

**SEVENTH FRAMEWORK PROGRAMME  
CAPACITIES**



**Research Infrastructures  
INFRA-2009-1 Research Infrastructures**

**OpenAIREplus**

**Grant Agreement 283595**

**“2nd-Generation Open Access Infrastructure for Research in  
Europe  
OpenAIREplus”**



**OpenAIRE Data Model Specification**

**Deliverable Code: D6.1**

## Document Description

### Project

Title:	OpenAIREplus, 2 <sup>nd</sup> Generation Open Access Infrastructure for Research in Europe
Start date:	1 <sup>st</sup> December 2011
Call/Instrument:	INFRA-2011-1.2.2
Grant Agreement:	<b>283595</b>

### Document

Deliverable number:	D6.1
Deliverable title:	OpenAIRE Data Model Specification
Contractual Date of Delivery:	30 <sup>th</sup> of April, 2012
Actual Date of Delivery:	21 <sup>st</sup> of April 2012
Editor(s):	Paolo Manghi
Author(s):	Paolo Manghi, Marko Mikulicic, Claudio Atzori
Reviewer(s):	Natalia Manola, Katerina Iatropulou, Antonis Lempesis, Jochen Schirrwagen, Mathias Loesch, Mateusz Kobos, Linda Reijnhoudt, Eko Indarto, Arjan Hogenaar, Wilko Steinhoff, Lars Nielsen
Participant(s):	Nikos Houssos, Keith Jeffery, Brigitte Joerg, Marko Mikulicic
Workpackage:	WP6
Workpackage title:	OpenAIREplus data model and content management services
Workpackage leader:	CNR
Workpackage participants:	NKUA, UNIBI, DTU/DataCite, CERN, UNIWARSAW, EKT-NHF/EuroCRIS, EBI-EMBL, KNAW-DNAS, STFC-BADC
Distribution:	Public
Nature:	Deliverable
Version/Revision:	1.0
Draft/Final:	Final
Total number of pages:	37

(including cover)

File name:

Key words: data model

## Disclaimer

This document contains description of the OpenAIREplus project findings, work and products. Certain parts of it might be under partner Intellectual Property Right (IPR) rules so, prior to using its content please contact the consortium head for approval.

In case you believe that this document harms in any way IPR held by you as a person or as a representative of an entity, please do notify us immediately.

The authors of this document have taken any available measure in order for its content to be accurate, consistent and lawful. However, neither the project consortium as a whole nor the individual partners that implicitly or explicitly participated in the creation and publication of this document hold any sort of responsibility that might occur as a result of using its content.

This publication has been produced with the assistance of the European Union. The content of this publication is the sole responsibility of the OPENAIRE consortium and can in no way be taken to reflect the views of the European Union.

The European Union is established in accordance with the Treaty on European Union (Maastricht). There are currently 27 Member States of the Union. It is based on the European Communities and the member states cooperation in the fields of Common Foreign and Security Policy and Justice and Home Affairs. The five main institutions of the European Union are the European Parliament, the Council of Ministers, the European Commission, the Court of Justice and the Court of Auditors. (<http://europa.eu.int/>)



OpenAIREplus is a project funded by the European Union

## Table of Contents

<b>Document Description</b>	<b>2</b>
<b>Disclaimer</b>	<b>4</b>
<b>Table of Contents</b>	<b>5</b>
<b>Table of Figures</b>	<b>6</b>
<b>Summary</b>	<b>7</b>
<b>Log of Changes</b>	<b>8</b>
<b>1 Scenario</b>	<b>9</b>
1.1 Data Model: Information Space Entities and Relationships	9
1.2 Extending Data Model to Support Data Inference	11
1.3 Outline	12
<b>2 Main data model entities</b>	<b>13</b>
2.1 CERIF Semantic Layer	13
2.2 Entities description	15
2.3 Entity-Relationship model	22
<b>3 Data population</b>	<b>28</b>
3.1 Information packages	29
3.2 Population Workflows	29
3.3 Identity of original entities	30
<b>4 Data inference</b>	<b>33</b>
4.1 Inference actions	34
<b>5 General data management issues</b>	<b>38</b>
5.1 Administrative entity properties	38

## Table of Figures

Figure 1 – Entity data flow .....	12
Figure 2 – E-R model: semantic layer entities .....	15
Figure 3 – E-R model: main, linked, static and structural entities.....	22
Figure 4 - E-R model: Result entities.....	23
Figure 5 – E-R model: provenance relationships .....	23
Figure 6 – Entity layers: original entities .....	28
Figure 7 – Information Packages: ingestion workflows (AM = Access Method, DS = Data source typology, PET = Primary Entity Type, F = information package Format): data sources of the same typology export the same primary entity of the same type through different “raw” information package format structures. ....	30
Figure 8 – Assigning unique identifiers to sub-entities: primary entity scope.....	32
Figure 9 – Assigning unique identifiers to sub-entities: data source scope.....	32
Figure 10 – Extension to data model to cope with inferred entities: explicit relationships become entities.....	36

## Summary

The OpenAIREplus web site will offer functionalities for administrators, anonymous and registered users to manage an Information Space of publications, together with their connections with funding projects (from the EC and national agencies) and research datasets. The aim of this document is to describe the conceived structure and semantics of this Information Space, i.e., the *OpenAIREplus data model*, by providing an abstract definition of its main entities and the relationships between them.

In this definitional process, the interaction with the EuroCRIS initiatives, several scientific institutes (i.e., KNAW-DANS, EBI-EMBL, BADC) as well as inspiration from DataCite and LinkedData play an important role in the specification of project data, i.e., how project data should be described, stored and exported in OpenAIRE, dataset metadata, how datasets should be described, and in the specification of how such interconnected entities can be made available and consumable by third-party systems.

The data model will be subject to changes in the future and therefore result in further versions. Such changes will be described in the following Section, in order to summarize to the reader the differences from the previous versions.

## Log of Changes

<b>Deliverable Version</b>	<b>Date</b>	<b>Changes</b> (description, section and pages)

## 1 Scenario

The OpenAIREplus web site will offer functionalities for administrators, anonymous and registered users to manage an Information Space of open access and non-open access publications, together with their connections to funding projects (from the EC and national agencies) and research datasets. In particular:

- *Anonymous users* will be able to search and consult the Space;
- *Registered users* will be able to:
  - Give feedbacks to improve the quality of the Space;
  - Claim publications or datasets into the Information Space;
- *Administrators (data curators)* will have full access and rights to such Space to:
  - Collect data from data sources;
  - Edit properties of publications, datasets, persons, etc in the Information Space;
  - Validate/invalidate insertions/deletions/updates suggested by registered users or by automatic inference processes.

### 1.1 Data Model: Information Space Entities and Relationships

In our reasoning we generalize the concept of datasets and publications to that of **project result**, so as to be able of including further kinds of research outputs. OpenAIRE initially proposes two kinds of results: **datasets** (e.g., experimental data, software products) and **publications**. But others can be added in the future (e.g., patents). Besides, project results are always associated to one or more *instances* of the results, in the sense that different “physical representations” of the same result may exist. For example, the same publication may be kept in two different repositories, both exposing the payload file (e.g., PDF) at different internet locations (URLs). Moreover, an instance of a result is represented as a combination of one or more *web resources* relative to the sub-parts of the result and of the internet data sources from which such resources are made available.<sup>1</sup>

Similarly, we extend the notion of authors of publications or datasets to that of **persons**, to include in the same set people connected to project fundings or organizations. For example “authorship” relationships between results and persons, which represent the fact that a given *person* has (co-)authored a given *result* while being affiliated with a given *organization*.

**Organizations** include companies, research centers or institutions involved as project partners or as responsible of operating data sources. Information about *organizations* will be initially collected from CORDA and CRIS systems, as being related to projects, or be ingested by users, for example to complete authorships information in the database.

Of crucial interest to OpenAIREplus is also the identification of the **funding programmes** which co-funded the **projects** that have led to a given result:

---

<sup>1</sup> The purpose of the project result-instance-web resource model is to capture a list of internet pointers relevant to the project result and not that of capturing the compound object structure which some results may have. If a project result is a compound object (e.g. ORE aggregation), its instances will likely be associated to web resources embodying its compound object nature (e.g., ORE resource maps).

- EC FP7 programme projects data will be fetched from the authoritative EC CORDA database, together with the *organizations* or *persons* which are *participants* of such projects.
- Data relative to National funding schemes and relative projects will be instead fetched from CRIS systems, together with other entities which may be typically kept within a CRIS system (e.g., publications, datasets, projects, organizations, people, etc).

Finally, OpenAIRE entity instances are created out of data collected from various **data sources** of different kinds, such as publication repositories, dataset archives, CRIS systems, etc. Data sources export **information packages** (e.g., XML records, HTTP responses, RDF data) which may contain information on one or more of such entities and possibly relationships between them. It is important, once each piece of information is extracted from such packages and inserted into the information space as an entity, for such pieces to be linked to the originating data source. This is to give visibility to the data source, but also to enable the reconstruction of very the same piece of information if problems arise. Initially, information relative to repositories will be collected by the OpenDOAR data source, which will act as main entity registry for (literature) repositories in Europe, but other data sources may join OpenAIRE in the future. The same centralized directory is instead not available for dataset archives and CRIS systems, whose managers/administrators will have to provide to OpenAIREplus while registering their data sources.

Entity instances, relative to persons, projects, organizations, results, and data sources will be instantiated from information packages collected, inferred, feed-backed, claimed from two main collection workflows: automated fetching from registered data sources and expert-provided information.

- *OpenAIREplus* compliant data sources: these include all data providers willing to authoritatively provide content to OpenAIRE,
  - *CRIS systems*: CERIF compatible data sources;
  - *Repositories*: institutional and thematic;
  - *Dataset Archives*: intended here as dataset providers from several subject areas;
  - *Data Source Aggregators*: intended here as systems federating several data sources of the same kind (e.g., repository federations);
  - *Entity Registries*: intended here as external data sources whose purpose is to provide a unique identity and reference for given entities (e.g., OpenDOAR for repositories, ORCID for persons)
- *OpenAIREplus expert-validated entity pool*: entities may reach the information space through human-driven workflows, which deliver entities into an expert-validated entity pool. Such a pool acts like a special data source, from which experts of various kinds (e.g., authors of results, project coordinators, data curators) can inject authoritative information into the OpenAIRE information space:
  - Authors (or others on their behalf) can “claim publications” into OpenAIREplus through the portal by:
    - Providing CrossREF DOIs and enriching them with project and license information;

- Searching publication metadata through external information spaces (e.g., BASE search engine) and selecting/enriching (project and license information) the ones they claim to be related to a project;
- Registered users can provide end-user feedbacks through the portal, to suggest data corrections or enrichments through “editing” actions to be validated by OpenAIRE data curators;
- EC project coordinators can confirm relationships between publications and EC projects automatically inferred by dedicated services;
- Data curators can validate guesses made by end-users through feedbacks or validate *inference actions* (see below) to make them persistent;
- Data curators can perform edit actions, such as entity addition, removal, or updates.

Such data sources remove, delete, update in the OpenAIRE information space information about one or more of the entities, as well as relationships between them. For example, some archives provide dataset metadata which also includes links to publications relevant for the dataset or vice versa. Such cross-entity and cross-sources data integration brings in data inference issues, which have mainly to do with information absence, duplication, and versioning (intended as replicas of the same entity). For example, many relationships (instances of) may not be available from data sources or may not be considered at all as valuable (e.g., not of interest to the specific research domain). Moreover, the same publication metadata may be collected from several resources, including repositories or CRIS systems. These issues will push into the data model a number of entities, properties, and relationships whose aim is to deliver to data curators the tools to maintain a clean, uniform, and consistent information space.

## 1.2 Extending Data Model to Support Data Inference

As planned in the DoW the OpenAIRE infrastructure will feature a number of services capable of curating the information space by disambiguating and enriching its entities, namely:

- *Duplicate inference*: different records representing the same entity (e.g., Result, Persons, Organizations, Projects) may be merged to disambiguate the information space.
- *Relationship inference*: new relationships between entities (e.g., citations, similarity semantics) may be inferred and added to the information space.
- *Attribute inference*: attribute values, such as titles or author names, will be inferred from the original full-text or digital files.

The OpenAIRE data model will therefore require to be extended in order to capture the entity information needed to cope with the distinction between **original entities**, which are either collected from data sources or provided by experts, and **inferred entities**, which can be instead re-calculated any time starting from the former set of entities and therefore have a “lower level of trust”. As shown in Figure 1, end-users and application access an information space which consists of the pool of original entities as collected from data sources and provided by experts, deduplicated and enriched by the data inference process through a layer of inferred information. We shall see that in some cases, inferred entities may become original entities, when an expert validates them and inserts them into

the expert-validated entity pool. In this case, such entities enter the information space as original entities and therefore become input to the data inference process.

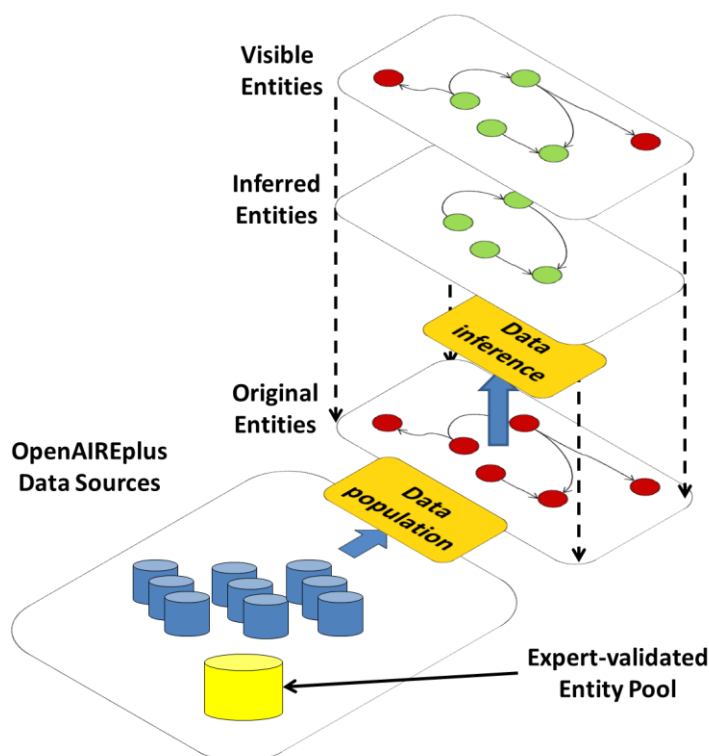


Figure 1 – Entity data flow

### 1.3 Outline

In the following, Section 2 provides a detailed description of the main entities that come into play by providing the relative Entity Relationship model. Section 3 describes the workflow of data population that is how data is collected from data sources and ingested into the information space according to the data model. Section 3 describes the issues encountered when introducing the concept of inferred entities into the data model and presents the entailed model changes. Finally, Section 5 describes the changes to be applied to the data model in order to cope with data management, from the an administrative perspective.

The data model will be subject to changes in the future and therefore result in further versions.

## 2 Main data model entities

This section provides a detailed description of the OpenAIREplus data model. Since the data model marries the notion of “semantic layer” as proposed by the CERIF model,<sup>1</sup> we shall first describe its abstraction mechanisms and then provide the details of OpenAIRE entities, their properties and their relationships.

### 2.1 CERIF Semantic Layer

According to this notion, (i) “horizontal” classification of entities (e.g., by vocabularies of terms) is not modeled through properties associated to given controlled vocabularies and (ii) semantic relationships between entities are not modeled by adding dedicated relationships. In both cases, CERIF introduces a flexible modeling mechanism which allows injecting classification semantics into “semantics-agnostic” entities and relationships. The mechanism is obtained by introducing two entities *Schemes* and *Classes* such that (see Figure 2):

**Class** A *Class* represents one term of a classification, e.g., vocabulary, taxonomy. As such it is characterized by the following properties: a *Code*, which represents the persistent identifier associated to the term (e.g., real-world classifications, such as ISO vocabularies for countries, have a standard identification code for terms), a *name*, an *acronym*, a *description*, a *StartDate*, and an *EndDate*. A *Class* is characterized by the following relationships with other entities:

- (i) *scheme*: the set of *Schemes* to which the *Class* belongs to (typically a *Class* describes a term which belongs to one *Scheme*, but there are cases where the same term can be shared across several *Schemes*);
- (ii) *related Classes* (inverse of relationships *Class1* and *Class2* from *Classes\_Classes* entities): the *Classes* related with the *Class* through the relationships entities *Classes\_Classes*; the semantics of these associations is specified in the *Classes\_Classes* entities (e.g., “partOf”, “parent”, “child”);

**Scheme** A *Scheme* identifies the existence of a classification scheme, which is modeled as a set of interrelated *Class* entities. A *Scheme* is characterized by the following properties: a *Code*, which represents the persistent identifier associated to the *Scheme* (e.g., real-world schemes, such as taxonomies, may have a standard identification code), a *name*, an *acronym*, a *description*, a *StartDate*, and an *EndDate*. A *Scheme* is characterized by the following relationships:

- (i) *related Classes* (inverse of the relationship *schemes* of *Class* entities): the *Classes* associated to the *Scheme*;
- (ii) *entryPoints*: the *Classes* at the first level of the *Scheme*.

**Class\_Class** Such entities represent associations between different terms (*Classes*), they are characterized by the following properties: a *StartDate* and an *EndDate*. They are characterized by the following relationships:

---

<sup>1</sup> CERIF data model: <http://www.eurocris.org/Index.php?page=CERIFreleases&t=1>

- (iii) *Class1* and *Class2*: which identify the two *Class* entities to be associated;
- (iv) *semanticsClass* and *semanticsScheme*: which respectively specify the semantics of this association through a given classification scheme.

**Scheme\_Scheme** Such entities represent associations between different Schemes, they are characterized by the following properties: a *StartDate* and an *EndDate*. They are characterized by the following relationships:

- (i) *Scheme1* and *Scheme2*: which identify the two *Scheme* entities to be associated;
- (ii) *semanticsClass* and *semanticsScheme*: which respectively specify the semantics of this association through a given classification scheme.

The mechanism allows the adoption of new classification schemes of arbitrary complexity: flat structures, such as vocabularies of terms for country (ISO 3166-1), tree structures, such as the EC FP7 funding scheme (<http://cordis.europa.eu/fp7>) or the WoRMS taxonomy of marine species (<http://www.marinespecies.org>), and graph structures such as the gene ontology (<http://www.geneontology.org>).

Table 1 – Properties and relationships for *Class*, *Scheme*, *Classes\_Classes*, and *Schemes\_Schemes* entities

<p><b>Class</b></p> <ul style="list-style-type: none"> <li>code</li> <li>name</li> <li>acronym (optional)</li> <li>description (optional)</li> <li>startDate</li> <li>endDate</li> </ul> <p>→ class1<sup>-1</sup> (0 or N Class_Class)</p> <p>→ class2<sup>-1</sup> (0 or N Class_Class)</p> <p>→ schemes (1 or N Scheme)</p> <p><i>Notes:</i> For simplicity, inverse relationships with other entities are not reported</p>	<p><b>Scheme</b></p> <ul style="list-style-type: none"> <li>code</li> <li>name</li> <li>acronym</li> <li>description</li> <li>startDate</li> <li>endDate</li> </ul> <p>→ entryPoints (0 or N Class)</p> <p>→ schemes<sup>-1</sup> (1 or N Class)</p> <p>→ scheme1<sup>-1</sup> (0 or N Scheme_Scheme)</p> <p>→ scheme2<sup>-1</sup> (0 or N Scheme_Scheme)</p> <p><i>Notes:</i> For simplicity, inverse relationships with other entities are not reported</p>
<p><b>Class_Class</b></p> <ul style="list-style-type: none"> <li>startDate</li> <li>endDate</li> </ul> <p>→ semanticsClass (1 Class)</p> <p>→ semanticsScheme (1 Scheme)</p> <p>→ class1 (1 Class)</p> <p>→ class2 (1 Class)</p>	<p><b>Scheme_Scheme</b></p> <ul style="list-style-type: none"> <li>startDate</li> <li>endDate</li> </ul> <p>→ semanticsClass (1 Class)</p> <p>→ semanticsScheme (1 Scheme)</p> <p>→ scheme1 (1 Scheme)</p> <p>→ scheme2 (1 Scheme)</p>

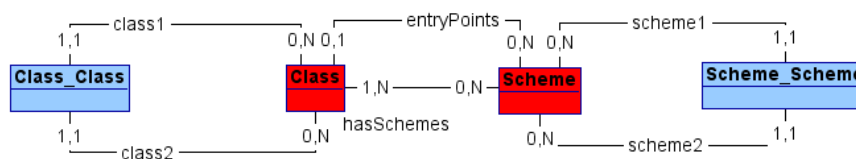


Figure 2 – E-R model: semantic layer entities

The OpenAIREplus data model introduces semantic-agnostic relationships between *publications-results*, *publications-publications*, *datasets-datasets*, *publications-datasets*, and *organizations-projects*. Their intended semantics will be injected thanks to a *Class* entity of a *Scheme* entity. Similarly, whenever entities need to be classified based in a property value (e.g., nationality of a person), property and values are modeled by an association to a *Class* (e.g., *nationalityClass*) and one to the relative *Scheme* (e.g., *nationalityScheme*). The benefit of the approach is that applications can be written in such a way they cope with the dynamic addition, removal, or deletion of *Classes* and *Schemes*. This is indeed the case in OpenAIREplus, where the intended entities and relationships will be subject to changes based on the results of Joint Research Activities.

## 2.2 Entities description

The entities in the data model can be grouped in the following way:

- *Main* entities: the entities whose information is continuously and incrementally fed to the information space; namely Result (Publication and Dataset), Person, Organization, DataSource (Repository, Dataset Archive, CRIS, Aggregator, Entity Registry), Projects;
- *Structural* entities: the entities added to the model to represent complex information about an entity; namely Instances, WebResources, Titles, Dates, Identities, and Subjects;
- *Static* entities: entities whose content is inserted in the information space at some point in time; namely Funding, Class, and Scheme;
- *Linked* entities (CERIF notation): relationship entities, used to connect in a semantic-agnostic way two or more main entities; namely, those denoted by an Entity1\_Entity2 notation.

### 2.2.1 Main Entities

Main entities are characterized by a provenance relationship *collectedFrom*, which indicates the DataSource entity (if it exists) from which entity information was collected. The conceptual representation of the schema is illustrated in Figure 3, where main entities, static entities and linked entities are represented, in Figure 4, where Result entities and their surroundings are represented, and in Figure 5, which groups all entities involved in the *collectedFrom* relationship.

**Result** A *Result* is here intended as the (metadata) description/representation of a scientific output (possibly) resulting out of one or more projects. A result is characterized by the following properties: a *year of acceptance*, a *publisher* (optional), a *description* (optional), and an *embargo end date* (empty if the *licenseClass* does not imply an embargo). A Result is characterized by the following relationships with other entities:

- Titles* (mandatory): the *Titles* of the Result, represented as *Class* entities of a Scheme entity, e.g., alternative, subtitle, etc.

- (ii) *Creators* (mandatory): the creators of the *Result*, which are *Person* entities connected to the result through relative *Person\_Result* entities;
- (iii) *Instances* (mandatory): the *Instances* of the *Result*, which represent the locations (*DataSource* entities) where the *Result* files (*web resources* entities, e.g., DOIs) can be found;
- (iv) *languageClass* and *languageScheme* (optional): the language used in the description or body of the *Result*, specified according to a given classification of languages, respectively described as a *Class* entity and a *Scheme* entity;
- (v) *collectedFrom* (mandatory): the *DataSource* entity from which the information relative to the *Result* entity was collected;
- (vi) *licenseClass* and *licenseScheme* (mandatory): the license of the *Result* according to a given classification of *Dataset* licenses, respectively described as a *Class* entity and a *Scheme* entity; the *Classes* should include "Unknown";
- (vii) *otherDates* (optional): a list of dates relevant to the *Result*;
- (viii) *fundingProjects* (optional): the *Projects* which co-funded the research underlying the *Result*;
- (ix) *subjects* (optional): the scientific disciplines (represented as *Class* entities of a *Scheme* entity) covered by the *Result*;

Note that the difference between *collectedFrom* and *hostedBy* is introduced to encode the peculiar notion of *Aggregation* data sources, whose entities are obtained by federating a set of *DataSource* entities. In OpenAIRE aggregators are *DataSources* from which entities are "collected", while the *DataSources* they aggregate "host" instead such entities. As such, Aggregators differ from the data sources they aggregate, but play an equally important role in delivering the entities to OpenAIRE, and should therefore be given visibility. In the case of other *DataSources*, e.g., repositories, *collectedFrom* and *hostedBy* refer to the same *DataSource*.

**Dataset** A *Dataset* is a *Result* further characterized by the following optional properties (ref. DataCite initiative v2.2): *resource type*, *size*, *format*, *version*, *last metadata update*, and *metadata version number*. A *Dataset* is characterized by the following relationships:

- (i) *the set of Datasets it is related with* (inverse of relationships *dataset1* and *dataset2* of *Dataset\_Dataset* entities): the semantics of such relationships is injected in *Dataset\_Dataset* entities through a *semanticsClass* and *semanticsScheme* relationships);
- (ii) *the set of publications it is related with* (inverse of relationships *dataset* of *Publication\_Dataset* entities): the semantics of such relationships is injected in *Publication\_Dataset* entities through a *semanticsClass* and *semanticsScheme* relationships);
- (iii) *resourceTypeClass* and *resourceTypeScheme*: the type of the *Dataset* according to a given classification of *Dataset* types, respectively described as a *Class* entity and a *Scheme* entity; The *resourceTypeScheme* "DataCite\_resource" defines the following *typeClass* values: *Collection*, *Dataset*, *Event*, *Film*, *Image*, *InteractiveResource*, *Model*, *PhysicalObject*, *Service*, *Software*, *Sound*, *Text*.

**Publication** A *Publication* is a *Result* further characterized by the following properties. A *Publication* is characterized by the following relationships:

- (i) *the set of publications it is related with* (inverse of relationships *publication1* and *publication2* of *publication\_publication* entities): the semantics of such relationships is injected in *publications\_publications* entities through a *semanticsClass* and *semanticsScheme* relationships);
- (ii) *the set of datasets it is related with* (inverse of relationships *Publication of Publication\_Dataset* entities): the semantics of such relationships is injected in *Publication\_Dataset* entities through a *semanticsClass* and *semanticssScheme* relationships)
- (iii) *typeClass* and *typeScheme*: the type of the *Publication* according to a given classification of publication types, respectively described as a *Class* entity and a *Scheme* entity.

**Person** A *Person* is characterized by the following properties: *firstName*, *secondNames*, *infixName*, *fax*, *email*, *phone*, *title*, and *gender*. A *Person* is also characterized by the following relationships:

- (i) *persistent identifiers* (optional): which is the list of unique and persistent identifiers used to identify the person together with the relative identification agency, e.g., ORCID;
- (ii) *creations* (optional): the creations of the *Persons*, which are *Results* entities connected to the result through relative *Person\_Result* entities;
- (iii) the participations of the *Person* to *Projects* (optional): inverse of *contactPerson* relationships of *Participants* entities;
- (iv) *nationalityClass* and *nationalitycheme*: the nationality of the *Persons* according to a given classification of nationalities, respectively described as a *Class* entity and a *Scheme* entity.
- (v) *collectedFrom*: the *DataSource* entity from which the information relative to the *Person* entity was collected;

**Project** A *Project* is characterized by the following properties: an *persistent identifier* (which is the unique and persistent identifier used to identify the project by its funding agency, e.g., grant agreement number for the EC), a *title*, an *acronym*, a *web site* (e.g., in the case of EC projects the project page at CORDIS), a *start\_date*, an *end\_date*, a *duration* (derived by start and end dates), a *project call identifier*, and a list of *keywords*. It also features a special flag *EC\_SC39*, indicating whether or not an EC project is subjected to clause 39. A *Project* is also characterized by the following relationships:

- (i) *participants*: which is the set of *Participants* participating to the *Project*;
- (ii) the set of *Results* whose research was co-funded by the *Project* (inverse of the relationship *fundingProjects* of the *Result* entities);
- (iii) *fundedBy*: the set of *Fundings* entities indicating which grants co-funded the *Project*;
- (iv) *collectedFrom*: the *DataSource* entity from which the information relative to the *Project* entity was collected.

**Organization** An *Organization* is characterized by the following properties: an *original identifier* (persistent identifier, if any, made available by the source providing organization information), a *legal short name*, a *legal name*, an *URL of its web site*, and an *URL of the logo* (if available). An organization is further characterized by a number of flags, which are always present in the case of *Organization* information provided by the EC: *legal body*, *legal person*, *non profit*, *research organization*, *higher education*, *international organization with Eur interests*, *international organization*, *enterprise*, *SME validated* flag, and *NUTS code* (Nomenclature of Territorial Units for Statistics). An *Organization* is characterized by the following relationships:

- (i) the set of *Persons* which have created a *Result* while affiliated with the *Organization* (the inverse of the relationship *creatorAffiliation* of *Person\_Result*);
- (ii) *dataSources*: the *DataSource* entities under the responsibility of the *Organization*;
- (iii) the participations of the *Organization* to *Projects* (inverse of *respOrganization* relationships of *Organization\_Project* entities);
- (iv) *countryClass* and *countryScheme*: the country of the *Organization* according to a given classification of countries, respectively described as a *Class* entity and a *Scheme* entity;
- (v) *collectedFrom*: the *DataSource* entity from which the information relative to the *Organization* entity was collected.

**DataSource** A *DataSource* is characterized by a *persistent identifier* (if available, as released by the agency or organization providing *DataSource* information, e.g., OpenDOAR for repositories), an *official name*, an *English name* (if available), a *URL of its web site*, an *URL of the logo*, a *description*, a *contact email* (if available), and an *access information package*, which contains API information useful to access the data source. It is characterized by the following relationships:

- (i) the set of *Result Instances* it is hosting (inverse of the relationship *hostedBy* of *Instance* entities);
- (ii) the set of *Result Instances* which were collected from it (inverse of the relationship *collectedFrom* of *Instance* entities)
- (iii) the set of *Organizations* responsible for the *DataSource* (inverse of the relationship *respOrganization* of the *Organization* entities)
- (iv) *collectedFrom*: the *DataSource* entity from which the information relative to the *DataSource* entity was collected.

**Repository** A *Repository* is a *DataSource* characterized by the further properties inherited from OpenDOAR description: *numberOfItems*, *numberOfItemsDate*, *subjects*, *policies*, *languages*, and *contentTypes*.

**CRIS system/Data Archive** Both entities are *DataSources*.

**Aggregator** An *Aggregator* is a *DataSource* characterized by the following relationships:

- (i) *typeClass* and *typeScheme*: the type of *Aggregator* according to a given classification of *Aggregator* types (i.e., the kind of data sources they aggregate, such as repositories, dataset archives, etc), respectively described as a *Class* entity and a *Scheme* entity.

**Entity Registry** An *Entity Registry* is a *DataSource* characterized by the following relationships:

- (i) *typeClass* and *typeScheme*: the type of *Entity Registry* according to a given classification of *Registry* types (i.e., the kind of entities they contain), respectively described as a *Class* entity and a *Scheme* entity.

### 2.2.2 Structural entities

**Title** A *Title* is a title of a *Result* characterized by a *Name* and a type provided by a *titleTypeScheme* and a *titleTypeClass* properties. The *titleTypeScheme* "DataCite\_title" defines the following *typeClass* values: Alternative Title, Subtitle, Translated Title. A *Title* is characterized by the relationships *titles*<sup>-1</sup>, which refers to the *Result* to which it belongs.

**Date** A *Date* is a date different from the date of "acceptance" for both publications and datasets. It is characterized by a property *Date* and a type provided by a *dateTypeScheme* and a *TypeClass* properties. The *typeScheme* "DataCite\_date" defines the following *dateTypeClass* values: Available, Copyrighted, Created, EndDate, Issued, Start Date, Submitted, Updated, Valid. A *Date* is characterized by the relationships *otherDates*<sup>-1</sup>, which refers to the *Result* to which it belongs.

**Identity** An *Identity* is a combination of a *unique identifier* used to uniquely refer to the individual together with the relative *identifier scheme* (issuer), e.g., ORCID, ISNI. An *Identity* is characterized by the relationships *instances*<sup>-1</sup>, which refers to the *Person* to which is assigned.

**Instance** An *Instance* represents the combination of the *Web Resources* (i.e., URLs relative to files or file locations, such as "splash pages") associated with a *Result* and the *DataSource* where such *Web Resources* are hosted. An *Instance* also contains the persistent identifier of the *Result*, if available. As such, an *Instance* is characterized by the following relationships:

- (i) the *Result* of which it is an *Instance* (the inverse of the relationship *instances* of *Result* entities);
- (ii) *webResources*: the set of *WebResources* associated to the *Instance*;

**Web Resource** A *WebResource* is characterized by *its unique URL* and by its relationship with the associated *Instance* (inverse of the relationship *webResources* of *Instance* entities).

**Subject** *Subjects* are scientific disciplines associated to one *Result*, together with a given semantics. As such they are characterized by the following relationships:

- (i) the *Result* associated to the subject (the inverse of the relationship *subject* of *Results* entities);
- (ii) *semanticsClass* and *semanticsScheme*: the semantics of the Subject according to a given classification of Subjects, respectively described as a *Class* entity and a *Scheme* entity.

### 2.2.3 Static entities

**Funding** A *Funding* entity represents funds that can be granted to *Projects*. As such it is characterized by the following properties: *name*, *description*, *keywords*, and a *persistent identifier* (if available). Moreover, by the following relationships:

- (iii) the *Projects* granted the funding (inverse of relationship of *Project* in *Projects\_Funding* entities);
- (iv) *semanticsClass* and *semanticsScheme*: the semantics of the Fundings (e.g., in FP7: Supporting and Coordination Actions, I3, STREP) according to a given classification of Fundings semantics (e.g., EC-FP7 contract schemes), respectively described as a *Class* entity and a *Scheme* entity.

#### 2.2.4 Linked entities

**Publication\_Publication** *Publication\_Publication* entities represent relationships between *Publication* entities, together with the duration and the semantics of the relationship. As such it is characterized by the properties *start date* and *end date* and by the following relationships:

- (i) *Publication1* and *Publication2*: the two interrelated *Publication* entities;
- (ii) *semanticsClass* and *semanticsScheme*: the semantics of the relationship according to a given classification of relationships between Publications, respectively described as a *Class* entity and a *Scheme* entity.

**Funding\_Funding** *Funding\_Funding* entities represent relationships between *Funding* entities, together with the duration and the semantics of the relationship. As such it is characterized by the properties *start date* and *end date* and by the following relationships:

- (i) *Funding1* and *Funding2*: the two interrelated *Publication* entities;
- (ii) *semanticsClass* and *semanticsScheme*: the semantics of the relationship according to a given classification of relationships between Fundings, respectively described as a *Class* entity and a *Scheme* entity.

**Project\_Funding** *Project\_Funding* entities represent relationships between *Projects* and *Fundings* entities, together with the duration and the semantics of the relationship. As such it is characterized by the properties *start date* and *end date* and by the following relationships:

- (i) *Project* and *Funding*: the two interrelated entities;
- (ii) *semanticsClass* and *semanticsScheme*: the semantics of the relationship according to a given classification of relationships between *Funding* and *Projects*, respectively described as a *Class* entity and a *Scheme* entity.

**Organizations\_Funding** *Organization\_Funding* entities represent relationships between *Organizations* and *Fundings* entities, together with the duration and the semantics of the relationship. For example, the EC is the organization behind FP7 funding (Class = funding organization). As such it is characterized by the properties *start date* and *end date* and by the following relationships:

- (i) *Organization* and *Funding*: the two interrelated entities;
- (ii) *semanticsClass* and *semanticsScheme*: the semantics of the relationship according to a given classification of relationships between *Funding* and *Organization*, respectively described as a *Class* entity and a *Scheme* entity.

**Organization\_Project** *Organizations\_Projects* entities represent relationships between *Publication* and *Organization* entities, together with the semantics of the relationship.

Examples of such semantics, as suggested by CERIF vocabularies, are: subcontractors, principal investigating, exploitation, coordinators, participant. Example is a formal beneficiary of fundings within a *Project*. It is characterized by an *Original Identifier* (e.g., for EC projects it is the participant number, a progressive integer 1,2,3,4, etc., where 1 is assigned to the project coordinator), a *start date* and an *end date*. A Participant is characterized by the following relationships:

- (i) the *Project* involved in the *Organization\_Project* entity;
- (ii) *respOrganization* involved in the *Organization\_Project* entity;
- (iii) *contactPerson*: the *Person* recorded as contact point for the responsible *Organization*;
- (iv) *semanticsClass* and *semanticsScheme*: the semantics of the relationship according to a given classification of relationships between Organizations and Projects, respectively described as a *Class* entity and a *Scheme* entity.

**Dataset\_Dataset** *Dataset\_Dataset* entities represent relationships between *Dataset* entities, together with the duration and the semantics of the relationship. As such it is characterized by the properties *start date* and *end date* and by the following relationships:

- (i) *dataset1* and *dataset2*: the two interrelated *Dataset* entities;
- (ii) *semanticsClass* and *semanticsScheme*: the semantics of the relationship according to a given classification of relationships between Datasets, respectively described as a *Class* entity and a *Scheme* entity.

**Publication\_Dataset** *Publication\_Dataset* entities represent relationships between a *Publication* and a *Dataset*, together with the duration and the semantics of the relationship. As such it is characterized by the properties *start date* and *end date* and by the following relationships:

- (i) *Publication* and *Dataset*: the two interrelated *Publication* and *Dataset* entities;
- (ii) *semanticsClass* and *semanticsScheme*: the semantics of the relationship according to a given classification of relationships between *Publications* and *Datasets*, respectively described as a *Class* entity and a *Scheme* entity.

**Person\_Result** *Person\_Result* entities represent relationships between a *Person* and a *Result*, together with the duration and the semantics of the relationship. As such it is characterized by the properties *start date* and *end date* and by the following relationships: *Person* and *Result*: the two interrelated *Person* and *Result* entities;

- (i) *semanticsClass* and *semanticsScheme*: the semantics of the relationship according to a given classification of relationships between *Persons* and *Results*, respectively described as a *Class* entity and a *Scheme* entity;
- (ii) *personAffiliation*: the *Organization* to which the *Person* creator of the *Result* was affiliated to at the time of *Result* creation.

**Result\_Project** *Result\_Project* entities represent relationships between a *Result* and a *Project*, together with the duration and the semantics of the relationship. As such it is characterized by the properties *start date* and *end date* and by the following relationships:

- (iii) *Result* and *Project*: the two interrelated *Result* and *Project* entities;

**Result\_Organization** *Result\_Organization* entities represent relationships between a *Result* and a *Project*, together with the duration and the semantics of the relationship. As such it is characterized by the properties *start date* and *end date* and by the following relationships:

- ## 2.3 Entity-Relationship model

The diagram illustrates the database schema for the Open Science Framework (OSF). It features several entities and their relationships:

- Entities (Yellow boxes):**
  - DataSource**: A central entity with a dashed line indicating a complex relationship with a group of entities (EntityRegistry, Aggregator, DataArchive, CRIS, Repository).
  - Organization 1/2**: A base organization entity.
  - Person**: A base person entity.
  - Project**: A base project entity.
  - Result**: A base result entity.
  - Funding**: A base funding entity.
  - Organization 2/2**: A specialized organization entity.
  - Discipline**: A specialized result entity.
- Relationships (Blue boxes):**
  - Project\_Organization**: Connects Organization 1/2 and Organization 2/2.
  - Result\_Organization**: Connects Organization 1/2 and Result.
  - Person\_Result**: Connects Person and Result.
  - Result\_Project**: Connects Result and Project.
  - Organization\_Funding**: Connects Organization 2/2 and Funding.
  - Funding\_Funding**: Connects Funding and Funding.
- Cardinalities and Relationship Names:**
  - DataSource** to **Organization 1/2**: 0,1 to 0,N (dataSources).
  - Organization 1/2** to **Project\_Organization**: 1,1 to 0,N.
  - Organization 1/2** to **Result\_Organization**: 1,1 to 0,N.
  - Organization 1/2** to **Person**: 1,1 to 0,N (personAffiliation).
  - Person** to **Person\_Result**: 0,N to 1,1.
  - Person** to **Result**: 0,N to 0,N (creations).
  - Result** to **Result\_Organization**: 1,1 to 1,1.
  - Result** to **Result\_Project**: 0,N to 1,1 (projects).
  - Result** to **Discipline**: 0,N to 1,1 (disciplines).
  - Result** to **Person\_Result**: 0,N to 1,1 (creators).
  - Project** to **Project\_Organization**: 1,1 to 1,1 (involvedOrg).
  - Project** to **Result\_Project**: 0,N to 1,1 (results).
  - Project** to **Organization\_Funding**: 1,1 to 1,1 (funds).
  - Organization\_Funding** to **Funding**: 1,1 to 0,N (funders).
  - Funding** to **Funding\_Funding**: 0,N to 1,1.
  - Funding** to **Organization\_Funding**: 0,N to 1,1 (funding).

Figure 3 – E-R model: main, linked, static and structural entities

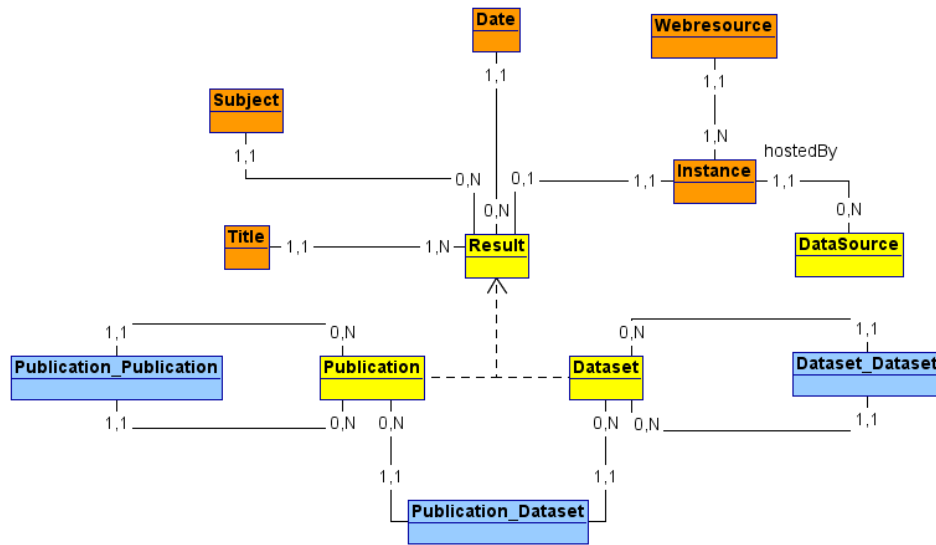


Figure 4 - E-R model: Result entities

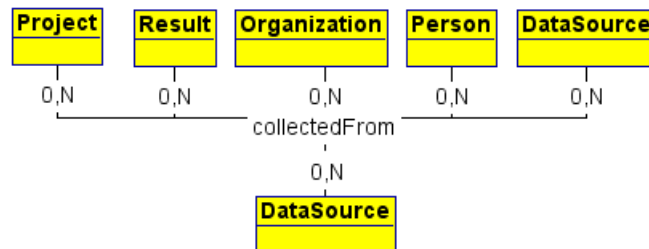


Figure 5 - E-R model: provenance relationships

Table 2 – E-R Schema: class properties

<b>Result</b> <ul style="list-style-type: none"> <li>title</li> <li>subtitle (optional)</li> <li>dateOfCreation (optional)</li> <li>description</li> <li>→ creators (0 or N Person_Result)</li> <li>→ result<sup>-1</sup> (0 or N Result_Organization)</li> <li>→ instances (0 or N Instance)</li> <li>→ subjects (0 or N Subject)</li> <li>→ collectedFrom (0 or 1 Data Source)</li> <li>→ languageClass (1 Class)</li> <li>→ languageScheme (1 Scheme)</li> <li>→ licenseClass (1 Class)</li> <li>→ licenseScheme (1 Scheme)</li> </ul>	<b>Person</b> <ul style="list-style-type: none"> <li>persistentIdentifier</li> <li>firstName</li> <li>secondNames</li> <li>fax</li> <li>email</li> <li>phone</li> <li>creations (0 or N Person_Result)</li> <li>• contactPerson<sup>-1</sup> (0 or N Project_Organization)</li> <li>• nationalityClasses (0 or 1 Class)</li> <li>• nationalityScheme (0 or 1 Scheme)</li> <li>→ collectedFrom (0 or 1 DataSource)</li> </ul> <p><i>Constraints:</i></p> <ul style="list-style-type: none"> <li>This.nationalityClass must be associated to This.nationalityScheme</li> </ul>	<b>Dataset</b> <b>(isA Result)</b> <ul style="list-style-type: none"> <li>device (optional)</li> <li>→ dataset1<sup>-1</sup> (0 or N Dataset_Dataset)</li> <li>→ dataset2<sup>-1</sup> (0 or N Dataset_Dataset)</li> <li>→ dataset<sup>-1</sup> (0 or N Publication_Dataset)</li> <li>→ typeClass (1 Class)</li> <li>→ typeScheme (1 Scheme)</li> </ul> <p><i>Constraints:</i></p> <ul style="list-style-type: none"> <li>This.typeClass must be associated to This.typeScheme</li> </ul>
---	--	--

<p><b>Publication</b> <b>(isA Result)</b></p> <ul style="list-style-type: none"> <li>• publisher (optional)</li> <li>• embargoEndDate (optional)</li> <li>→ publication<sup>1</sup> (0 or N Publication_Publication)</li> <li>→ publication2<sup>1</sup> (0 or N Publication_Publication)</li> <li>→ publication<sup>1</sup> (0 or N Publication_Dataset)</li> <li>→ typeClass (1 Class)</li> <li>→ typeScheme (1 Scheme)</li> </ul> <p><i>Constraints:</i></p> <ul style="list-style-type: none"> <li>• This.typeClass must be associated to This.typeScheme</li> <li>• This.licenseClass must be associated to This.licenseScheme</li> </ul>	<p><b>Project</b></p> <ul style="list-style-type: none"> <li>• persistentIdentifier</li> <li>• webSiteURL (optional)</li> <li>• acronym</li> <li>• title</li> <li>• start_date</li> <li>• end_date</li> <li>• call_identifier (optional)</li> <li>• keywords (optional)</li> <li>• duration (derived: from start_date and end_date)</li> <li>• EC_SC39 (optional)</li> <li>• organizations (1 or N Project_Organization)</li> <li>→ funds (0 or N Project_Funding)</li> <li>→ collectedFrom (0 or 1 DataSource)</li> </ul> <p><i>Notes</i></p> <p>For EC projects <i>OriginalIdentifier</i> contains the grant agreement number.</p>	<p><b>Organization</b></p> <ul style="list-style-type: none"> <li>• persistentIdentifier</li> <li>• legal short name</li> <li>• legal name</li> <li>• webSiteURL</li> <li>• logoURL (optional)</li> <li>• EC_LegalBody (boolean) (optional)</li> <li>• EC_LegalPerson (boolean) (optional)</li> <li>• EC_NonProfit (boolean) (optional)</li> <li>• EC_ResearchOrganization (boolean) (optional)</li> <li>• EC_HigherEducation (boolean) (optional)</li> <li>• EC_InternationalOrganization EurInterests (boolean) (optional)</li> <li>• EC_Enterprise (boolean) (optional)</li> <li>• EC_SMEValidated (boolean) (optional)</li> <li>• EC_NUTScore (optional)</li> <li>→ organization<sup>1</sup> (0 or N Result_Organization)</li> <li>→ personAffiliation<sup>1</sup> (0 or N Person_Result)</li> <li>→ dataSources (0 or N Data Sources)</li> <li>→ funder<sup>1</sup> (0 or N Organization_Funding)</li> <li>→ involvedOrg<sup>1</sup> (0 or N Project_Organization)</li> <li>→ countryClass (0 or 1 Class)</li> <li>→ countryScheme (0 or 1 Scheme)</li> <li>→ collectedFrom (0 or 1 DataSource)</li> </ul> <p><i>Constraints:</i></p> <p>This.countryClass must be associated to This.countryScheme</p>
<p><b>Data Source</b></p> <ul style="list-style-type: none"> <li>• persistentIdentifier</li> <li>• officialName</li> <li>• englishName (optional)</li> <li>• webSiteURL</li> <li>• logoURL</li> <li>• contactEmail</li> <li>• accessInfoPackage</li> <li>→ hostedBy<sup>1</sup> (0 or N Instances)</li> <li>→ collectedFrom<sup>1</sup> (0 or N Instances)</li> </ul>	<p><b>Instance</b></p> <ul style="list-style-type: none"> <li>• persistentIdentifier</li> <li>→ instances<sup>1</sup> (1 Results)</li> <li>→ hostedBy (1 DataSource)</li> </ul>	<p><b>Web Resource</b></p> <ul style="list-style-type: none"> <li>• Web Resource URL</li> <li>→ webResources<sup>1</sup> (1 Instances)</li> </ul>

<ul style="list-style-type: none"> <li>→ dataSources<sup>-1</sup> (0 or N Organizations)</li> <li>→ collectedFrom (0 or N DataSource)</li> <li>→ collectedFrom<sup>-1</sup> (0 or N Person, DataSource, Organization, Project, Instance)</li> </ul>		
<b>Repository</b> <b>(isA DataSource)</b> <ul style="list-style-type: none"> <li>• OD_description (optional)</li> <li>• OD_numberOfItems (optional)</li> <li>• OD_numberOfItemsDate (optional)</li> <li>• OD_subjects (optional)</li> <li>• OD_policies (optional)</li> <li>• OD_languages (optional)</li> <li>• OD_contentTypes (optional)</li> </ul>	<b>CRIS</b> <b>(isA DataSource)</b>	<b>DataArchive</b> <b>(isA DataSource)</b>
<b>Aggregator</b> <b>(isA DataSource)</b> <ul style="list-style-type: none"> <li>→ typeClass (1 Class)</li> <li>→ typeScheme (1 Scheme)</li> </ul> <i>Constraints:</i> <ul style="list-style-type: none"> <li>• This.typeClass must be associated to This.typeScheme</li> </ul>	<b>Subject</b> <ul style="list-style-type: none"> <li>→ semanticsClass (1 Class)</li> <li>→ semanticsScheme (1 Scheme)</li> <li>→ subjects<sup>-1</sup> (1 Result)</li> </ul>	<b>Funding</b> <ul style="list-style-type: none"> <li>• persistentIdentifier</li> <li>• name</li> <li>• description</li> <li>• keywords</li> <li>• funding1<sup>-1</sup> (0 or N Funding_Funding)</li> <li>• funding2<sup>-1</sup> (0 or N Funding_Funding)</li> <li>• funders (0 or N Organization_Funding)</li> <li>• funding<sup>-1</sup> (0 or N Project_Funding)</li> </ul>
<b>Person_Result</b> <ul style="list-style-type: none"> <li>• startDate</li> <li>• endDate</li> <li>→ creators<sup>-1</sup> (1 Person)</li> <li>→ creations<sup>-1</sup> (1 Result)</li> <li>→ creatorAffiliation (0 or 1 Organizations)</li> <li>→ personRoleClass (1 Class)</li> <li>→ personRoleScheme (1 Scheme)</li> </ul> <i>Constraints:</i> <ul style="list-style-type: none"> <li>• This.roleClass must be associated to This.roleScheme</li> </ul>	<b>Publication_Publication</b> <ul style="list-style-type: none"> <li>• startDate</li> <li>• endDate</li> <li>→ publication1 (1 Publication)</li> <li>→ publication2 (1 Publication)</li> <li>→ semanticsClass (1 Class)</li> <li>→ semanticsScheme (1 Scheme)</li> </ul> <i>Constraints:</i> <ul style="list-style-type: none"> <li>• This.semanticsClass must be associated to This.semanticsScheme</li> </ul>	<b>Organization_Project</b> <ul style="list-style-type: none"> <li>• participantNumber (optional)</li> <li>• startDate</li> <li>• endDate</li> <li>→ participants<sup>-1</sup> (1 Projects)</li> <li>→ contactPerson (0 or 1 Persons)</li> <li>→ respOrganization (1 Organizations)</li> <li>→ semanticsClass (1 Class)</li> <li>→ semanticsScheme (1 Scheme)</li> </ul>

		<p><i>Notes</i></p> <ul style="list-style-type: none"> <li>For EC projects, <i>OriginalIdentifier</i> is the participant number (1 for project coordinators)</li> <li>For other projects, unless they support the concept of participant, the <i>OriginalIdentifier</i> is empty.</li> </ul>
<p><b>Dataset_Dataset</b></p> <ul style="list-style-type: none"> <li>startDate</li> <li>endDate</li> <li>→ dataset1 (1 Dataset)</li> <li>→ dataset2 (1 DataSet)</li> <li>→ semanticsClass (1 Class)</li> <li>→ semanticsScheme (1 Scheme)</li> </ul> <p><i>Constraints:</i></p> <ul style="list-style-type: none"> <li>This.semanticsClass must be associated to This.semanticsScheme</li> </ul>	<p><b>Publication_Dataset</b></p> <ul style="list-style-type: none"> <li>startDate</li> <li>endDate</li> <li>→ publication (1 Publication)</li> <li>→ dataset (1 DataSet)</li> <li>→ semanticsClass (1 Class)</li> <li>→ semanticsScheme (1 Scheme)</li> </ul> <p><i>Constraints:</i></p> <ul style="list-style-type: none"> <li>This.semanticsClass must be associated to This.semanticsScheme</li> </ul>	<p><b>Funding_Funding</b></p> <ul style="list-style-type: none"> <li>startDate</li> <li>endDate</li> <li>→ funding1 (1 Funding)</li> <li>→ funding2 (1 Funding)</li> <li>→ semanticsClass (1 Class)</li> <li>→ semanticsScheme (1 Scheme)</li> </ul> <p><i>Constraints:</i></p> <ul style="list-style-type: none"> <li>This.semanticsClass must be associated to This.semanticsScheme</li> </ul>
<p><b>Organization_Funding</b></p> <ul style="list-style-type: none"> <li>startDate</li> <li>endDate</li> <li>→ funder (1 Organization)</li> <li>→ funders<sup>-1</sup> (1 Funding)</li> <li>→ semanticsClass (1 Class)</li> <li>→ semanticsScheme (1 Scheme)</li> </ul> <p><i>Constraints:</i></p> <ul style="list-style-type: none"> <li>This.semanticsClass must be associated to This.semanticsScheme</li> </ul>	<p><b>Project_Funding</b></p> <ul style="list-style-type: none"> <li>startDate</li> <li>endDate</li> <li>→ funding (1 Funding)</li> <li>→ funds<sup>-1</sup> (1 Projects)</li> <li>→ semanticsClass (1 Class)</li> <li>→ semanticsScheme (1 Scheme)</li> </ul> <p><i>Constraints:</i></p> <ul style="list-style-type: none"> <li>This.semanticsClass must be associated to This.semanticsScheme</li> </ul>	<p><b>Project_Organization</b></p> <ul style="list-style-type: none"> <li>startDate</li> <li>endDate</li> <li>→ involvedOrg (1 Organization)</li> <li>→ organizations<sup>-1</sup> (1 Project)</li> <li>→ contactPerson (1 Person)</li> <li>→ semanticsClass (1 Class)</li> <li>→ semanticsScheme (1 Scheme)</li> </ul> <p><i>Constraints:</i></p> <ul style="list-style-type: none"> <li>This.semanticsClass must be associated to This.semanticsScheme</li> </ul>
<p><b>Result_Project</b></p> <ul style="list-style-type: none"> <li>startDate</li> <li>endDate</li> <li>→ project (1 Project)</li> <li>→ result (1 Result)</li> <li>→ semanticsClass (1 Class)</li> <li>→ semanticsScheme (1 Scheme)</li> </ul> <p><i>Constraints:</i></p>		

<ul style="list-style-type: none"><li>• This.semanticsClass must be associated to This.semanticsScheme</li></ul>		
--	--	--

### 3 Data population

As shown in Figure 1 the OpenAIREplus infrastructure collects entity information from an expert-validated entity pool and from data sources of different typologies, such as repositories, CRISs, dataset archives, aggregators, entity registries. All these contain information relative to different interrelated entities of the OpenAIRE data model. In particular, as shown by **Error! Reference source not found.**, data sources of the same typology may deliver metadata records which contain information relative to different entities. For example, an OpenAIRE compliant repository delivers information packages (i.e., metadata records) which contain information about the publication result, the persons who created such result, the projects funding such result, and the instances relative to the result; while a DRIVER compliant repository does not contain information about project entities.

In the following we shall call **original entities** the entities collected from data sources, hence from “authoritative” providers of data, onto the Information Space. The layer of original entities includes entities and relationships between them as collected from the different data sources (i.e., mapped from their original structure onto the one of the OpenAIRE data model). The layer is “stateless”, in the sense that entities have “reproducible” identifiers derived by combining identifiers of the entities in the original data sources with a data source identifier assigned by OpenAIRE. In other words, if the same entity or relationships is collected more than once from the same data source, it will be transparently overridden.

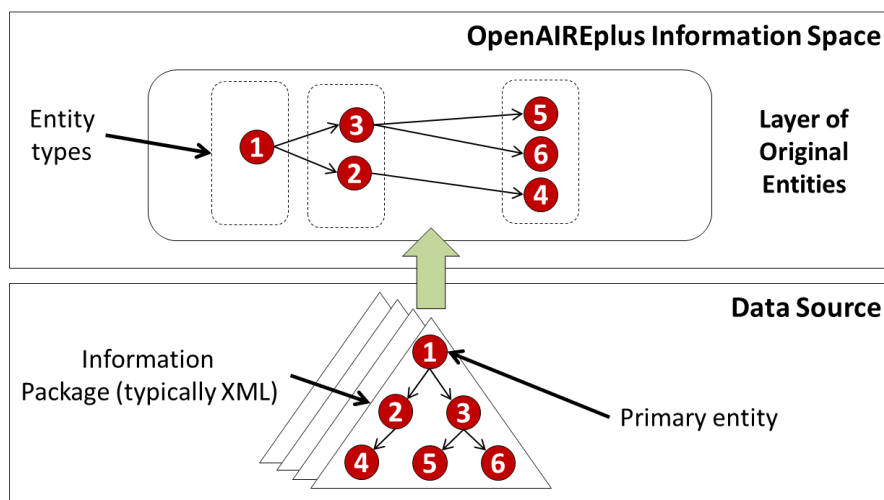


Figure 6 – Entity layers: original entities

With respect to data population, hence with the process of collecting data from data sources and map them onto the information space, this section we will introduce the notions of:

- **Information packages** and **population workflows**: entity ingestion from data sources into the information space;
- How to assign a **stateless (and permanent) identifiers** to the entities when they enter the information space;

For each of these aspects, we shall provide an explanation of the problem, an extension of the data model to handle the problem, and, where necessary, a solution to the problem using the updated model.

### 3.1 Information packages

In OpenAIRE entities are collected from external data sources in the form of “information packages”. This notion aims at generalizing the OpenAIRE scenario of bibliographic metadata records import from repository data sources to other data source typologies and other types of (primary) entities. In particular, we shall call *information package* a file in some interpretable format (e.g. XML), which contains *identifier* and *information* (e.g., properties) relative to one entity, called *primary entity*, of a given entity type. An information package may contain information (but not necessarily the identifier) relative to other entities (of likely different entity types), called *sub-entities*, which must be directly or indirectly associated with the package primary entity. Figure 6 shows an example of an information package whose primary entity is 1: for example, an information package from OpenDOAR is relative to a repository data source and can be identified by the relative OpenDOAR identifier. Its sub-entities are those from 2 to 6: for example, an OpenDOAR package also contains information about the organization responsible for the repository data source.

### 3.2 Population Workflows

Original entities are collected from information packages originating from various data sources. We call *population workflow* the process that takes the information package from a data source, extracts its primary entity and its related entities, and stores them into the OpenAIRE information space. A workflow is therefore dependent on:

- The *data source typology* (including the expert-validate entity pool);
- The *access method*, namely (i) protocol required to get the data (e.g., OAI-PMH, JDBC, FTP) and (ii) relative access configuration (e.g., entry point, parameters, etc.);
- The *primary entity type* of the information packages;
- The *XML structure of the information packages* at hand, which depends on the primary entity type.

Note that data sources of the same typology may deliver information packages relative to different primary entity types (in general we can assume they will do it from different access points). For example, CRIS systems may expose through OAI-PMH both publication or project primary entities.

**Information package structure (OpenAIREplus guidelines)** The OpenAIREplus infrastructure will include services capable of handling automated collection of entities from data sources according to given population workflows. To this aim, the OpenAIRE guidelines will describe which XML information packages structure should be expected for each population workflow triple

*<datasource typology, access method, primary entity type> → XML information package structure*

available to the system. WP6 will develop services to automatically process the information packages and insert the relative entities onto the information space.

**Information package heterogeneity and harmonization** Unfortunately, the “raw” information packages exported by data sources will likely not match the information

package structures to be identified in the previous step. For example, CRIS systems generally support OAI-PMH harvesting of information, but may export information packages relative to the same entities (e.g., projects, publications) in different XML formats. To this aim, WP6 will update its transformation services in order to map the specific structures exposed by a data source through a given workflow so that they match the expected information package structure.

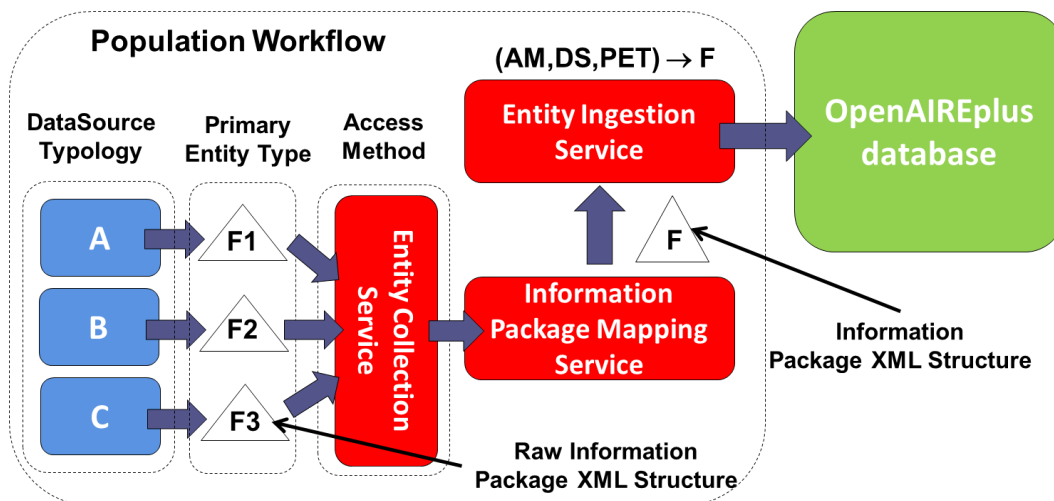


Figure 7 – Information Packages: ingestion workflows (AM = Access Method, DS = Data source typology, PET = Primary Entity Type, F = information package Format): data sources of the same typology export the same primary entity of the same type through different “raw” information package format structures.

### 3.3 Identity of original entities

Original entities reach the information space from different workflows. Once they enter the information space they must be assigned a unique “stateless” identifier. The data sources of such entities are not under the OpenAIRE infrastructure control and may in any moment decide to delete, update, or add new entities or relationships between them. Hence, it is particularly important to make sure such identifiers are generated from the incoming information packages in a stateless and stable way that is “if the same entity enters the information space at different times, it will be assigned the same identifier”.

To this aim, the OpenAIRE infrastructure constructs identifiers for primary entities and sub-entities in an information package by combining three levels of scope: data sources, relative primary entities, and sub-entities of such primary entities. More specifically:

- *Infrastructure scope*: all data sources are registered and assigned a unique identifier in OpenAIRE;
- *Data source scope*: information packages from the same data source contain one primary entity with an identifier which is unique in the context of the data source;
- *Primary entity scope*: information packages may contain a number of sub-entities relative to the primary entity; unlike primary entities, sub-entities may not necessarily come with an identifier (data source scope)<sup>1</sup> and can be generally

<sup>1</sup> Absence of identifiers may be due to the fact they do not have an identifier on the original data source, e.g., authors of publications in bibliographic metadata records, or to the structure of the information package export format, which focuses on the main entity only.

uniquely identified in the scope of the primary entity based on their properties. The process of identification of such “unique information” is very much dependent on the given information package structure.

**Primary entity identifiers** The process of generation of stateless identifiers for primary entities is based on a *data source scope strategy*. Independently of the workflows, the type of entity, and the data source kind, primary entity identifiers are always obtained by concatenating the data source identifier with the primary entity identifier (using and underscore):

*datasourceID\_mainEntityID*

Although one may consider an *infrastructure scope strategy*, where the assumption is that all primary entities identifiers are persistent identifiers, therefore unique across several data sources, in OpenAIRE this is not generally the case, hence we adopt a common and safe strategy of identifier generation.

**Sub-entity identifiers** Assigning identifiers to sub-entities can be performed following different strategies, some more “optimistic” and some more “pessimistic” about the ability of inferring unambiguous “unique information” for sub-entities from their properties in the information package. For example, one may assume that:

- *Infrastructure scope strategy*. sub-entities with the same “unique information” are collapsed in the same entity across different data sources (infrastructure scope). Entity splitting will identify and solve possible entity “overloads” in a second stage.

*uniqueInformation*

- *Data source scope strategy*. Sub-entities with the same “unique information” are collapsed in the same entity but only within the same data source scope. Entity splitting will identify and solve possible entity “overloads” in a second stage.

*datasourceID\_uniqueInformation*

- *Primary entity scope strategy*. Sub-entities with the same “unique information” are collapsed in the same entity but only within the same primary entity scope (see Figure 8). De-duplication of entities will solve redundancy in a second stage.

*datasourceID\_mainEntityID\_uniqueInformation*

Assigning identifiers to sub-entities follows different strategies depending on the specific workflow, hence the relative information packages structure.

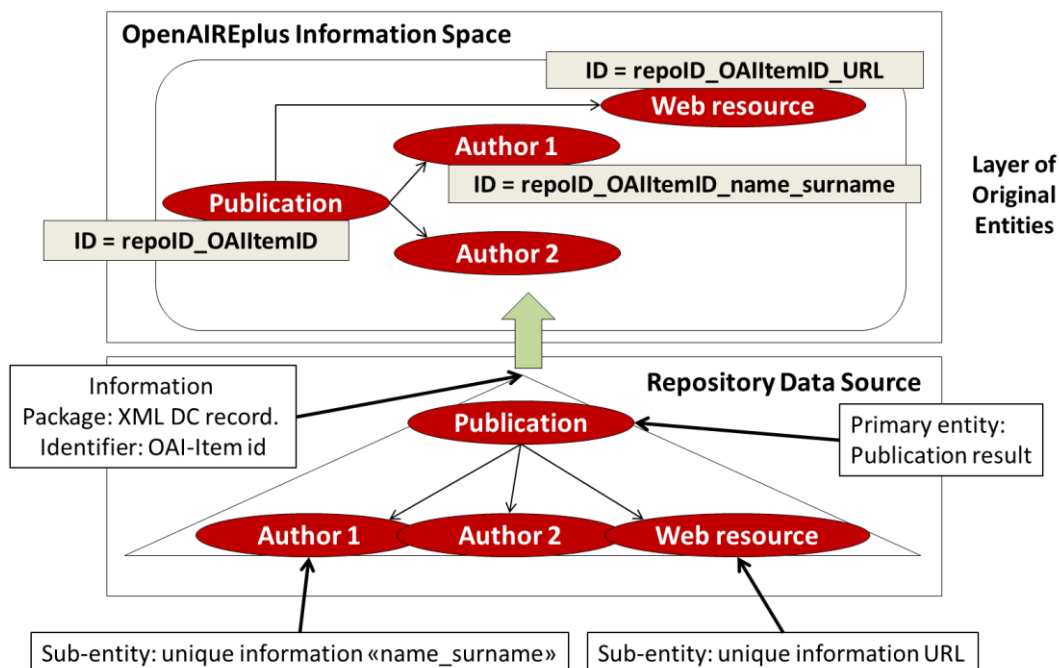


Figure 8 – Assigning unique identifiers to sub-entities: primary entity scope

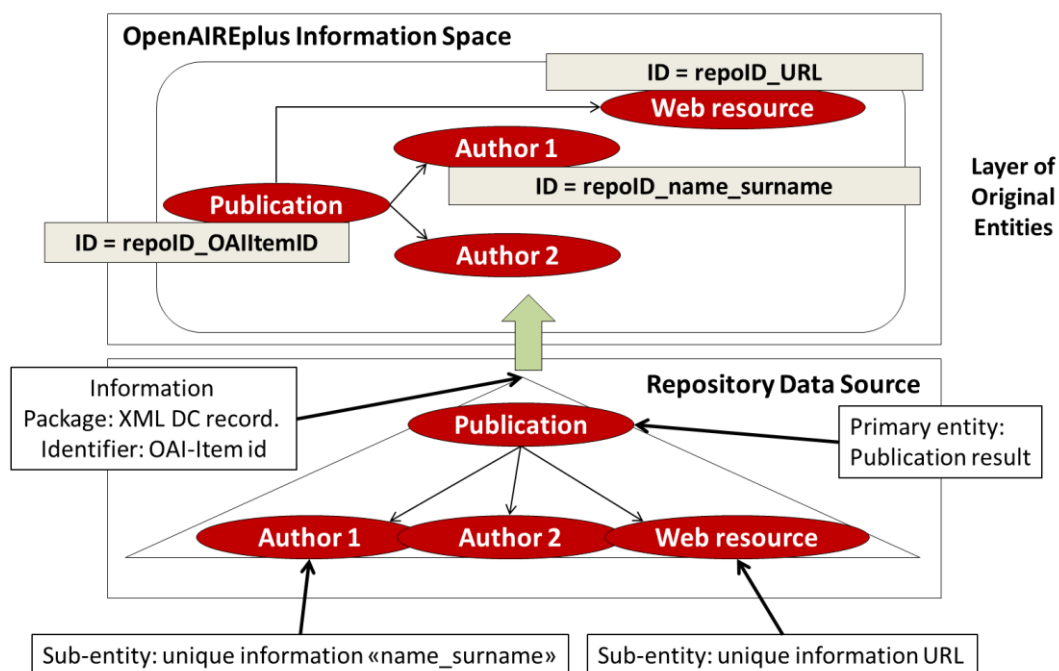


Figure 9 – Assigning unique identifiers to sub-entities: data source scope

## 4 Data inference

As highlighted in Section 1, the process of mapping **information packages** onto **original entities** is not enough to ensure a high quality information space. Indeed, the autonomy of the data sources brings in three main issues:

- A certain degree of information duplication: information about the same real-world entities will be collected from several data sources;
- A certain degree of entity disconnection: data sources deliver entities and relationships between them, which are generally limited to the boundaries of the data source; cross-data source relationships are generally missing;
- A certain degree of inaccuracy: data sources may deliver only a portion of the entity information required by the OpenAIREplus information space.

To tackle such issues, original entities serve as input to a data inference process, which yields a layer of **inferred entities** resulting from resolving entity inaccuracy and disconnection through a mix of human intervention and information inference services. Inferred entities logically “override” original entities in the information space to provide an enhanced view of the information space, namely the layer of **visible entities**, which will be effectively accessible to end-users and applications. The data inference process, which leads from original entities to the new set of inferred entities, consists of four main phases, at the end of which the set of inferred entities is actually transferred into the information space and used to give life to a new set of visible entities. The phases are executed over a clone of the OpenAIRE information space, containing a copy of the current set of original entities. The phases are:

1. *De-duplication of original entities*: the phase returns a first set of inferred entities which delivers a set of disambiguated visible entities;
2. *Data inference actions over the set of visible entities*: the phase returns a second set of inferred entities, which enriches the first set and therefore further updates the set of visible entities; this second set enriches the information space of new entities but may introduce further duplicates;
3. *De-duplication of the set of visible entities*: the phase returns a third and final set of inferred entities which delivers a set of disambiguated visible entities;
4. *Transfer the final set of inferred entities into the information space*: the set of inferred entities is moved from the information space clone into the production information space in order to re-calculate the set of visible entities.

The whole data inference process is always executed over original entities and can be re-executed anytime to generate a new visible entity layer. As such, there is no need to keep in the information space a history of inference actions, since both original entity and inferred entity layers are intended to be always “refreshed” by executing a data population operation or a data inference operation. Still, the information space data model needs to be extended to enable the distinction between original and inferred entities and the identification of visible entities as a logical overlap of the former two.

In this section we shall first introduce the data inference actions which can be executed in OpenAIRE and comment on how such actions are to be encoded in the OpenAIRE data model.

## 4.1 Inference actions

OpenAIREplus original entities can be subject to the following *data inference* actions, performed by data mining and inference services and by humans:

- *Merging of entities;*
- *Adding new inferred entities (relationship between main entities);*
- *Updating of properties of an original entity;*
- *Removing original entities.*

Such actions may be caused by different data inference workflows, which may or may not involve humans. Examples in the first OpenAIRE service settings are:

- *End-user feedbacks*, when approved by data curators;
- *Data curator* edit activities;
- *Automatic inference algorithms*: (i) similarity relationships between result entities, (ii) information extraction from full-text of publications, which may infer title, authors, organizations, and emails of publication entities, (iii) citation management, which may lead to the introduction of publication results not available to the information space, (iv) subject classification, which may lead to the introduction of new entities, etc.
- *End-user claim/validation actions* through the OpenAIRE portal, e.g., project coordinators validating publication-project relationships.
- *Deduplication services* identifying a set of entities of the same type which all correspond to the one and the same real-world entity.

In the following we shall focus on how to modify the OpenAIRE data model in order to represent the outcome of data inference actions on the information space and to resolve data inference conflicts which may arise between original and inferred entities. For each of these aspects, we shall provide an explanation of the problem, an extension of the data model to handle the problem, and, where necessary, a solution to the problem using the updated model.

### 4.1.1 Trust and Inference Provenance of Entities

Inference actions, hence the entities they bring onto the information space, may have different **levels of trust**, depending on their service or agent which generated them, namely **entity inference provenance**. Main entities and linked entities in the model are therefore enriched with three properties:

- *Inferred*: a flag which is set to TRUE when the entity is created as a consequence of an inference action, including de-duplication.
- *deletedByInference*: a flag which is set to TRUE when some inference process establishes that the entity should be removed from the visible entities.
- *Trust*: a 0 to 1 value which establishes the level of credibility of the information. By definition, original entities are assigned a OE level of trust, while entities from the expert-validated entity pool a value  $EVE > OE$ . Inferred entities are assigned a value  $IE < OE$  which depends on the quality and reliability of the automatic inference mechanisms or of the humans bringing the information.
- *Inference Provenance*: a value from a controlled dictionary which encodes the algorithm/process/service which brought the inferred entity in the information

space, e.g., “de-duplication”, or deleted, i.e., made invisible, an original entity from the information space. In the case of original entities the property is initially set to NULL.

### 4.1.2 Merging of entities

The de-duplication service operates over the set of original entities of the same main type (i.e., Person, Organization, Result, DataSource, Project) and identifies groups of entities which are redundantly describing the same entity. For each group, the service identifies a *representative entity*, which by default is the one with the higher level of trust among the duplicates (if more entities have the same level of trust, the one with the shorter identifier is selected). The remaining duplicate entities (*merged entities*) will not be visible to end-users and will point to the representative entity. However, in the case merged entities bear values for properties which are unavailable (empty values) for the representative entity, such values will be used for indexing and visualization purposes (“shadow entity” strategy). Hence, the ordering of the entities plays an important role in establishing their usage within the information space.

The second de-duplication phase occurs after a number of inference actions which may introduce further duplication. In this context, if a group of duplicates includes a representative entity, then all entities merged into the latter are included in the new group and a new representative entity is elected.

Merging introduces further issues, which regard the strategy to be adopted to cope with linked entities, static entities, and structural entities associated with entities that were merged and should therefore be excluded from the set of visible entities. In other words, in order to offer a coherent view of the space, all entities “surrounding” merged entities must be altered in order not to be associated with the representative entity. For example, merging two Publication p1 and p2 into p1 implies that Project, Authors, Instance, and Subject entities of p2 will be linked to p1 (avoiding duplicated associations in the case p1 already points to the same entities). Note that Title and Date structural entities are considered part of a single Result and will not be linked to the representative entity.

To this aim the data model needs to represent the possibility of adding a new relationship with a status *Inferred*, which encodes such re-linking, without losing the original relationship (in order to “clean” the inferred information and return to the original layer of entities). As shown in Figure 10, this requires a major change, since explicit relationships between entities, such as *collectedFrom*, become **explicit relationship entities** (in green in the Figure) which include properties capable of capturing such status. Specifically, the model will include:

for each main entity type:

- A relationship *mergedWith*, which points from one entity to the group-representative entity;

for all explicit relationships (not modelled through linked entities):

- Properties *Inferred*, *deletedByInference*, *Trust*, and *Inference Provenance*

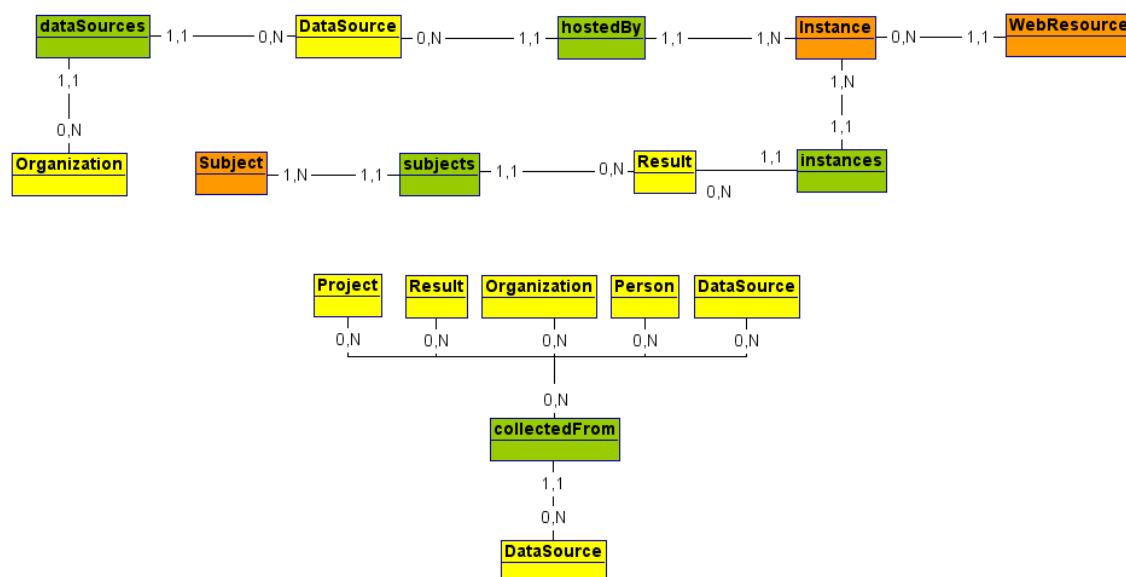


Figure 10 – Extension to data model to cope with inferred entities: explicit relationships become entities

When merging a set of entities  $e_1, \dots, e_K$  into a representative entity  $e_j$  two main strategies are adopted, which regard explicit relationship entities and linked entities. As an example we shall consider Result entities.

**Explicit relationship entities** Result entities are surrounded by the following explicit relationship entities: *Instances*, *CollectedFrom*, and *Subjects*.

*Instances* entities associate Results with a set of Instance entities (and in turn a set of WebResource entities). In this case, all Instance entities of merged entities are to be connected to the representative entity, so that all files of all Results are visible. To this aim, for each Result in  $r_2, \dots, r_K$  all *instances* relationships are set the flag *deletedByInference* to TRUE. All such relationships are cloned, i.e., the same values for all properties, but with the flag *inferred* set to TRUE, *deletedByInference* set to FALSE, *InferenceProvenance* set to "de-duplication", *Trust* set to a given value, and as target Result entity the representative entity.

The same strategy is adopted for the other explicit relationship entities surrounding Result entities, namely *CollectedFrom* and *Subjects*. Existing explicit relationship entities for merged entities are made logically invisible (set the flag *deletedByInference* to TRUE) and clone relationship entities are created which point to the representative entity and have the flag *inferred* set to TRUE. In the case the source entity of the relationships to be cloned is already linked with the representative entity, the relationship is not created (e.g., merged entities and representative entity share a subset of Subject entities).

**Linked entities** Result entities are surrounded by the following linked entities:

- Relationship *projects* to Result\_Project linked entities
- Relationship *creators* to Person\_Result linked entities
- (If *Publication* entity) relationship *pub1* and *pub2* to Publication\_Publication linked entities
- (If *Publication* entity) relationship *publication*<sup>-1</sup> Publication\_Dataset linked entities
- (If *Dataset* entity) relationship *dataset1* and *dataset2* to Dataset \_ Dataset linked entities

- (If *Dataset* entity) relationship *dataset*<sup>-1</sup> Publication\_Dataset linked entities

For each linked entity involving merged entities, the same two steps highlighted above take place. Existing linked entities for merged entities are made logically invisible (set the flag *deletedByInference* to TRUE) and clone linked entities are created which point to the representative entity and have the flag *inferred* set to TRUE. In the case the source entity of the linked entity to be cloned is already linked with the representative entity, the relationship is not created (e.g., merged entities and representative entity share a subset of Project entities).

**Note:** as an optimization of the merging process, which leads to cloning of entities hence to a possible increase of the entities in the information space, we could keep track of the changes inferred by merge actions within linked entities and explicit relationship entities. For each linked entity type and explicit relationship entity type connecting the entities A and B, the model would feature two associations *original\_<entityA>* and *original\_<entityB>* as an addition to the associations *<entityA>* and *<entityB>* which such entities contain. The former pair keeps memory of the original pointers *<entityA>* and *<entityB>* to the entities in A and B before the merge took place; the inferred pointers will substitute the original ones and be kept in *<entityA>* and *<entityB>*.

#### 4.1.3 Adding an inferred entity

Some inference processes may lead to the introduction of new main entities together with linked entities or explicit relationship entities. Consider for example the citation inference process, which returns a set of Publication\_Publication entities with semantics of type "citation" and may add new Publication entities relative to the citations that are not in the information space. In this case, all such entities are added with the *inferred* flag set to TRUE, *deletedByInference* set to FALSE, *InferenceProvenance* set to "citation\_inference", *Trust* set to a given value.

#### 4.1.4 Updating an original entity

Updating an entity consists in selecting one main entity and updating its property values or its relationships with other entities. Updates are modeled as additions of entities obtained from cloning the updated entity. The new entity is tagged as *inferred* and has a level of trust UE, to be used in the second de-duplication phase. Depending on the value, the entity may become the representative entity of a group which includes its original clone or be used as "shadow" entity in such a group.

##### 4.1.1 Removing an entity

An entity can be removed from the information space by setting its flag *inferredAsDeleted* to TRUE. As any other inference action, the action is accompanied by its inference provenance information and level of trust.

## 5 General data management issues

### 5.1 Administrative entity properties

All entities feature the following administrative properties:

- *Timestamp of creation* (linked entities have a start date and an end date inherited by the CERIF semantic layer);

All entities but Structural and Static entities:

- *Timestamp of inference action, i.e., deleteByInference, mergedWith*
- A flag *visible*, set to TRUE if the entity is a visible entity, to FALSE otherwise