

**SEVENTH FRAMEWORK PROGRAMME  
CAPACITIES**



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

**OPENAIRE**

**Grant Agreement 246686**

“Open Access Infrastructure for Research in Europe”



**OpenAIRE Data Model Specification**

Deliverable Code: D5.1

## Document Description

### Project

Title:	OPENAIRE, Open Access Infrastructure for Research in Europe
Start date:	1 <sup>st</sup> December 2009
Call/Instrument:	INFRA-2007-1.2.1
Grant Agreement:	<b>246686</b>

### Document

Deliverable number:	D5.1
Deliverable title:	OpenAIRE Data Model Specification
Contractual Date of Delivery:	31 <sup>st</sup> of March 2010
Actual Date of Delivery:	15 <sup>th</sup> of May 2010
Editor(s):	Paolo Manghi
Author(s):	Paolo Manghi
Reviewer(s):	Michele Artini, Magchiel Bijsterbosch, Wolfram Horstmann, Samuele Kaplun, Natalia Manola, Salvatore Mele, Jochen Schirrwagen, Tim Smith
Participant(s):	
Workpackage:	WP5
Workpackage title:	OpenAIRE System Back-end: storage and mediation services
Workpackage leader:	CNR
Workpackage participants:	NKUA, ICM, CERN, UNIBI
Distribution:	Public
Nature:	Deliverable
Version/Revision:	v 23
Draft/Final:	Final
Total number of pages: (including cover)	
Web Resource name:	



Key words:

## Disclaimer

This document contains description of the OPENAIRE 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/>)



OPENAIRE 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>1 Scenario</b> .....	<b>8</b>
<b>2 Data definition</b> .....	<b>10</b>
2.1 Detailed description .....	10
2.2 Entity-Relationship model .....	13
2.3 Relational Schema .....	15
<b>3 Implications for the definition of OpenAIRE Guidelines</b> .....	<b>19</b>
<b>4 Collaborative refinement of the Information Space</b> .....	<b>21</b>

## Table of Figures

Figure 1 - CORDIS Classification Scheme: Specific programmes .....	13
Figure 2 – OpenAIRE Entity-Relationship model .....	14
Figure 3 – OpenAIRE Relational Schema .....	16
Figure 4 – Controlled Vocabularies Relational Schema: “vocabulary tables” .....	16

## Summary

The OpenAIRE web site will offer functionalities for administrators, anonymous and registered users to manage an Information Space of FP7-funded open access publications. The aim of this document is to describe the conceived structure and semantics of this Information Space, i.e., the *Open AIRE data model*, by providing an abstract definition of its main entities and the relationships between them.

In this definitional process, the intended interaction (Task 7.2) between the OpenAIRE Information Space and the EC Participant Portal software systems (e.g., CORDIS, NEF, SESAM) plays an important role in the specification of project data, i.e., how project data should be described, stored and exported in OpenAIRE. At this stage of interaction, we shall consider as authoritative sources for project data, the SESAM system and the CORDIS web site, which contain information on publications, reports, projects and participants in FP7. As such, we shall store project information adopting a schema inspired by SESAM and CORDIS. Furthermore, with regard to data export, the results of projects such as CERIF and KE-CRIS-OAR are taken into account as inspiring guidelines.

Finally, the document will analyze the impact of the OpenAIRE data model on the repository managers willing to expose, through OAI-PMH interfaces, metadata conforming to the OpenAIRE Information Space, so as to avoid double ingestion work to their researchers.

## 1 Scenario

The OpenAIRE web site will offer functionalities for administrators, anonymous and registered users to manage an Information Space of FP7-funded open access publications. In particular:

- *anonymous users* will be able to search and consult the Space;
- *registered users* will be able to insert and update content in the Space;
- *administrators* will have full access and rights to such Space and in particular are in charge of validating/invalidating insertions/deletions/updates from registered users.

In our reasoning we generalize the concept of publication to that of *project result*, so as to be able of including further kinds of research outputs whose production is “economically supported” by a EC project. Examples of project results may range from traditional publications to experimental data, software products, patents, books, ORE aggregation (i.e., compound object web representations). Similarly, we extend the notion of author to that of *person*. Other typologies of person may be required in the future, so as to describe future scenarios of interest.

Currently, OpenAIRE registered users can ingest information relative to *project results of type publication, persons, organizations* (companies, research centers or institutions) in order to link them together with so-called *authorships*. An authorship represents the fact that a given *person* has (co-)authored a given *publication project result* while being affiliated with a given *organization*. 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). Accordingly, an instance of a project result is represented as a combination of zero or more *web resources* relative to the project result and of the internet *data source* (e.g., OAI-PMH accessible repositories, FTP sites, web sites) from which such resources are accessible/available. **Note that:** 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 that 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).

Of crucial interest to the OpenAIRE Information Space is also the identification of the European *projects* which co-funded the research that has led to a given publication. In particular, data describing European *projects* will be fetched from authoritative EC databases (currently SESAM but others are being contacted, such as CORDIS and NEF), together with the *organizations* or *persons* (ERC projects) which are *participants* of projects. While project data will authoritatively originate in EC databases, information about organizations may be also ingested by users, for example to complete authorships information in the database.

**Project data models** Focus studies on the results of CERIF project and the KE-CRIS-OAR initiative are being carried out in WP7. In particular, with regard to project, organization and publication data, it will be considered to adopt the description schemes proposed by such projects as models of the correspondent entities in the OpenAIRE data model. In this document, however, we shall proceed taking inspiration from the choices made by the EC databases on that respect, in particular by SESAM and CORDIS web site.



**Outline** In the following, Section 2 provides a detailed description of the entities that come into play by providing an Entity Relationship model and a relational database representation, Section 3 provides the requirements of input to the definition of the OpenAIRE Guidelines for Repository Managers (to be produced in WP3-T3.3), Section 4 describes how the data model could be extended so as to support collaborative refinement of the data.

## 2 Data definition

### 2.1 Detailed description

**Project Result (or Result)** A project Result is described by a result *kind* (e.g., publication, research data, software product), a *type*, which depends on the given kind (e.g., for publications we have "article", "book", "manual", etc.), a *date of creation*, a *description*, a *publisher*, a *language*, a list of *keywords*, an *access mode* (e.g., license to access the result) and an *embargo end-date* (empty if the access kind does not imply an embargo). A Result is associated with (i) the set of its Authorships, (ii) the set of its Instances, if any (which can be located at different Data Sources), and (iii) the Projects which co-funded the research underlying the Result.

*Note:* the property *access mode* can have one of the following values: "open access", "restricted" or "embargo". If it is "embargo", the field *embargo end-date* contains the date after which the publication becomes open access.

In SESAM publication results are described as follows:

I_ID	NOT NULL	NUMBER (10)
I_PROJECT_ID		NUMBER (10)
I_NUMBER		NUMBER (10)
ST_TITLE		VARCHAR2 (100)
ST_AUTHOR		VARCHAR2 (100)
ST_TITLE_PERIODIC		VARCHAR2 (100)
ST_PERIODIC_NUMBER		VARCHAR2 (100)
ST_PUBLISHER		VARCHAR2 (100)
ST_PLACE_PUB		VARCHAR2 (100)
ST_RELEVANT_PAGES		VARCHAR2 (30)
ST_PERMANENT_ID		VARCHAR2 (20)
ST_ACCESS_OPEN		VARCHAR2 (1)
ST_STATUS		VARCHAR2 (50)
I_VERSION		NUMBER (5)
DA_DATE_PUB	NOT NULL	DATE
I_EMBARGO_MONTHS		NUMBER (3)
DA_EMBARGO_END		DATE

Note that the attributes that were found relevant to a resource are a subset or can be derived from those used in SESAM. If required, in the future we can devise ways to import publication data from SESAM to OpenAIRE.

**Organizations** An Organization is described in the SESAM database as follows:

PROJECT_REF	NOT NULL	NUMBER (10)
LEGAL_SHORT_NAME		VARCHAR2 (4000)
LEGAL_NAME		VARCHAR2 (4000)
LEGAL_STATUS		VARCHAR2 (4000)
CECID		NUMBER
SME		VARCHAR2 (4000)
CTRY_CODE		VARCHAR2 (20)
CTRY_DESCR		VARCHAR2 (200)
THIRDCTRY		VARCHAR2 (1)
CPF_PARTICIPANT_NO		NUMBER
COST_MODEL		VARCHAR2 (20)
COST_MODEL_DESCR		VARCHAR2 (200)
TOTAL_COST		VARCHAR2 (4000)
TOTAL_FUNDING		VARCHAR2 (4000)
FUNDING_MAX_NOE		VARCHAR2 (4000)
TOTAL_COST18		VARCHAR2 (4000)
TOTAL_FUNDING18		VARCHAR2 (4000)
SEAR_MAL		VARCHAR2 (4000)



SEAR\_FEM  
DOCT\_MAL  
DOCT\_FEM



VARCHAR2 (4000)  
VARCHAR2 (4000)  
VARCHAR2 (4000)

Of such descriptive fields, we consider in OpenAIRE the following: *legal short name*, *legal name*, *legal status*.<sup>1</sup> To these properties, we add further ones, relevant to search, statistics and visualization: an *URL of its web site*, an *URL of the logo*, a *country of origin*, and a *location* (longitude, latitude, and time zone). An Organization is associated with (i) the set of Data Sources, if any, of which the Organization is responsible for and (ii) if any, with the Persons authoring a Result while being affiliated with such Organization. An Organization may be the Participant of one or more Projects.

**Persons** A Person is described by a *name*, a *surname*, a *nationality*, and a *gender*. A Person is associated to zero (if not an author of Results in OpenAIRE) or more Authorships and may be the Participant of one or more Projects.

**Participants** A Participant is either a Participant Person or a Participant Organization which is benefiting from the funding of a project. It is described by an *EC participant number*, provided the first time the participant is granted a project by the EC, and to be used in future participation to projects. A Participant is always associated to one or more Projects.

**Authorships** An Authorship describes the fact that one Person has authored one Result while being affiliated with one Organization, if available. As such, an Authorship is associated to the Person, to the Result and, possibly, to the Organization.

**Instances** An Instance represents the combination of the Web Resources associated with a Result and the Data Source where such Web Resources are stored. As such, an Instance is described by the *original unique identifier* (e.g. DOI if the Result is a publication) of the Result at the original Data Source. Furthermore, an Instance is associated with the Result, to the relative Data Source and, if any, to a list of Web Resources. Constraint: if the Result is of kind "publication", the Instance must be associated to at least one Web Resource, relative to the payloads of the publication.

**Data Sources** A Data Source is described by an *official name*, an *English name*, a *URL of its web site*, an *URL of the logo*, an email of the technical contact, a *typology of source* (e.g., OAI-PMH, OAI-ORE, FTP, Web Site), *location* (longitude, latitude and time zone), and an *access information package*, i.e., an XML text bearing the information needed to access the specific source: for example, an OAI-PMH URL, if source typology is OAI-PMH, or an FTP address with login and password, if source typology is FTP. A Data Source is associated with (i) the set of Instances relative to Results available to OpenAIRE that are located at the Data Source and (ii) with the set of Organizations responsible for the Data Source.

**Web Resources** A Web Resource is described by *its unique URL*. A Web Resource is associated to the Instance of the Result of which it contains relevant content.

**Projects** A Project is described in the SESAM database as follows:

ST\_TEL  
ST\_FAX  
ST\_EMAIL

VARCHAR2 (255)  
VARCHAR2 (255)  
VARCHAR2 (255)

<sup>1</sup> We are still waiting for a document from SESAM describing the meaning of all table attributes. So for example, we can only assume that CPF\_participant is the unique identifier assigned by the EC to organizations.

ST_WEB_SITE		VARCHAR2 (255)
DA_PERIOD_END		DATE
DA_PERIOD_START		DATE
ST_FUNDING_SCHEME	NOT NULL	VARCHAR2 (255)
ST_GRANT_AGREEMENT_NUMBER		VARCHAR2 (255)
ST_PROJECT_ACRONYM		VARCHAR2 (255)
ST_PROJECT_TITLE		VARCHAR2 (255)
DA_PROJECT_START_DATE		DATE
DA_PROJECT_END_DATE		DATE
ST_SCI_REP		VARCHAR2 (255)
ST_COORDINATOR_ORGANISATION		VARCHAR2 (255)
ST_COORDINATOR_NAME		VARCHAR2 (255)

Of such descriptive fields, we consider in OpenAIRE the following: *web site*, *EC project web site* (e.g., the project page at CORDIS), *grant\_agreement\_number*, *acronym*, *title*, *start\_date* and *end\_date*. To these, we add a *project call identifier* and a list of *Keywords*. A Project is associated with (i) the set of Participants participating to the Project, (ii) the set of Results whose research was co-funded by the Project, (iii) the Contract Types (e.g., Coordination and Supporting Action, I3) it conforms to, and (iv) the Subdivisions (e.g., FP7-ENERGY, FP7-INFRASTRUCTURES) under which it is funded.

### ***EC Projects classification***

The entities described below mirror the CORDIS classification scheme (see Figure 1), accessible at the address [http://cordis.europa.eu/fp7/info-programmes\\_en.html](http://cordis.europa.eu/fp7/info-programmes_en.html).

A **Funding Programme** is an EC funding programme (e.g., FP7), here characterized by an *identifier*, a *name* and an *acronym* as provided by the EC. A Funding Programme is associated with one or more **Specific Programmes** (e.g., Capacities, Cooperation, Ideas and People, JRC), here characterized by an *identifier*, a *name* and an *acronym*. In turn, a Specific Programme is associated with one or more **Subdivisions** (e.g., research infrastructures), here characterized by an *identifier*, a *name* and an *acronym*. Finally, **Contract Types** (e.g., supporting and coordination actions, eContentPlus, I3, IP), characterized by an *identifier*, a *name* and an *acronym*, are associated to zero, one or more Projects.

- *Funding Programmes*: e.g., FP7
  - *Specific Programmes*: e.g., Capacities, Cooperation, Ideas, People, Euratom direct, Euratom indirect
    - *Subdivisions*: e.g., research infrastructures, Energy, Initial Training

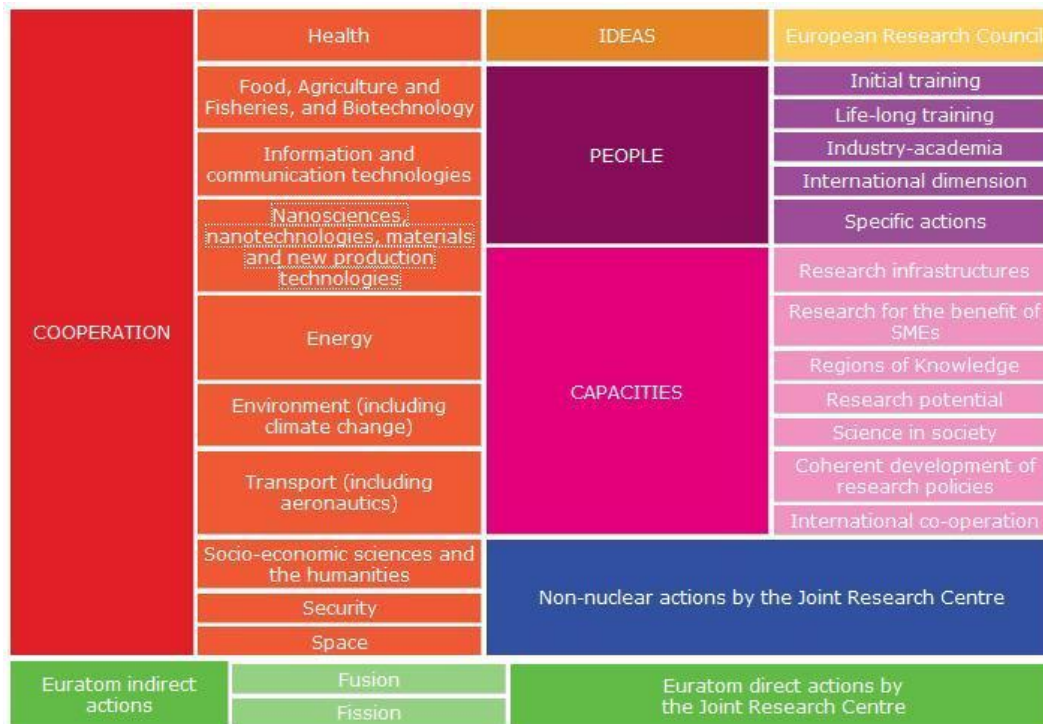


Figure 1 - CORDIS Classification Scheme: Specific programmes

## 2.2 Entity-Relationship model

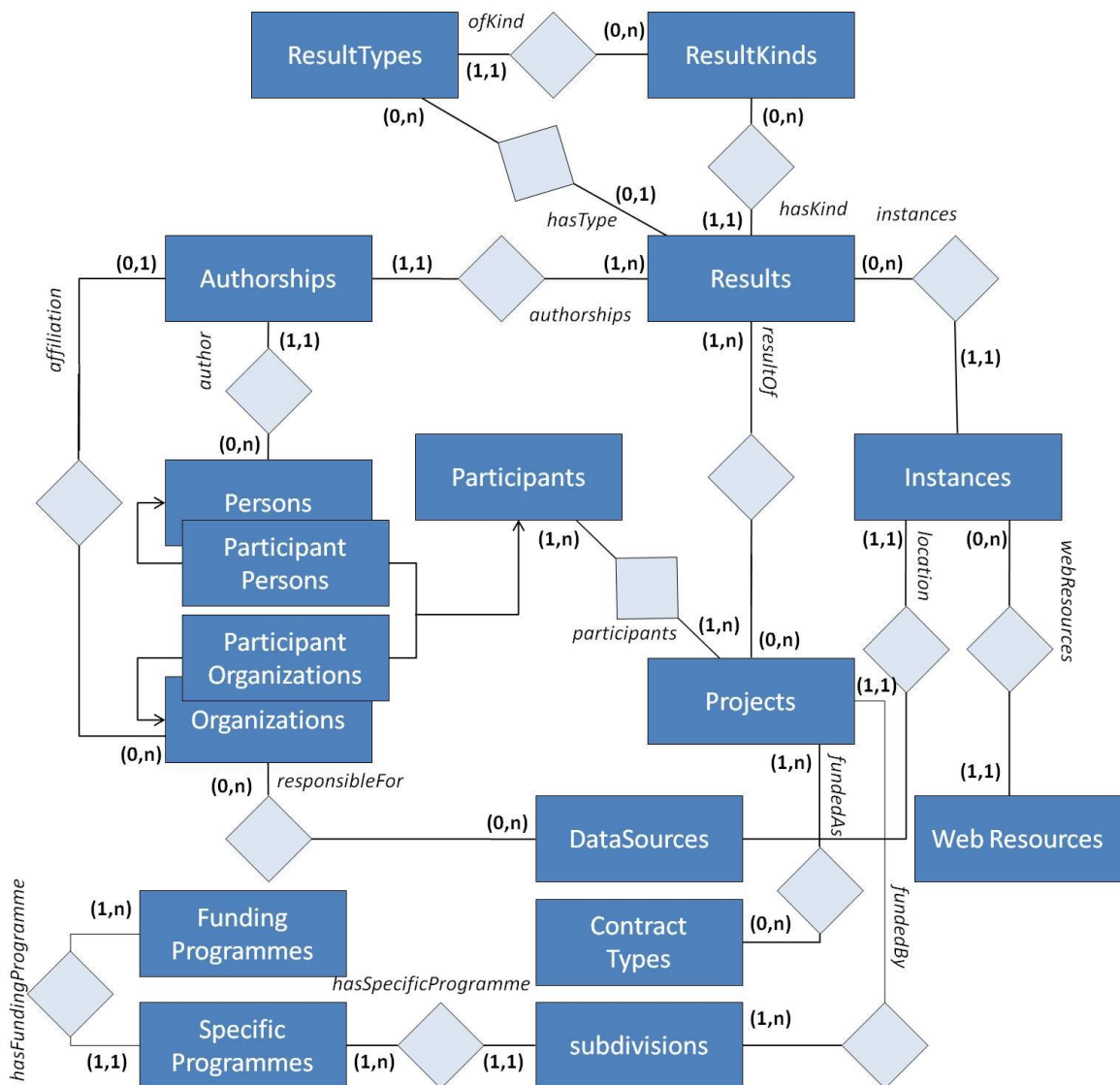


Figure 2 – OpenAIRE Entity-Relationship model

Table 1 reports the attributes and the named associations for the classes in the E-R schema in Figure 2.

Table 1 – E-R Schema: class properties

<p><b>Persons</b></p> <ul style="list-style-type: none"> <li>Name</li> <li>surname</li> <li>Nationality</li> <li>→ author<sup>-1</sup> (0 or N Authorships)</li> </ul> <p><i>Has specializations:</i></p> <ul style="list-style-type: none"> <li><b>Participant Persons</b></li> </ul>	<p><b>Authorships</b></p> <ul style="list-style-type: none"> <li>→ authorships<sup>-1</sup> (1 Results)</li> <li>→ author (1 Persons)</li> <li>→ affiliation (0 or 1 Organizations)</li> </ul>	<p><b>Instances</b></p> <ul style="list-style-type: none"> <li>Unique_identifier (URI)</li> <li>→ instances<sup>-1</sup> (1 Results)</li> <li>→ location (1 Data Source)</li> <li>→ webResources (0 or N Web Resources)</li> </ul>
<p><b>Organizations</b></p> <ul style="list-style-type: none"> <li>Legal short name</li> <li>Legal name</li> <li>Legal status</li> <li>Web site URL</li> <li>Logo URL</li> </ul>	<p><b>Results</b></p> <ul style="list-style-type: none"> <li>Title</li> <li>Date of publication (optional)</li> <li>Description</li> <li>Publisher (optional)</li> <li>Language</li> </ul>	<p><b>Projects</b></p> <ul style="list-style-type: none"> <li>Web site</li> <li>Grant_agreement_number</li> <li>Call_identifier (optional)</li> <li>Acronym</li> <li>Title</li> </ul>

<ul style="list-style-type: none"> <li>Country of origin</li> <li>Longitude, Latitude, Time zone</li> <li>→ affiliation<sup>-1</sup> (0 or N Authorships)</li> <li>→ responsibleFor (0 or N Data Sources)</li> </ul> <p><i>Has specializations:</i></p> <ul style="list-style-type: none"> <li><b>Participant Organizations</b></li> </ul>	<ul style="list-style-type: none"> <li>Access mode</li> <li>Embargo end-date (optional)</li> <li>Keywords</li> <li>→ hasKind (1 ResultKinds)</li> <li>→ hasType (1 ResultType) (optional)</li> <li>→ authorships (1 or N Authorships)</li> <li>→ instances (0 or N Instances)</li> <li>→ resultOf (1 or N Projects)</li> </ul> <p><i>Constraints:</i></p> <ul style="list-style-type: none"> <li>The <i>hasType</i> entity, if specified, must be in the association <i>relatedTypes</i> of the entity <i>hasKind</i></li> </ul>	<ul style="list-style-type: none"> <li>Start_date</li> <li>End_date</li> <li>Keywords</li> <li>→ participants (1 or N Participants)</li> <li>→ resultOf<sup>1</sup> (0 or N Results)</li> <li>→ fundedBy (1 Subdivisions)</li> <li>→ fundedAs (1 or N Contract Types)</li> </ul> <p><i>Derived properties:</i></p> <ul style="list-style-type: none"> <li>Duration (from start_date and end_date)</li> <li>EC_Project_Website(from grant_agreement_number)</li> </ul>
<p><b>Participants</b></p> <ul style="list-style-type: none"> <li>EC_participant_number</li> <li>→ participants<sup>-1</sup> (1 or N Projects)</li> </ul> <p><i>Has specializations:</i></p> <ul style="list-style-type: none"> <li><b>Participant Organizations</b></li> <li><b>Participant Persons</b></li> </ul>	<p><b>Data Sources</b></p> <ul style="list-style-type: none"> <li>Official name</li> <li>English name (optional)</li> <li>Web Site URL</li> <li>Logo URL</li> <li>Contact email</li> <li>Longitude, Latitude, Time zone</li> <li>Typology (e.g., OAI-PMH, OAI-ORE, FTP)</li> <li>Access Info Package (XML)</li> <li>→ location<sup>-1</sup> (0 or N Instances)</li> <li>→ responsibleFor<sup>-1</sup> (0 or N Organizations)</li> </ul>	<p><b>Web Resources</b></p> <ul style="list-style-type: none"> <li>Web Resource URL</li> <li>→ webResources<sup>-1</sup> (1 Instances)</li> </ul>
<p><b>Funding Programmes</b></p> <ul style="list-style-type: none"> <li>Identifier</li> <li>Name</li> <li>Acronym</li> <li>→ hasFundingProgramme<sup>-1</sup> (1 or N Funding Programmes)</li> </ul>	<p><b>Specific Programmes</b></p> <ul style="list-style-type: none"> <li>Identifier</li> <li>Name</li> <li>Acronym</li> <li>→ hasSpecificProgramme<sup>-1</sup> (1 or N Subdivisions)</li> <li>→ hasFundingProgramme (1 Funding Programme)</li> </ul>	<p><b>Subdivisions</b></p> <ul style="list-style-type: none"> <li>Identifier</li> <li>Name</li> <li>Acronym</li> <li>→ hasSpecificProgramme (1 Specific Programmes)</li> <li>→ fundedBy<sup>-1</sup> (1 or N Projects)</li> </ul>
<p><b>Contract Types</b></p> <ul style="list-style-type: none"> <li>Identifier</li> <li>Name</li> <li>→ fundedAs<sup>-1</sup> (0 or N Projects)</li> </ul>	<p><b>Result Kinds</b></p> <ul style="list-style-type: none"> <li>Name</li> <li>→ ofKind<sup>-1</sup> (0 or N ResultTypes)</li> </ul>	<p><b>Result Types</b></p> <ul style="list-style-type: none"> <li>Name</li> <li>→ ofKind (1 ResultKinds)</li> </ul>

## 2.3 Relational Schema

The relational schema in Figure 3 is obtained as a normalized transformation from the E-R class schema in Figure 2.

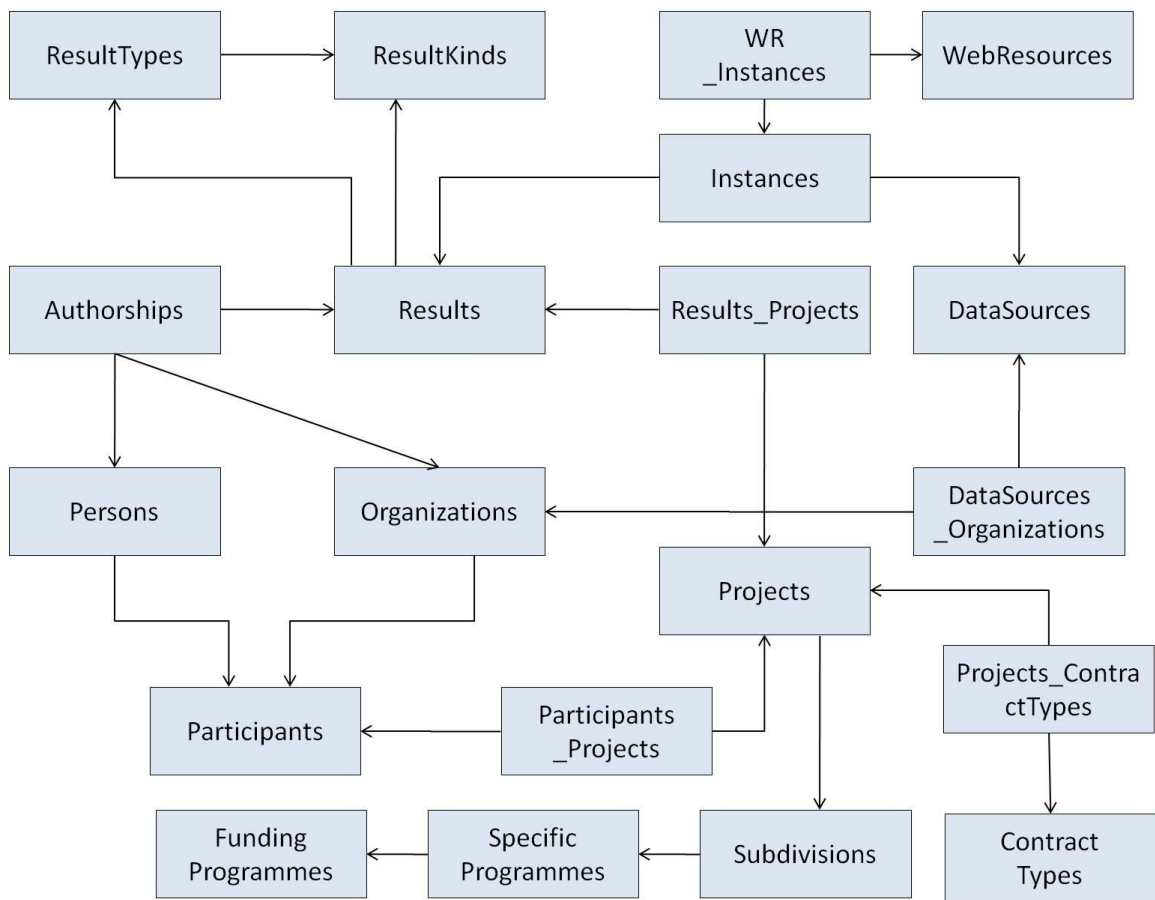


Figure 3 – OpenAIRE Relational Schema

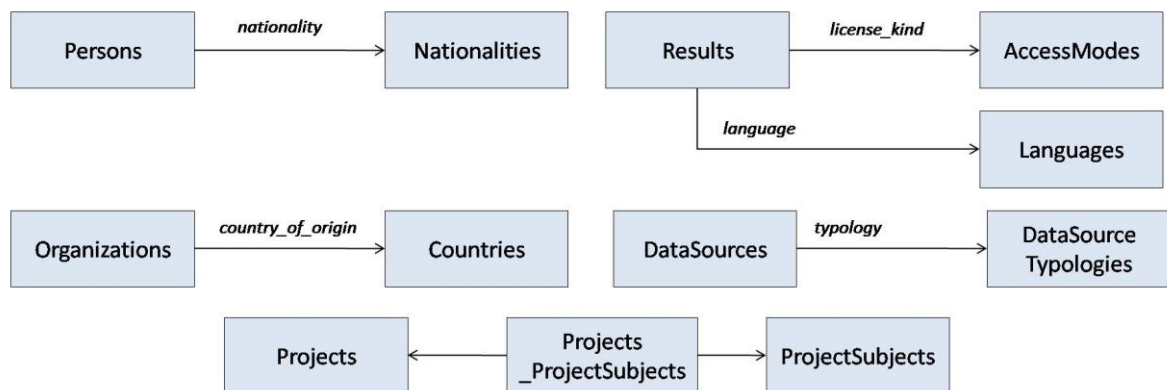


Figure 4 – Controlled Vocabularies Relational Schema: "vocabulary tables"

Figure 4 introduces the tables representing controlled vocabularies in the E-R schema, namely *vocabulary tables*. A vocabulary table has a structure of the form (key, value) is introduced whenever a controlled vocabulary needs to be managed. Currently, vocabularies are: Nationalities (of people), Result Kinds (categories of results, e.g., publication, research data, software product), Result Types (dependent on ResultKinds, for publication for example: article, book, manual, etc), Countries and Languages (possibly taken from ISO standard), Access Modes (described in the introduction), Data Source typologies (e.g., institutional repository, web site) and Project Subjects (taken from an authoritative list of subjects, to be provided by the EU: for example the seven research areas of the Pilot).



It is to be considered whether “vocabulary tables” should referred to from other tables by only using the primary key or by using the pair (key,value) in order to optimize searches (*tech*: no join query needed) and enable full-text search on values.

Finally, the “vocabulary tables” defined in this document may be subject to changes in the future: other attributes may be added to such tables, thereby extending the notion of vocabulary to that of an authority file, or new vocabularies be added (e.g., publisher). The trade-off stands in the balance between having clean and precise data and the effort required to keep a consistent and meaningful authoritative vocabulary of terms.

Table 2 – Relational schema: tables structure

<p><b>Persons</b></p> <ul style="list-style-type: none"> <li>• personID: PK</li> <li>• name: string</li> <li>• surname: string</li> <li>• nationality: FK Nationality(nationalityID) (optional)</li> <li>• participantID: FK Participants(participantID) (optional)</li> </ul>	<p><b>Authorships</b></p> <ul style="list-style-type: none"> <li>• authorshipID: PK</li> <li>• result: FK Results(resourceID)</li> <li>• author: FK Persons(personID)</li> <li>• affiliation: FK Organizations(organizationID) (optional)</li> </ul>	<p><b>Instances</b></p> <ul style="list-style-type: none"> <li>• instanceID: PK</li> <li>• unique_identifier: String</li> <li>• datasource: FK datasources(DataSourceID)</li> <li>• resource: FK Results(ResultID)</li> </ul>
<p><b>Web Resources</b></p> <ul style="list-style-type: none"> <li>• webResourceID: PK</li> <li>• file_URL: “URL” String</li> </ul>	<p><b>Web Resources_Instances</b></p> <ul style="list-style-type: none"> <li>• webResourceID: PK FK Web Resources(fileID)</li> <li>• instanceID: PK FK Instances(manifestationID)</li> </ul>	<p><b>Organizations</b></p> <ul style="list-style-type: none"> <li>• organizationID: PK</li> <li>• legal_short_name: String</li> <li>• legal_name: String</li> <li>• legal_status: String</li> <li>• web_site_URL: String</li> <li>• logo_URL: String (optional)</li> <li>• country_of_origin: FK Countries(countryID)</li> <li>• participantID: FK Participants(participantID) (optional)</li> </ul>
<p><b>Results</b></p> <ul style="list-style-type: none"> <li>• resultID: PK</li> <li>• title: String</li> <li>• publication_date: Date (optional)</li> <li>• description: String</li> <li>• publisher: String (optional)</li> <li>• hasKind: FK ResultKinds(resourceKindID)</li> <li>• hasType: FK ResultTypes(resourceTypeID) (optional)</li> <li>• language: FK Languages(languageID)</li> <li>• access_mode: FK Access_Modes(accessModeID)</li> <li>• embargo_end_date: Date</li> <li>• keywords: String</li> </ul>	<p><b>Projects</b></p> <ul style="list-style-type: none"> <li>• projectID: PK</li> <li>• web_site: “URL” String</li> <li>• EC_project_website: “URL” String</li> <li>• grant_agreement_number: String</li> <li>• call_identifier: String (optional)</li> <li>• acronym: String</li> <li>• title: String</li> <li>• start_date: Date</li> <li>• end_date: Date</li> <li>• fundedBy: FK Subdivisions(subdivisionID)</li> </ul>	<p><b>Participants</b></p> <ul style="list-style-type: none"> <li>• beneficiaryID: PK</li> <li>• EC_participant_number: String (unique)</li> </ul>
<p><b>DataSources</b></p> <ul style="list-style-type: none"> <li>• datasourceID: PK</li> <li>• official_name: String</li> <li>• English_name: String (optional)</li> <li>• web_site_URL: String</li> <li>• logo_URL: String</li> <li>• contact_email: “email” String</li> </ul>	<p><b>Participants_Projects</b></p> <ul style="list-style-type: none"> <li>• participant: PK FK Participants(participantID)</li> <li>• project: PK FK Projects(projectID)</li> </ul>	<p><b>DataSources_Organizations</b></p> <ul style="list-style-type: none"> <li>• datasource: PK FK DataSources(datasourceID)</li> <li>• organization: PK FK Organizations(organizationID)</li> </ul>

<ul style="list-style-type: none"> <li>• longitude: Number</li> <li>• latitude: Number</li> <li>• time zone: String</li> <li>• typology: FK DataSourceTypologies(datasourceTypologyID)</li> <li>• access_info_package: "XML" String</li> </ul>		
<p><b>Results_Projects</b></p> <ul style="list-style-type: none"> <li>• result: PK FK Results(resultID)</li> <li>• project: PK FK Projects(projectID)</li> </ul>	<p><b>FundingProgrammes</b></p> <ul style="list-style-type: none"> <li>• fundingProgrammeID: PK String</li> <li>• programme_name: String</li> <li>• programme_acronym: String</li> </ul>	<p><b>SpecificProgrammes</b></p> <ul style="list-style-type: none"> <li>• specificProgrammeID: PK String</li> <li>• specificProgramme_name: string</li> <li>• specificProgramme_acronym: String</li> <li>• hasFundingProgramme: FK FundingProgrammes(fundingProgrammeID)</li> </ul>
<p><b>Subdivisions</b></p> <ul style="list-style-type: none"> <li>• subdivisionID: PK String</li> <li>• subdivision_name: string</li> <li>• subdivision_acronym: String</li> <li>• hasSpecificProgramme: FK SpecificProgrammes(specificProgrammeID)</li> </ul>	<p><b>ContractTypes</b></p> <ul style="list-style-type: none"> <li>• contractTypeID: PK String</li> <li>• contractType_name: string</li> </ul>	<p><b>Projects_ContractTypes</b></p> <ul style="list-style-type: none"> <li>• project: PK FK Projects(projectID)</li> <li>• contractType: PK FK ContractTypes(contractTypeID)</li> </ul>
<p><b>ResultTypes</b></p> <ul style="list-style-type: none"> <li>• resultTypeID: PK String</li> <li>• name: String</li> <li>• ofKind: FK ResultKinds(resultKindID)</li> </ul>	<p><b>ResultKinds</b></p> <ul style="list-style-type: none"> <li>• resultKindID: PK String</li> <li>• name: String</li> </ul>	<p><b>DataSourceTypologies</b></p> <ul style="list-style-type: none"> <li>• datasourceTypologyID: PK String</li> <li>• name: String</li> </ul>
<p><b>ProjectSubjects</b></p> <ul style="list-style-type: none"> <li>• projectSubjectID: PK String</li> <li>• name: String</li> </ul>	<p><b>Projects_ProjectSubjects</b></p> <ul style="list-style-type: none"> <li>• project: PK FK Projects(projectID)</li> <li>• project_subject: PK FK Project_Categories(projectCategoryID)</li> </ul>	<p><b>Nationalities</b></p> <ul style="list-style-type: none"> <li>• nationalityID: PK String</li> <li>• name: String</li> </ul>
<p><b>Countries</b></p> <ul style="list-style-type: none"> <li>• countryID: PK String</li> <li>• name: String</li> </ul>	<p><b>AccessModes</b></p> <ul style="list-style-type: none"> <li>• accessModeID: PK String</li> <li>• name: String</li> </ul>	<p><b>Languages</b></p> <ul style="list-style-type: none"> <li>• languageID: PK String</li> <li>• name: String</li> </ul>

### 3 Implications for the definition of OpenAIRE Guidelines

This section analyzes the OpenAIRE data model and identifies some of the requirements the model may impose to OpenAIRE repository managers. As such, it will serve as input to the definition of OpenAIRE Guidelines for Repository Managers.

Based on the OpenAIRE data model and on the OpenAIRE system requirements, in order to facilitate the automatic ingestion process of metadata records relative to publications from the repositories to OpenAIRE, the following requirements were identified:

- 1) The repository should implement an OAI-PMH interface (OpenAIRE provides a suite for managing harvesting and aggregation of a federation of OAI-PMH conformant repositories);
- 2) The repository should export an OAI-PMH set (possibly with an OpenAIRE-flavored name, such as "EC\_Project\_Results") whose metadata records correspond to the publication results of interest to OpenAIRE; if such a set is not available, an alternative way of selecting records of interest to OpenAIRE should be provided (e.g., a special metadata field allowing for direct or indirect identification of such records)
- 3) Information about funding projects and licenses of the publications should come along with the bibliographic metadata: repository managers should find ways to "enrich" their Dublin Core records with new fields, in particular:
  - *access\_mode*, mandatory (possibly values from the OpenAIRE controlled vocabulary "AccessKinds")
  - *embargo\_end\_date*, mandatory when *access\_mode* = "embargo"
  - *project*, mandatory sequence (values from the OpenAIRE vocabulary obtained by projection of the attribute *projectID* of the table Projects)
  - *unique identifier*, optional (the unique identifier of the publication at its original repository, different from the value expected in *dc:identifier*, see below)
- 4) Some of the Dublin Core fields should be mandatory or optional based on their correspondence with the OpenAIRE properties for Results of typology "publication". Similarly, values of the DC fields should possibly, not mandatorily, adhere to the OpenAIRE vocabularies to which the corresponding properties adhere. Alternatively, if values of such fields are provided according to different vocabularies, a mapping onto the matching OpenAIRE vocabulary values should be provided. Some examples are:
  - *dc:title*, mandatory
  - *dc:creator*, mandatory sequence
  - *dc:description*, mandatory
  - *dc:identifier*, mandatory sequence (the URL of the file, in the style of DRIVER Guidelines)
  - *dc:subject*, optional sequence
  - *dc:publisher*, optional

- *dc:language, optional (possibly values from the OpenAIRE controlled vocabulary "Languages")*
- *dc:type, optional (possibly values from the OpenAIRE controlled vocabulary "ResourceTypes" of ResourceKind="Publication")*

## 4 Collaborative refinement of the Information Space

The OpenAIRE Web Site enables registered users to access and modify the Information Space in order to add, remove clean or enrich its content. The Information Space structure is conceived to enable a growing process of data quality improvement, so as to construct and grow authority files for persons, results, organizations, data sources and projects. This process is based on trust of registered users, whose corrective actions are double-checked by administrators, and supported by automatic tools for identification of duplicates.

In this scenario, actions from registered users should be recorded and limited to those that are reversible by administrators; actions from administrators should also be recorded, but can be irreversible. Different technical solutions are possible and will be discussed at the time of implementing the Information Space.