

Technical Integration of Data Repositories: status and challenges

28/01/2020

Isabel Bernal (DIGITAL.CSIC)

Slava Tykhonov (DANS-KNAW)

Fernando Aguilar (CSIC)



Outline



- Repositories
- Previous Integration initiatives
- Technical barriers and challenges
- Potential integration

Repositories: DIGITAL.CSIC enabling Open Science

- CSIC Green Open Access/Open Data Mandate since April 2019
- Repository's Data Policy since 2014: nearly 12,000 datasets, training to CSIC researchers and librarians, support to comply with funders and journals data sharing policies
- Involvement in several EOSC / FAIR Data initiatives
- DOI assignation to datasets, software, notebooks, preprints and other outputs via DataCite
- Member of DataCite Metadata Working Group
- Current aggregation with OpenAIRE, CORE, SHARE, DataCite Search, BASE...
- Registered in Re3Data and Repository Finder

Welcome to DIGITAL.CSIC, the institutional repository of the Spanish National Research Council.
 DIGITAL.CSIC organizes, preserves and provides open access to CSIC research outputs.
[DIGITAL.CSIC Annual Reports](#)
[CSIC Mandate FAQs](#)

DIGITAL.CSIC

OPEN SCIENCE

Share your Open Access story
 Send us your works

BIBLIOTECAS CSIC

POR LA CIENCIA ABIERTA

AGENDA GLOBAL 	AGENDA EUROPEA
AGENDA NACIONAL 	AGENDA INSTITUCIONAL
COMPETENCIAS PROFESIONALES 	INFRAESTRUCTURAS
NUEVOS SERVICIOS 	AGENTES DE TRANSFORMACIÓN

<http://hdl.handle.net/10261/19954>

Media Gallery Your research in images

Highlights

- 5th Conference of CSIC Libraries and Archives Network [11/12/2019]
 Last November 28-29 CSIC Libraries and Archives Network held its 5th Conference in Madrid. Under the motto "CSIC Libraries and Archives for Open Science: present and future" a wide range of topics were lively discussed, in particular what role CSIC libraries and archives are playing in the shift to an Open Science paradigm and what is on the horizon. A special collection gathers all presentations and related materials.
- New educational resources [07/11/2019]
 In October we participated in two training workshops for CSIC researcher and technical consultants. Impact and measurement of Open Access analyzes the volume of CSIC and global scientific publications that are already made available in Open Access within a context of emerging evaluation systems and several strategies to accelerate the transition to a new scholarly communications model. For its part, Open Science: challenges and solutions explores the pillars of Open Science framework as well as opportunities for CSIC.
- New features for DIGITAL.CSIC users [20/09/2019]
 On September 18 DIGITAL.CSIC migrated from ORACLE to PostgreSQL. We have seized the opportunity to make improvements on browsing and search functionalities. In the new version, users can browse and search contents based on a greater variety of filters and it is also possible to limit searches to a specific entry (items, profiles, institutes/collections). A new module with access to all usage statistics since 2008 is also available.

[All News](#) RSS

!

GOOD PRACTICES, RESOURCES AND SUPPORT

OA

OPEN ACCESS MANDATES

DIGITAL.CSIC in figures

- 189.041 records available
- 61,32 % open access
- 153 Institutes and 1.383 Collections

www.eosc-synergy.eu

DSpace-CRIS state of art

- Publications and datasets (organized in communities and collections) supported in standard DSpace
- DSpace-CRIS extends this functionality and involves the other entities that are part of the research landscape:
 - Researchers
 - Projects
 - Organization Units (Groups, Departments)
 - Second Level Dynamic Objects
- DSpace 5.x has useful software integrations like CKAN and ORCID and data integrations (Archivematica) and can be deployed from Docker images

SSHOC Dataverse

Makes use of Dataverse software developed by Harvard IQSS

4 ERICs: DARIAH, CLARIN, EHRIS and CESSDA

Building mature infrastructure based on requirements of involved EOSC communities (Docker and Kubernetes)

Investigating sustainable governance models

Training Service Providers and institutes how to use Dataverse as a service

Different levels of repositories integration

Metadata integration

- aggregation by OAI-PMH and ResourceSync

Software integration

- data repository provides functionality delivered by another services (DSpace-CRIS with CKAN and ORCID support)

Data integration

- controlled vocabularies support (COAR, OpenAire, LOC, FundRef)

Data archiving

- Datasets with files can be moved from one repository to another by SWORD protocol (for example, from DSpace-CRIS to the long term archive like Archivematica)

Previous Integration initiatives

- OAI-PMH, Resource sync (THIS ONE IS EMERGING, TRYING TO REPLACE OAI)
- Aggregators
- OpenAIRE
- ...

Previous Integration initiatives - Metadata

OAI-PMH

Open Archives Initiative Protocol for Metadata Harvesting

6 Actions. Any metadata schema (pre-defined). Simple but enough.

Used by aggregators. Information about the resource, but not the resource itself. (there is a further development which is OAI-ORE)

Previous Integration initiatives - Metadata

ResourceSync

ResourceSync supports synchronization of both Resources and Metadata about Resources with the relationships clearly indicated.

Tracks the status of a resource (updated, deleted).

Access to the resource.

Aggregators - Integration of repos

OpenAIRE explorer

CORE

Google Dataset Search

DataCite Search

SHARE



OpenAIRE has grown through a series of project phases funded by the European Commission: from the DRIVER projects to link Europe's repository infrastructure, to the first OpenAIRE project aimed to assist the EC in implementing its initial pilot for Open Access (OA) to publications.

Services:

- Explorer (aggregator)
- Monitor (compliance with EC Open Access/Open Data Mandate)
- Scholix - Scholexplorer (links papers-datasets)
- AMNESIA
- DMPs

WP3.3 FAIR data integration



- Thematic services usually keeping data inside of its own database, not storing in some data repository with persistence identifiers. Data derivatives often located in the temporary storage, aren't sustainable and can be cleaned (Findable)
- Metadata schema varies in the different scientific communities, should be supported by all data repositories, however citation block should be common (Accessible)
- controlled vocabularies (CV)/ontologies are different, should be standardized and preferably provided by common CV services maintained by EOSC (Interoperable)
- most of metadata contain descriptions in the native language of researcher, should be translated to English and linked to common CV (Interoperable)
- most of Thematic services don't have provenance information exposed as PROV-O or other standard required to be stored in datasets together with metadata and data files (Reusable)
- different types of licenses for data access should be supported by all EOSC data repositories, sensitive data should be handled differently (Reusable)

WP3 Software integration challenges

- The maintenance of the distributed applications with external services is very difficult and expensive
- requires the highest level of service maturity
- increasing the **code coverage** does not necessarily lead to more **functionality coverage**
- writing integration tests even more important than adding more unit tests
- it's almost not possible to run distributed services without help from community

We need Docker to increase the maturity of the infrastructure

Docker advantages

- Faster development and deployments
- Isolation of running containers allows to scale up apps
- Portability saves time to run the same image on the local computer or in the cloud
- Snapshotting allows to archive Docker images state
- Resource limitation can be adjusted
- Increasing reproducibility

SQA Service maturity requires Docker on Kubernetes!

Still, there are some technical barriers...

- The maturity of software and services is different, common baseline in the development state (WP3)
- Maintenance of the ecosystem consisting of variety of tools and services is complicated and consumes a lot of human resources (upgrade of servers, bug fixing, security updates)
- fast technological development requires continuous training and knowledge transfer

Service for Quality Assurance



Selenium IDE - Dataverse*

Project: Dataverse*

Executing ▾

Search tests*

	Command	Target	Value
1	open	/	
2	set window size	1680x962	
3	click	linkText=Arts and Humanities	
4	click	id=j_idt414:searchBasic	
5	type	id=j_idt414:searchBasic	news
6	send keys	id=j_idt414:searchBasic	\$(KEY_ENTER)

Command //

Target

Value

Description

Runs: 1 Failures: 1

Log	Reference	
3. click on linkText=Arts and Humanities OK		16:23:57
4. click on id=j_idt414:searchBasic OK		16:23:59
5. type on id=j_idt414:searchBasic with value news OK		16:24:04
6. sendKeys on id=j_idt414:searchBasic with value \$(KEY_ENTER) Failed: { "code": -32000, "message": "DOM Error while querying" }		16:24:06

Selenium IDE allows to create and replay all UI tests in your browser

Shared tests can be reused by community to increase reproducibility

SQA for the service maturity = unit tests + integration tests

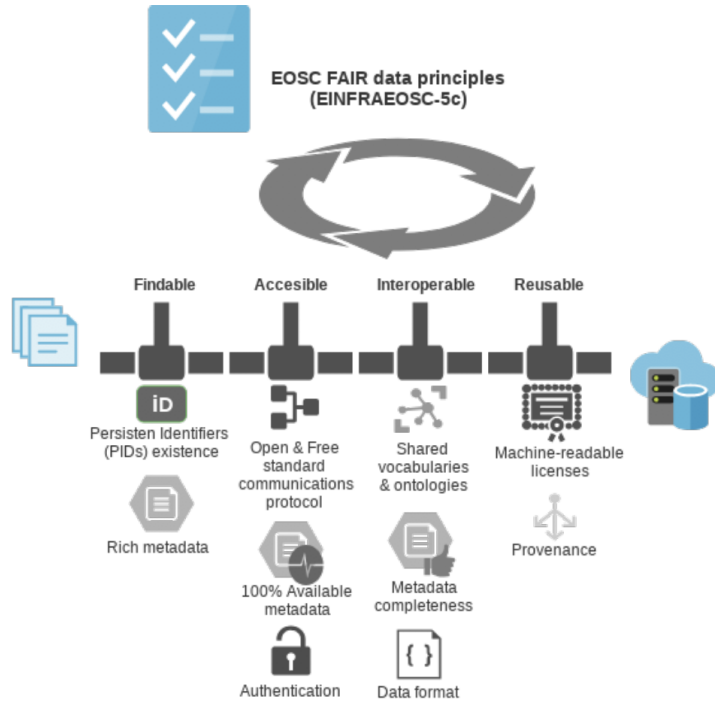
WP3.3 FAIR baseline on the technical integration

Integration of data repositories in EOSC means:

- Findable: data have to be findable in a standard way by the users and other core and thematic services
- Accessible: data accessible via standard interfaces supported in EOSC
- Interoperable: data combinable with other data in EOSC repositories
- Reusable: data exploitable by the EOSC services enabling data analysis to be performed using resources and services from the EOSC infrastructure providers, keeping datasets versions and provenance.

(from EOSC-Synergy proposal)

WP3.3 FAIR adoption



Evaluation of FAIRness of data:

- check if PID exists (F)
- metadata accessible via standard protocols (A)
- shared vocabularies used (I)
- machine readable data contain provenance information (R)

Outcome: “EOSC-ready” badge should be released by SQA service for all other services suitable to be delivered under EOSC

WP3.1 Maturity of thematic services

Thematic services should:

- follow common SQA baseline for Software and Services
- increase the maturity by adding unit and integration tests, preferably by community
- improving Cloud infrastructure both horizontally and vertically requires a good testing strategy

Potential Integration

~~Integration of repositories metadata — data~~

Enrich repositories services



Focused on:

- FAIRness
- Metadata and data formats flexibility
- Connection to software
- Connection to computing

Examples - FAIRness evaluation

WP3 - WP4

Automatic or semi-automatic FAIRness evaluation

Based on recommendations (e.g. FAIRsFAIR)

Some solutions available: #FAIRevaluator



Import MI Tests

Import Maturity Indicators Tests as YAML [smartAPI](#) interface annotation

Get started



Create collections

Assemble Maturity Indicators Tests into community centered collections

Get started



Evaluate resources

Evaluate resources FAIRness against Collections of Maturity Indicator Tests

Get started

Example - Jupyter integration

- Search and retrieve data/software from repositories.
- Ensure reproducibility.
- Workflows
- Publication of results.
- Etc.