



## D5.1 FAIR Research Data Management Tool Set

<b>Lead Partner:</b>	CINES
<b>Authors</b>	N. Cazenave (CINES), L. Candela (CNR-ISTI), L. Berberi (KIT), Jos van Wezel (KIT), A. Hashibon (Fraunhofer), Yann Le Franc (CINES)
<b>Version:</b>	2.0
<b>Status:</b>	Final version
<b>Dissemination Level:</b>	Public
<b>Document Link:</b>	Portal users: <a href="https://repository.eosc-pillar.eu/index.php/f/35824">https://repository.eosc-pillar.eu/index.php/f/35824</a> Public URL (PDF): <a href="https://repository.eosc-pillar.eu/index.php/s/MBgkN5kexKMbMoq">https://repository.eosc-pillar.eu/index.php/s/MBgkN5kexKMbMoq</a>

### Deliverable Abstract

This document is accompanying the delivery of the bundle of service instance(s) that are the output of T5.1 and T5.2 activities. It provides a short summary of the work and the list of services and how to access them. The tool-set aims at offering solutions for Research Data Management promoting the implementation of FAIR principles and practices. This tool-set is updated every 6 months after the initial release. A major release is planned in June 2021 (PM 24).



## COPYRIGHT NOTICE



This work by Parties of the EOSC-Pillar is licensed under a Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>). The EOSC-Pillar project is co-funded by the European Union Horizon 2020 programme under grant number 857650.

## DELIVERY SLIP

	<b><i>Name</i></b>	<b><i>Partner/Activity</i></b>	<b><i>Date</i></b>
<b>From:</b>	N. Cazenave L. Candela	CINES CNR-ISTI	
<b>Moderated by:</b>	Y. Le Franc	CINES	
<b>Reviewed by:</b>			
<b>Approved by:</b>			

## TERMINOLOGY

<https://eosc-portal.eu/glossary>

<b><i>Terminology/Acronym</i></b>	<b><i>Definition</i></b>
<b>API</b>	Application Programming Interface
<b>FAIR Data Point</b>	A software enabling the implementation of a metadata repository providing access to metadata according to the FAIR principles.
<b>Federated FAIR Data Space</b>	A unifying data space aggregating datasets from separated data sources and repositories with the aim to give access to them according to the FAIR principles.
<b>FDP</b>	see FAIR Data Point
<b>FFDS</b>	Federated FAIR Data Space
<b>Virtual Research Environment</b>	A web-based working environment conceived to provide a community of practice with services and data of interest;
<b>VRE</b>	see Virtual Research Environment

# Contents

1	Introduction .....	5
2	Implementations and initial results.....	6
2.1	Generic F2DS solution .....	6
2.2	D4Science solution.....	8
3	Access to the bundle of services.....	10
4	Concluding remarks.....	12
	References .....	13

## Executive summary

The goal of the EOSC-Pillar WP5 “The Data layer: establishing FAIR data services at the national and transnational level”, is to create the settings for an effective sharing, exploitation and reuse of data across initiatives and communities partaking to EOSC-Pillar and beyond. To pursue this challenging goal, the project leverages and builds upon results from previous and ongoing projects as well as on the experience of the partners in the project. This combined expertise will provide data providers and data consumers with a dedicated set of services (and accompanying training) supporting the creation of a data space where multiple datasets from separate locations are virtually joined by combining their metadata and are subsequently published in accordance with the FAIR principles (Wilkinson et al, 2016). EOSC-Pillar performs a pilot into the federation of data sources presenting solutions that differ depending on the requirements of the use cases promoted by researchers.



# 1 Introduction

The aim of T5.1 and T5.2 is to provide a set of services for creating a Federated FAIR data space (F2DS). This F2DS should provide tools for data producers to make their data more compliant with the FAIR principles and any other specific policies (T5.1) and to integrate them with other data coming from multiple disciplines that could then be accessed and reused by data consumers through dedicated interfaces (T5.2).

The advantage of such federated FAIR data space as piloted in EOSC-Pillar, which implementation involves development as well as support and collaboration with domain scientists, is manifold. With a F2DS, researchers will be able to search, find and retrieve data using a single access point and tool set. Not only does this save working time because it masks the different access methods of different sources and offers a single access point through user and programming interfaces (UI and API), but given the proper combination of selected search criteria and content of one or more data-sets is explored in unison it can truly deliver new insights.

Prerequisite to the unification of (meta)data sources is on one hand the normalisation through the use of existing common standards, starting from the access method and authentication, through authorisation and metadata and data formats. This normalisation i.e. the adoption of such common standards and format is the only path to improve the FAIRness of the data. Although it is considered the most straightforward approach for unification and shows many results, the extent of the scientific domains, the dynamics of data collection over time and ultimately the rapidly changing information technology makes normalisation a challenge. On the other hand technology is needed that can adapt to the changing and developing data types and at the same time presents a stable and flexible connection to the data sources for automatic parsing and data exploration. This is one of the goals in the implementation of, for example, the [FAIR data point](#)<sup>1</sup> software.

The solutions delivered in WP5 and discussed in this document consist of a variety of tools and services. Some tools are optimised for either solution, some maybe used in another context or are offered as an independent service. Several (components of) the services are able to use existing resource offerings available in EOSC e.g. virtual computing and data services. Eventually many of the WP5 services can be easily deployed using for example the framework<sup>2</sup> developed in the [Deep Hybrid DataCloud](#) project and its predecessors. Most important, the presented solutions offer the best possible adaptations to the specific area i.e. use-case, they are developed for.

The tools presented implement the EOSC-Pillar Federated FAIR Data Space (F2DS), i.e. a unifying **data space** that is built by aggregating and enriching datasets from a set of multidisciplinary repositories, i.e. data sources, with the aim to facilitate data discovery and re-use. Although datasets are the primary focus of the resulting **data space**, other items are managed including repositories and data sources, APIs, metadata schemas and ontologies.

---

<sup>1</sup> <https://github.com/FAIRDataTeam/FAIRDataPoint-Spec>

<sup>2</sup> available in the EOSC marketplace

## 2 Implementations and initial results

The implementation of the envisaged Federated FAIR Data Space concept is leveraging on existing tools and services that have been developed for similar requirements and have matured independently. Because of this it was decided to follow a pragmatic and exploratory-oriented approach for the first release of the tool-set leading to the development and testing of two independent solutions along side. Such an approach promotes discussion and exchange between two possible implementations sharing some commonalities and technologies yet proposing slightly diverse work flows, technical approaches and delivery strategies. This target driven optimisation in two lines of development is to be evaluated at PM24 at which further development of either one of the solutions must be halted for reasons of efficiency. Even without ongoing development the service will function, enhance FAIRness of the presented dataset and continue to be useful for specific use cases. The two proposed solutions are not completely disjoint. Rather, it is likely that during their operation and exploitation they will be assembled into a single entry access.

The two implementations of the F2DS for EOSC-Pillar are the Generic F2DS (g-F2DS) and the D4Science-based F2DS. The D4Science data space is based on an existing service whereas the g-F2DS implementation is built from the ground up by integrating existing services for FAIRifying data and the latest of technologies and protocols. While both offer the researcher a shared data space, there are some fundamental technical differences between the generic F2DS and D4Science based solutions. E.g. the g-F2DS approach uses a common API description and a simple metadata mapping to automatically and intelligently harvest, convert and publish metadata describing data sets in a single format ([DCAT](#)<sup>3</sup>). D4Science is based on the [CKAN](#)<sup>4</sup> technology which is also used as the underlying framework of the EUDAT [B2FIND](#)<sup>5</sup> service and offers the possibility to harvest metadata using various formats ([CSW](#)<sup>6</sup>, [OAI-PMH](#)<sup>7</sup>,...). Both approaches offer capabilities to search datasets either through the rich search interface of D4Science solution, or the [SPARQL](#)<sup>8</sup> engine built-in the g-F2DS.

### 2.1 Generic F2DS solution

The g-F2DS solution has been designed to use existing state-of-the-art tools developed by various EU stakeholders for FAIRifying data and to integrate them together into a coherent, scalable and innovative solution that can be easily deployed on any cloud infrastructure. To achieve this objective, it was therefore chosen to take advantage of container technologies and deployments on [kubernetes](#)<sup>9</sup>.

---

<sup>3</sup> <https://www.w3.org/TR/vocab-dcat-2/>

<sup>4</sup> <https://ckan.org/>

<sup>5</sup> <http://b2find.eudat.eu/>

<sup>6</sup> <https://www.ogc.org/standards/cat>

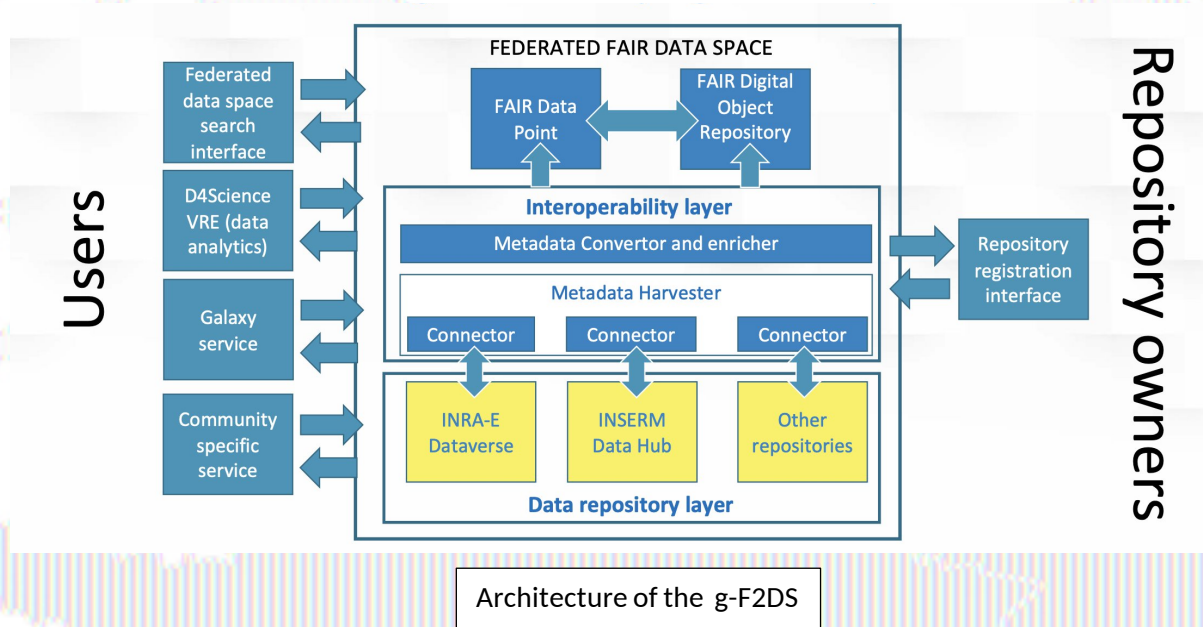
<sup>7</sup> <https://www.openarchives.org/pmh/>

<sup>8</sup> <https://www.w3.org/TR/sparql11-query/>

<sup>9</sup> <https://kubernetes.io>



The solution is based on the Fair Data Point technology<sup>10</sup> (FDP) linked to an access API registry (OpenAPI<sup>11</sup>/smartAPI<sup>12</sup> technology). The generic architecture of this solution is shown below.



The g-F2DS solution includes two key layers: the data repository layer and the interoperability layer and offers two specific interfaces: one for data producers to register their repositories within the F2DS and another interface to access the federated space for searching and using the data.

The first layer is the **Data repository layer** which contains the metadata description of the different repositories (which will be made compliant with the EOSC Portal Service Description Template, or “profile”) as well as technical information to access the repository content and the description of each repository’s API in a common format.

The Data Repository Layer is then connected to the **Interoperability layer** which is based on two components: the *metadata harvester* and the *metadata converter and enricher*. The metadata harvester, also called *smarHarvester*, uses the API descriptions to automatically build a dedicated client and the appropriate queries to gather the metadata stored in each repository: these clients make automatic links and regularly test the state of the metadata (changes, deletions, modifications, additions, ...). Next, the *metadata converter and enricher* parses, converts all the metadata into the DCAT model and serialises it as [RDF](https://www.w3.org/RDF/)<sup>13</sup>. It uses pre-defined elements (catalogs/datasets/distributions) of the FDP. Thus all metadata will be described according to the same data model, which promotes interoperability and reuse of data.

<sup>10</sup> <https://github.com/FAIRDataTeam/FAIRDataPoint-Spec>

<sup>11</sup> <https://swagger.io/specification/>

<sup>12</sup> <https://smart-api.info/>

<sup>13</sup> <https://www.w3.org/RDF/>

The heart of the F2DS is composed of a metadata storage called **FAIR data point (FDP)**. FDP is developed within the [GOFAIR](https://www.go-fair.org/) initiative<sup>14</sup> and in the [FAIRsFAIR project](https://www.fairsfair.eu/)<sup>15</sup>. On one side, this service enables data owners to expose datasets in a FAIR manner and on the other side, allows data users to discover properties of the datasets through the exposed metadata and to access the data for download depending on the license condition.

In this version of the tool set, we have been focusing on the metadata harvesting. The metadata enricher should be added in a later version of the service bundle and should be built using the FAIRifier<sup>16</sup> service which allows metadata enrichment with ontologies (leveraging the work of T5.5).

During the course of the EOSC-Pillar project, the FDP graph database is planned to be connected with the FAIR Digital Object Framework ([FDO-F](https://www.go-fair.org/today/fair-digital-framework/))<sup>17</sup>.

Finally, this solution is being tested with the repositories from INRA-E (T6.3) and IRD Data Terra (T6.2) and should integrate more datasets in the upcoming releases.

## 2.2 D4Science solution

The second technical solution to build the federated FAIR data space is characterised by the exploitation of the D4Science service (Assante et al. 2019a, 2019b) offering and operational settings. In particular, it has a catalogue as the key component that is capable to harvest items from data sources as well as to enable users to publish new items according to user/community defined profiles. Catalogue instances are integrated into working environment that is specifically created to serve the needs of their designated communities. This means that every working environment can be customised by selecting the data space and the set of tools to be made available. Hence researchers use this custom environment for further developing and consuming the specific catalogue content.

---

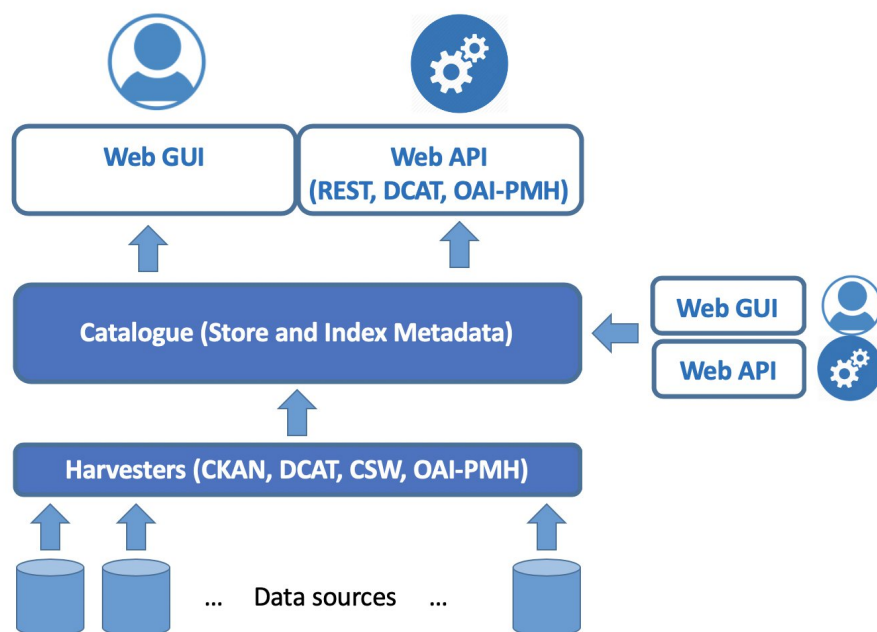
<sup>14</sup> <https://www.go-fair.org/how-to-go-fair/fair-data-point/>

<sup>15</sup> <https://www.fairsfair.eu/>

<sup>16</sup> <https://github.com/FAIRDataTeam/FAIRifier>

<sup>17</sup> <https://www.go-fair.org/today/fair-digital-framework/>





Architecture of the D4Science based F2DS

The **EOSC-Pillar Research Data Catalogue** has been created as instance of these tools. It is a virtual research environment and proof of concept of a working environment facilitating the development of the overall F2DS. The environment offers the basic services enabling researchers to collaborate and share material (e.g. in a social networking area and in a workspace for storing files of interest) and publish data in a catalogue (enacting authorized users to publish new items and manage the published items). Instances of other tools can be added depending on forthcoming needs. The catalogue is planned to be populated with datasets and any other digital artifacts stemming from the EOSC-Pillar use cases represented in WP6. Currently the catalogue contains more than 80K items, mainly resulting from harvesting content from the INRAE Data Repository ([data.inrae.fr](https://data.inrae.fr)) via the OAI-PMH protocol and the IFREMER Catalogue (<https://sextant.ifremer.fr/geonetwork>) via the CSW protocol.

### 3 Access to the bundle of services

To ease the access to the current state of work, we provide in the table below the list of services together with the partners hosting the services, the access URL and a short description

#### DISCLAIMER

Please be aware that these services might not be available continuously because they are periodically updated. In the case where links are not working please contact by email the related contact person to obtain information on the status of the service i.e.:

- Nicolas Cazenave (cazenave@cines.fr) for CINES hosted services,
- Leonardo Candela (leonardo.candela@d4science.org) for the D4Science hosted services, and
- Lisana Berberi (lisana.berberi@kit.edu) for the service at KIT.

Name of the service	Hosted by	Access URL	Description
FAIR Data Point API	CINES	<a href="http://ffds.eosc-pillar.eu">http://ffds.eosc-pillar.eu</a>	API to access programmatically the content of the FAIR Data Point
FAIR Data Point User Interface	CINES	<a href="http://ffds.eosc-pillar.eu/fdp">http://ffds.eosc-pillar.eu/fdp</a>	User Interface to access the content of the FAIR Data Point
FFDS Register front	CINES	<a href="http://ffds.eosc-pillar.eu/front">http://ffds.eosc-pillar.eu/front</a>	User Interface for repository registration
FAIR Data Point SPARQL search	CINES	<a href="http://ffds.eosc-pillar.eu/blazegraph">http://ffds.eosc-pillar.eu/blazegraph</a>	Explore metadata in FDP with SparQL query language
D4Science data catalogue	CNR - ISTI (by D4Science)	<a href="https://eosc-pillar.d4science.org/group/eoscpillarresdataactlg">https://eosc-pillar.d4science.org/group/eoscpillarresdataactlg</a>	The virtual research environment is available at (for authorized users)
D4Science data catalogue- public use	CNR - ISTI (by D4Science)	<a href="https://eosc-pillar.d4science.org/web/eoscpillarresdataactlg/catalogue">https://eosc-pillar.d4science.org/web/eoscpillarresdataactlg/catalogue</a>	The publicly available version of the catalogue

---

CORDRA instance	KIT	<a href="http://fc05f45d-4bab-441b-8d8b-b96a523576d6.ka.bw-cloud-instance.org:8080/">http://fc05f45d-4bab-441b-8d8b-b96a523576d6.ka.bw-cloud-instance.org:8080/</a>	Repository compliant with Digital Object specification to be updated with upcoming FAIR Digital Object specification
-----------------	-----	---	--

## 4 Concluding remarks

This deliverable describes an initial tool-set supporting the implementation of the Federated FAIR Data Space (F2DS) concept. In particular, the concept of F2DS has been formulated as a unifying space of datasets (and other typologies of items) stemming from several data repositories and data sources as well as from several communities. The aim of the Federated FAIR Data Space is to make heterogeneous and distributed datasets compliant to the FAIR principles.

These tools are being built in collaboration with the various communities involved in the project (WP6 use-cases) thus collecting their feedback to consolidate and extend the service offering of the EOSC-Pillar Federated FAIR Data Space.

The current implementations have both ingested metadata from diverse data sources and communities. Further data-sets will be added in the coming months to showcase the benefits to other communities as well. Each of the implementations will thereby equally evaluated in the field. Ultimately the implementations will be turned into real services and be included in the local or EOSC wide services catalogue.

This first bundle of services focuses specifically on the services involved in creating a federated data space. In the next iteration, facilities for data enrichment will be added as well as changes resulting from the feedback collected from the use cases will be incorporated.

As dictated in the Description of Work, another major release will be done in Project Month 24.



## References

Assante, M. et al. (2019a) Enacting open science by D4Science. *Future Gener. Comput. Syst.* 101: 555-563 DOI: [10.1016/j.future.2019.05.063](https://doi.org/10.1016/j.future.2019.05.063)

Assante, M. et al. (2019b) The gCube system: Delivering Virtual Research Environments as-a-Service. *Future Gener. Comput. Syst.* 95: 445-453 DOI: [10.1016/j.future.2018.10.035](https://doi.org/10.1016/j.future.2018.10.035)

Wilkinson, M., Dumontier, M., Aalbersberg, I. et al. The FAIR Guiding Principles for scientific data management and stewardship. *Sci Data* 3, 160018 (2016). DOI: [10.1038/sdata.2016.18](https://doi.org/10.1038/sdata.2016.18)