2015
  22. Auer, S., Dietzold, S., Lehmann, J., Hellmann, S., Aumüller, D.: Triplify: light-weight linked data publication from relational databases. In: Proceedings of the 18th International Conference on World Wide Web, pp. 621–630. ACM (2009)
  23. Das, S., Sundara, S., Cyganiak, R.: R2RML: RDB to RDF mapping language (2012)
  24. RDB2RDF Implementations. <http://www.w3.org/2001/sw/rdb2rdf/wiki/Implementations>. Accessed Dec 2015
  25. Dimou, A., Vander Sande, M., Colpaert, P., Verborgh, R., Mannens, E., Van de Walle, R.: RML: a generic language for integrated RDF mappings of heterogeneous data. In: Proceedings of the 7th Workshop on Linked Data on the Web (LDOW2014), Seoul, Korea (2014)
  26. de Laborda, C.P., Conrad, S.: Relational.OWL: a data and schema representation format based on OWL. In: Proceedings of the 2nd Asia-Pacific Conference on Conceptual Modelling, vol 43, pp. 89–96. Australian Computer Society, Inc. (2005)
  27. Langegger, A., Wöb, W.: XLWrap—querying and integrating arbitrary spreadsheets with SPARQL. Springer, New York (2009)
  28. O'Connor, M.J., Halaschek-Wiener, C., Musen, M.A.: Mapping master: a flexible approach for mapping spreadsheets to OWL. In: The Semantic Web—ISWC 2010, pp. 194–208. Springer, New York (2010)
  29. Vertere RDF. <https://github.com/knudmoeller/Vertere-RDF>. Accessed Dec 2015
  30. Tarql: SPARQL for Tables. <https://github.com/tarql/tarql>. Accessed Dec 2015
  31. Lange, C.: Krextor—an extensible framework for contributing content math to the web of data. In: Intelligent Computer Mathematics, pp. 304–306. Springer, New York (2011)
  32. AstroGrid-D: A Transformation from XML to RDF via XSLT. <http://www.gac-grid.de/project-products/Software/XML2RDF.html>. Accessed Dec 2015
  33. Tripliser. <http://daverog.github.io/tripliser/>. Accessed Dec 2015
  34. Bischof, S., Decker, S., Krennwallner, T., Lopes, N., Polleres, A.: Mapping between RDF and XML with XSPARQL. *J. Data Semant.* **1**(3), 147–185 (2012)
  35. Connolly, D. et al.: Gleaning resource descriptions from dialects of languages (GRDDL). W3C, W3C Recommendation, p. 11 (2007)
  36. Scharffe, F., Atemez, G., Troncy, R., Gandon, F., Villata, S., Bucher, B., Hamdi, F., Bihanic, L., Képékian, G., Cotton, F. et al.: Enabling linked data publication with the datalift platform. In: Proceedings of AAAI Workshop on Semantic Cities (2012)
  37. DataTank: Transform Datasets into a RESTful API. <http://thedataank.com/>. Accessed Dec 2015
  38. OpenRefine. <http://openrefine.org/>. Accessed Dec 2015
  39. RDFizers. <http://wiki.opensemanticframework.org/index.php/RDFizers>. Accessed Dec 2015
  40. Virtuoso Sponger. <http://virtuoso.openlinksw.com/dataspace/doc/dav/wiki/Main/VirtSponger>. Accessed Dec 2015
  41. Szekely, P., Knoblock, C.A., Yang, F., Zhu, X., Fink, E.E., Allen, R., Goodlander, G.: Connecting the Smithsonian American art museum to the linked data cloud. In: The Semantic Web: Semantics and Big Data, pp. 593–607. Springer, New York (2013)
  42. Sip Creator. <https://github.com/delving/delving/tree/master/sip-creator>. Accessed Dec 2015
  43. Groovy: A Multi-Faceted Language for the Java platform. <http://www.groovy-lang.org/>. Accessed Dec 2015
  44. Lavoie, B.: Meeting the challenges of digital preservation: the oais reference model. *OCLC Newsl.* **243**, 26–30 (2000)
  45. Kondylakis, H., Doerr, M., Plexousakis, D.: Mapping Language for Information Integration. Technical Report ICS-FORTH, vol. 385 (2006)
  46. Thomas, D., Hunt, A.: Orthogonality and the DRY Principle (2010)
  47. Coburn, E., Light, R., McKenna, G., Stein, R., Vitzthum, A.: LIDO—lightweight information describing objects version 1.0. In: ICOM International Committee of Museums (2010)
  48. Smyraki, O.: Design and Implementation of a Semi-Automatic Tool for Mapping Source Schemas to Target Ontologies. Master's thesis, University of Crete, Voutes Campus, 70013 Heraklion (2013)
  49. Gregorio, J., Fielding, R., Hadley, M., Nottingham, M., Orchard, D.: RFC 6570: URI template. Internet Engineering Task Force (IETF) Request for Comments (2012)
  50. Eckerson, W.W.: Three tier client/server architectures: achieving scalability, performance, and efficiency in client/server applications. *Open Inf. Syst.* **3**(20), 46–50 (1995)
  51. ARIADNE: Advanced Research Infrastructure for Archaeological Dataset Networking in Europe, FP7 Research Infrastructures, 2013–2017. <http://www.ariadne-infrastructure.eu/>
  52. dFMRÖ: Digitale Fundmünzen der Römischen Zeit in Österreich. <http://www.oew.ac.at/antike/index.php?id=358> (2007)
  53. Felicetti, A., Gerth, P., Meghini, C., Theodoridou, M.: Integrating heterogeneous coin datasets in the context of archaeological research. In: Workshop for Extending, Mapping and Focusing the CRM—co-located with TPD'2015 (2015)
  54. Felicetti, A., Scarselli, T., Mancinelli, M., Niccolucci, F.: Mapping ICCD archaeological data to CIDOC-CRM: the RA schema. In: A Mapping of CIDOC CRM Events to German Wordnet for Event Detection in Texts, p. 11 (2013)



55. Simeoni, F., Candela, L., Kakaletris, G., Sibeko, M., Pagano, P., Papanikos, G., Polydoras, P., Ioannidis, Y., Aarvaag, D., Crestani, F.: A Grid-Based Infrastructure for Distributed Retrieval. Springer, New York (2007)
56. Doerr, M., Gradmann, S., Hennicke, S., Isaac, A., Meghini, C., van de Sompel, H.: The europeana data model (EDM). In: World Library and Information Congress: 76th IFLA General Conference and Assembly, pp. 10–15 (2010)
57. VRE4EIC: A Europe-wide interoperable virtual research environment to empower multidisciplinary research communities and accelerate innovation and collaboration. H2020 Research Infrastructures 2015–2018. <http://www.vre4eic.eu/>
58. Asserson, A., Jeffery, K.G., Lopatenko, A.: CERIF: past, present and future: an overview. euroCRIS (2002)
59. ITN-DCH: Initial Training Network for Digital Cultural Heritage, 2013–2017. <http://www.itn-dch.eu/>
60. LifeWatch Greece: National Strategic Reference Framework, 2012–2015. <https://www.lifewatchgreece.eu/>
61. Prud'Hommeaux, E., Seaborne, A. et al.: Sparql Query Language for rdf. W3C Recommendation, p. 15 (2008)