








# A Primer on Open Science-Driven Repository Platforms

Alessia Bardi , Paolo Manghi , Andrea Mannocci<sup>(✉)</sup> , Enrico Ottonello ,  
and Gina Pavone 

Institute of Information Science and Technologies, National Research Council, Pisa, Italy  
{alessia.bardi, paolo.manghi, andrea.mannocci, enrico.ottonello,  
gina.pavone}@isti.cnr.it

**Abstract.** Following Open Science mandates, institutions and communities increasingly demand repositories with native support for publishing scientific literature together with research data, software, and other research products. Such repositories may be thematic or general-purpose and are deeply integrated with the scholarly communication ecosystem to ensure versioning, persistent identifiers, data curation, usage stats, and so on. Identifying the most suitable off-the-shelf repository platform is often a non-trivial task as the choice depends on functional requirements, programming and technical skills, and infrastructure resources.

This work analyses four state-of-the-art Open Source repository platforms, namely Dryad, Dataverse, DSpace, and InvenioRDM, from both a functional and a software perspective. This work intends to provide an overview serving as a primer for choosing repository platform solutions in different application scenarios. Moreover, this paper highlights how these platforms reacted to some key Open Science demands, moving away from the original and old-fashioned concept of a repository serving as a static container of files and metadata.

**Keywords:** Repository · Repository platforms · Repository software · Dryad · Dataverse · DSpace · InvenioRDM · Open Science · FAIR · Research Data

## 1 Introduction

Following Open Science mandates, institutions and communities increasingly demand repositories with native support for publishing scientific literature together with research data, research software, and other scientific products. As envisioned by the COAR Next Generation Repositories Working Group [1] and other initiatives in the digital libraries' domain [2], repositories form “the foundation for a distributed, globally networked infrastructure for scholarly communication, on top of which layers of value-added services will be deployed, thereby transforming the system, making it more research-centric, open to and supportive of innovation, while also collectively managed by the scholarly community”. The increasing adoption of Open Science and Open Access mandates further urges this process, with many institutions, infrastructures, and research communities setting the ambitious goal of providing an “Open Science-driven” repository for their

community of users. Such repositories may focus on research datasets or become the holders of all kinds of research products, including publications, datasets, and software, and aim at supporting out-of-the-box features such as FAIRness, collaboration, access control, and data curation. Their instances should be able, to some extent, to adapt to community-specific requirements and ensure the degree of interoperability required to be part of scientific workflows (e.g., to fetch-analyse or generate-deposit products) and interact with third-party scholarly communication services, such as monitoring platforms (e.g., MakeDataCount), altmetrics, and PID registries (e.g., ORCID.org, ROR.org, Crossref.org, DataCite.org).

However, researchers' and policymakers' requirements often make choosing the suitable repository platform a non-trivial choice for organisations, as institutional and community demands may vary in terms of kinds of research products, metadata descriptions, and functionalities (e.g., data curation, collaboration). Typically, repository platforms may satisfy these requirements to some extent and then require adaptation to fit additional specific needs better [3, 4]. Some solutions are designed to maximise flexibility, e.g., enabling customisation of metadata descriptions and facilitating the modular integration of new functionalities, while others exhibit less flexibility in favour of a tailored, one-stop-shop product. Besides, local requirements, e.g., available skills and resources, may pose constraints and limitations to software choices and following customisation.

In the attempt to facilitate such a choice, this work analyses four state-of-the-art repository solutions, namely Dryad, Dataverse, DSpace, and InvenioRDM, from both a functional and a software perspective, addressing the selection needs of organisations willing to become repository providers (for an end-users-oriented survey, please refer to [13]). The investigation is intended as a “primer” for choosing repository platform solutions in different institutional and community application scenarios. Most importantly, the paper highlights how such known platforms reacted to Open Science demands [5], moving away from the old-fashioned concept of a repository as static containers of files and metadata.

## 2 Repository Platforms

The four Open Source repository platforms at the centre of this analysis were picked among others due to the following reasons. Firstly, because of their wide adoption by institutions and communities worldwide, a trend often followed by company uptake. Secondly, their design and functionalities evolved to address the demands of Open Science scientific workflows [6]. Finally, the authors have familiarity with the four platforms and could rely on feedback from adopters; this explains why software platforms such as ePrints<sup>1</sup> and Fedora/Islandora<sup>2</sup> have not been included. This paper aims to analyse the platforms to identify their specific reactions to such demands, moving away from old-fashioned repositories conceived as static containers of files and metadata.

For this analysis, we referred to the work of the COAR Next Generation Repositories Working Group<sup>3</sup>, for which a “next generation repository” should:

<sup>1</sup> [www.eprints.org](http://www.eprints.org).

<sup>2</sup> <https://duraspace.org/fedora>.

<sup>3</sup> <http://ngr.coar-repositories.org>.

- **Support a diversity of research products** and thus manage, preserve, version, curate, and provide access to a broad range of *research products*, including published articles, pre-prints, datasets, working papers, images, and software;
- **Support a diversity of research communities and institutions**, and thus be to some extent customisable to satisfy community/local requisites [7] in terms of products, metadata, and functionalities;
- **Be part of an ecosystem of repositories and scholarly communication services** by interlinking via persistent identifiers their resources to relevant entities, such as author identifiers (e.g., ORCID IDs), project identifiers (e.g., FundRef), organisation identifiers (e.g., ROR.org, ISNI), other resources (e.g., data, software, literature via DOIs, arXiv IDs);
- **Be machine-friendly and interoperable** by adopting standards [8] that enable a broader range of scientific services, such as scientific workflows, discovery, access, annotation, sharing, quality assessment, content transfer, analytics, provenance tracking, recommendations, and so on.

The following two sections point out the desiderata that can be derived from such a vision in terms of functionalities and software features. *Functional desiderata* capture the ability to address Open Science resources and workflows and the proactivity expected by modern repositories in the context of scholarly communication. *Software desiderata* instead frame a software project's maturity, flexibility, and modularity, pointing out the degree of customisation the platform can meet.

For the sake of space, we deliberately left out *operational desiderata*, which may depend on the individual installation policies and resources. The most prominent ones identified during the investigation are *long-term preservation* [12] (i.e., terms of commitment towards long-term storage of resources) and *free deposition strategies* (e.g., quota allowed per research product/user, which demands fees to be exceeded).

## 2.1 Functional Desiderata

After reviewing the selected repository platforms, the following relevant functional desiderata emerged (see Table 1 for a summary):

**Research Product Types.** Following the approach of resource modelling recommended by the European Open Science Cloud<sup>4</sup> (EOSC), research products can be classified into four meta-entities: *publications*, e.g., articles, theses, reports, presentations, *research data*, e.g., tables, images, archives, *research software*, e.g., code, *and others*, i.e., all products whose nature does not match one of the other entities. The metadata descriptions of such products may differ profoundly, ranging from bibliographic descriptions to provenance and community-specific tags. Moreover, metadata may include semantic links to other products to capture the entire research lifecycle for the sake of discoverability and reproducibility.

**Data (Metadata & File) Curation Functionalities.** The ability to engage scientists in data curation and validation processes is becoming prominent, as trust in research

---

<sup>4</sup> <https://eosc-portal.eu>.

data and software is undermined by a general lack of policies, practices, and tools for certification of quality [9–11]. Data curation functionalities (where “data” means everything that is metadata or files) regard two main aspects. The first is to offer collaboration and validation tools to a group of community curators to ensure the data matches the expectation of the community at hand in terms of quality, formats, and so on. The second is to ensure that end-users, i.e., the scientists, can establish virtuous interactions with the curators to make sure data is published with the expected quality.

**Integration with Entity Registries.** To adhere to Open Science demands and mandates, repositories use persistent identifiers (PIDs) for scholarly communication entities. On the one hand, they provide PIDs for the products uploaded by the users. On the other hand, they enable referencing to scholarly communication entities via PIDs by connecting to the related registries; examples are ORCID for authors and ROR.org for organisations. Integration with PID systems, i.e., registries, can be supported at two different degrees. The basic integration level is one where the repository metadata includes fields dedicated to interlinking with external entity registries, managing entity identities (via PIDs, cool URIs, handles, and so on), such as DataCite, Crossref, ORCID, ROR, Commons. The approach is subject to human mistakes in the format of PID, which may be “misspelt”, or in the referencing, i.e., an existing, yet wrong, PID may be used. A deeper and optimal level of integration is one where the insertion of a PID is supported by direct interaction with the related registry APIs, ensuring both format and PIDs are correct.

**Access Control.** Access control provides users with different levels of restriction options and granularity regarding research product access. Users may deposit research products and fine-tune access rights (e.g., restricted, open access, embargo) for metadata and/or files and to all users or groups of users (e.g., a research community).

**Table 1.** Functional desiderata.

<i>Desiderata</i>	<i>Description</i>
Research product types	Type of research products that users can deposit (e.g., publications, datasets, and/or software)
Data curation functionalities	Validation, rejection, and curation of metadata and/or files of research products
Integration with entity registries	Integration with PID systems to support and contribute to a non-ambiguous scholarly record
Access control	Users can rule access to metadata and files they deposit

## 2.2 Software Desiderata

Open Source repository platforms may be deployed by organisations with ICT capacity, whose requirements may derive from local technological constraints, peculiar functionalities, or, conversely, due to lack of ICT resources, by organisations that require

ready-to-go solutions. After reviewing the candidate repository platforms, the following relevant software desiderata were identified (see Table 2 for a summary):

**Software Project Sustainability.** The maturity, traction, and licencing of a software package are key requisites for an organisation to invest in a software product.

**Functionality Customisation.** Repository platforms are typically modular, meaning new functionalities can be easily plugged into the system, but to different degrees, with a trade-off between out-of-the-box and customisation.

**Metadata Model Customisation.** The repository can be more or less flexible concerning the metadata model, e.g., the attributes, the vocabularies, and references to external PID systems or registries in the scholarly communication infrastructure. Customisation of the metadata format measures the potential reuse across different communities and use cases. Still, as a counterpart, it impacts the out-of-the-box capabilities of a repository, which cannot be grounded on data model assumptions.

**Custom Storage Infrastructure.** Repository software must be configured and deployed to address potentially different scenarios, such as the cross-institutional, cross-country deployment setting, or community one. Different storage requirements may apply in terms of scalability, preservation, and availability of resources. Examples are Amazon S3 standard storage or simpler local storage solutions, typically provided by institutional data centres. The extent of customisation of the storage infrastructure is therefore relevant to making the right choice.

**Integration with Scientific Services.** Programmatic access enables third-party services to perform product depositions, metadata searches, exports, and downloads via APIs. The former allows for the implementation of scientific workflows capable of depositing into the repository on behalf and prior authorisation of the scientists. The latter ensures the repository can expose its content to other scholarly communication services, such as aggregators, ultimately enabling the realisation of custom UIs using the repository as a back-end (e.g., Zenodo.org).

**Persistent Identifiers.** Repositories must rely on persistent identifiers, typically issued at the record level, to uniquely refer to the pair metadata-files. Software platforms may be more or less flexible concerning the identifiers scheme to be used (e.g., handles, DOIs) by offering support to one or more specific PID Agencies (e.g., DataCite, EZID) and by enabling the integration with any PID agency.

**Usage Statistics.** Repository platforms are increasingly integrating with usage statistics infrastructures (e.g., IRUS-UK<sup>5</sup>, MakeDataCount<sup>6</sup>, OpenAIRE<sup>7</sup>) compliant with COUNTER Code of Practice<sup>8</sup>. Repositories, on their occurrence, centrally share views and downloads events of research products via the related PIDs, enabling aggregation of PID usage statistics across different repositories.

---

<sup>5</sup> <https://www.jisc.ac.uk/irus>.

<sup>6</sup> <https://makedatacount.org>.

<sup>7</sup> <https://www.openaire.eu>.

<sup>8</sup> <https://support.datacite.org/docs/counter-code-of-practice>.

**Table 2.** Software desiderata.

<i>Desiderata</i>	<i>Description</i>
Software project sustainability	Software project trust is measured by the engagement of developer communities
Functionality customisation	Modular design enabling the extension/customisation of functionalities
Metadata model customisation	Degree of research product types, metadata customisation, vocabularies, etc.
Custom storage infrastructure	Degree of customisation of storage infrastructure
Integration with scientific services	Ability to integrate with scientific services to publish products and/or provide access to products programmatically via APIs
Persistent Identifiers	Platform embeds functionality to mint PIDs for deposited research products
Usage statistics	Availability of modules to integrate with usage statistics infrastructures

### 3 Repository Platforms Analysis

#### 3.1 Dryad

“Dryad is an open source, community-driven project that takes a unique approach to data publication and digital preservation. Dryad focuses on search, presentation, and discovery and delegates the responsibility for the data preservation function to the underlying repository with which it is integrated”<sup>9</sup>. Research data are uploaded with metadata representing the dataset landing pages. Still, they are formatted as an online version of a data paper that can be downloaded as an individual PDF file or as part of the complete dataset download package, incorporating all data files for all versions.

**Functionality.** Dryad supports research data deposition only, although software deposition is possible via integration with Zenodo.org. For the official Dryad installation, curation is performed at the instance level, ensuring metadata is complete and both metadata and files comply with the platform recommendations<sup>10</sup>. Automatic validation tools are available for tables. The platform is integrated with ORCID, as it requires an ORCID to log in and supports ROR.org and FundRef IDs to refer to organisations and projects/grants funding, respectively. Metadata and files in Dryad are by the policy under CC0 waiver, so no fine-grained access control is supported.

**Software.** Dryad’s software is released on GitHub<sup>11</sup> under the MIT Licence and maintained by a community of 21 contributors (as of August 2022). The software is modular

<sup>9</sup> <https://datadryad.org>.

<sup>10</sup> <https://datadryad.org/stash/faq#files>.

<sup>11</sup> <https://github.com/CDL-Dryad/dryad-app>.

and based on the Stash software, organised into three modules: *Store* (deposition of metadata and files), *Harvest* (export of metadata to third-party services and a full-text Solr index), and *Share* (GeoBlacklight UIs). Its design enables the extension/customisation of metadata export protocols and the personalisation of metadata schemas. The Store module supports the DataCite format, but can be extended to include extra fields. It supports SWORD 2.0 to temporarily enable programmatic access from journal publishing platforms during article submission to deposit research data for peer review. The Harvest module supports OAI-PMH and ResourceSync and can be customised to support different export protocols.

### 3.2 Dataverse

“The Dataverse Project is an open-source web application to share, preserve, cite, explore, and analyse research data. It facilitates making data available to others and allows you to replicate others’ work more easily”<sup>12</sup>. A Dataverse repository hosts multiple archives called Dataverses, each intended as a collection of datasets consisting of descriptive metadata and data files. As a design choice, Dataverse collections may also be nested. Dataverse has been conceived to automate archivists’ and librarians’ tasks and to provide services for and to distribute credit to the data creator.

**Functionality.** The Dataverse platform focuses on publishing research data and related supplementary material (e.g., code and docs supporting the data). Dataverses are created by super-users that can assign to users nine different roles, establishing rights for publishing (draft, public, to be validated), accessing (read-only, update, delete, access to the record, access to files), and curating (right to update and publish). By creating a Dataverse and assigning the proper user roles, Dataverse installations can support custom data curation workflows. The platforms support user interfaces for manually reviewing tabular data.

Dataverse repositories are integrated with ORCID and ROR.org via APIs to ensure up-to-date references to authors and institutions. Access to research data can be controlled at the level of the Dataverse collection, at the fine-grain level of metadata records and individual files. The restricted mode, i.e., download upon request to the owner, is optional.

**Software.** Dataverse’s software is released on GitHub<sup>13</sup> under Apache Licence v2.0, maintained by a community of 144 contributors (as of August 2022). Designed as a ready-to-deploy package, which can be customised in some core capabilities such as the underlying storage system and the data model. Further customisation is possible via plugins, for example, to fetch vocabulary or entity data from external information systems such as registries.

---

<sup>12</sup> <https://dataverse.org>.

<sup>13</sup> <https://github.com/iqss/dataverse>.

Dataverse software supports three levels of metadata: metadata for citation (standard DDI), metadata for journal info (linking to external publications), and disciplinary metadata (provided with six default templates). The three come with a Dataverse schema, which can be further customised to include application-specific fields. It is possible to create new templates for the discipline metadata, to be shared with the community. The platform is also compliant with schema.org to support Google's data search crawlers.

Dataverse user interfaces can be extended via different tools developed by third parties<sup>14</sup>, e.g., for data previewing and data curation. Dataverse is designed to be integrated with other systems via SWORD protocol for data deposition. Examples are OJS<sup>15</sup> and publisher platforms, supporting publication and data submission workflows, but also integration with existing scientific workflows, e.g., Lab notebooks in RSpace<sup>16</sup>. The platform also supports OAI-PMH protocol standards for metadata harvesting. Its storage layer supports object and file system preservation via S3 or Swift and is configurable at the collection and dataset level.

### 3.3 DSpace

"DSpace is a web application allowing researchers and scholars to publish documents and data. [...] It is free and easy to install "out of the box" and completely customisable to fit the needs of any organisation"<sup>17</sup>. DSpace repositories are organised into *communities* (e.g., departments or institutions) that include *collections*, which are groups of *items*, i.e., data files and metadata descriptions.

**Functionality.** DSpace items can be set to model any kind of research product. For data curation, DSpace enables the definition of "Tasks" as plugins via a customisable curation framework<sup>18</sup>. Tasks enable checks and controls over metadata and files upon their deposition. The platform is integrated with ORCID for login and reference to authors. DSpace allows controlling read/write permissions at the instance level or per community, collection, item, and file. Administrative permissions per community or per collection can be delegated to users.

**Software.** DSpace is a mature project released on GitHub<sup>19</sup> under a BSD 3-Clause Licence and is maintained by a community of 166 contributors (as of August 2022), including companies doing business out of its custom extensions and installations.

The component pair UI/index (metadata store) is decoupled from the file storage. Files in DSpace can be stored either using a local filesystem (default) or a cloud-based solution, such as Amazon S3.

---

<sup>14</sup> <https://guides.dataverse.org/en/latest/admin/external-tools.html#inventory-of-external-tools>.

<sup>15</sup> <https://openjournalsystems.com>.

<sup>16</sup> <https://www.researchspace.com>.

<sup>17</sup> <https://dspace.lyrasis.org>.

<sup>18</sup> <https://wiki.lyrasis.org/display/DSDOC6x/Curation+System#CurationSystem-Tasks>.

<sup>19</sup> <https://github.com/dspace/dspace>.



DSpace comes with a suite of tools (e.g., batch ingest, batch export, batch metadata editing) and plugins for translating content into DSpace objects. By default, DSpace uses a Qualified Dublin Core (QDC) based metadata schema. Institutions can extend that base schema or add custom QDC-like schemas. DSpace can import or export metadata from other major metadata schemas, such as MARC or MODS. DSpace supports the Handle system by default but also integrates with DOI DataCite and EZID identifiers (ARK, DOIs). The platform offers custom APIs for the deposition of files and metadata and OAI-PMH for harvesting metadata.

### 3.4 InvenioRDM

“InvenioRDM is a turn-key research data management repository platform based on Invenio Framework and Zenodo”<sup>20</sup>. Its instantiations offer deposition and access functionalities to a set of communities (i.e., collections) of research products, which in turn can be of any kind.

**Functionality.** InvenioRDM includes communities to model collections of research products. Deposited research products can be managed by multiple users (i.e., shared submissions). Communities support curation/management workflows, where different users with different roles (i.e., curator, manager, reader, owner) are involved to ensure smooth, tracked deposition workflows. The deposition of metadata and files is structured as a “pull request” in software repositories, in which the submitter and curators (who can modify the metadata) are engaged in a discussion via an internal ticketing system. Workflows can be customised to include specific steps of approval at the community level: assigning roles of submitters subject to validation and curators notified of new submissions and in charge of the evaluation. Multiple curators can interact with the same submitter for the same submission. Also, requests for extra storage may be sent and handled by community managers.

Users can specify ORCID and ROR IDs for creators and related affiliations via the UIs. The selection of PIDs is enabled by ORCID and ROR APIs for validation; otherwise, textual values can be typed in by users.

Access can be controlled at the level of the community (restrictions to community and non-community members) or at the level of the record, at the granularity of the metadata and the files. The embargo function ensures that a record is made public at the expiring date without users performing any manual action.

**Software.** InvenioRDM is a rather young project, active since June 2019, released on GitHub<sup>21</sup> under an MIT licence, and maintained by a community of 32 contributors (as of August 2022). The software is developed as a specialisation of the Invenio Framework v3.0, glueing known Open Source tools such as Elasticsearch, OpenSearch, and Postgres, and based on JSON and DataCite format. The software comes with a ready-to-deploy configuration to deliver a repository instance similar to Zenodo.org.

---

<sup>20</sup> <https://inveniosoftware.org/products/rdm>.

<sup>21</sup> <https://github.com/inveniosoftware/invenio-app-rdm>.

The metadata data model implements the DataCite guidelines, but can be customised with extra fields to match community requirements. Interaction with vocabularies can be implemented by integrating external APIs into vocabulary systems or PID registries (e.g., ORCID, ROR).

Due to the flexibility of the Invenio Framework, the software is highly modular: storage can be of any kind (e.g., S3, file systems), and using different indexing and database systems is possible via programming efforts. The software supports OAI-PMH protocols and offers custom APIs for file and metadata deposition. The index offers indexing synchronisation functions, which mirror a new deposition in the InvenioRDM full-text index on external indexes.

## 4 Discussion

The four platforms offer ready-to-go mature or experience-based software solutions and can also benefit from companies for the configuration and installation of custom solutions. However, differences exist; Table 3 and Table 4 offer a high-level comparison of the functionality and software desiderata we have identified for the four platforms. In summary:

**Research products and metadata model customisation** beyond research data are addressed by InvenioRDM and DSpace. Instead, Dryad and Dataverse specialise in research data, offering dedicated and rich data management functionalities. All platforms generally offer a degree of customisability of the metadata descriptions. Dataverse, however, is designed to be extremely flexible in this respect, supporting a set of community profiles and a fully flexible metadata framework.

**Data Curation** is addressed in InvenioRDM by enabling interaction via UIs between data curators and end-users, integrating validation modules, and going beyond the validate-reject workflows. Similarly, but disregarding end users-curators interactions, DSpace offers a “Task” framework that developers can use to implement research data validation checks. Dryad and Dataverse offer manual validation and rejection procedures at the repository instance and the Dataverse level, respectively.

**Customisability of functionalities and storage** is equally addressed. InvenioRDM and DSpace seem to be the solutions that meet, at best, a scenario where the customisability of functionalities is a strong requirement. InvenioRDM’s software has been designed after the lesson learned in realising Zenodo.org using the Invenio Framework and meeting the requirements of Zenodo users and Invenio repository providers. Similarly, DSpace 7 has built on the core platform, and experience reached up to the release of DSpace 6 to bring a “single, modern user interface and REST API and integrates current technological standards and best practices”. As such, the platforms offer a good balance of out-of-the-box functionalities and flexibility of customisation and software extension. Dataverse also shows a high degree of customisability and a rich set of functional modules shared by the community. The four platforms allow the integration of different kinds of storage layers by modularly decoupling them from the user interfaces.

**Integration with entity registries** is well covered by Dryad, followed by InvenioRDM, and then DSpace and Dataverse, which only support ORCID. All platforms enable the integration of entity registries of reference via custom plugins.

**Persistent identifier minting** is thoroughly addressed by DSpace, supporting the open Handle system and optionally DOIs from DataCite and EZID identifiers. Dataverse follows with DOIs and Handles, while Dryad and InvenioRDM support DOI from DataCite. Plugins to other PID services are, in general, allowed.

**Integration with scientific services** is supported by all platforms for both deposition and harvesting. SWORD is provided by Dryad, DSpace, and Dataverse, while InvenioRDM implements a proprietary API. All platforms implement OAI-PMH, while DSpace offers ResourceSynch out-of-the-box.

**Usage statistics** are supported by all platforms via Make Data Count and COUNTER Code of Practice implementation.

**Table 3.** Functional desiderata comparison.

<i>Desiderata</i>	<i>Dryad</i>	<i>Dataverse</i>	<i>DSpace</i>	<i>InvenioRDM</i>
Research product types	Research Data	Research Data	All research products	All research products
Data curation functionalities	Manual rejection or approval at repository instance level	Manual rejection or approval at “dataverse” level	Customisable data “curation tasks” as validation controls over metadata and files upon deposition	Customisable data curation workflows, the interaction between submitters and collection managers/curators
Integration with entity registries	ORCID, ROR, FundRef	ORCID	ORCID	ORCID, ROR.org
Access control	Records and files are under CC0 waiver by default	At the granularity of “dataverses”, record, and files in the record	At the granularity of “sites”, “community”, collection, item and files per item	At the granularity of “communities”, record, metadata and files of the record; embargo functionality

**Table 4.** Software desiderata comparison.

<i>Desiderata</i>	<i>Dryad</i>	<i>Dataverse</i>	<i>DSpace</i>	<i>InvenioRDM</i>
Software project sustainability	21 contributors, MIT Licence	144 contributors, Apache Licence v2.0	166 contributors, BSD-3-Clause Licence	32 contributors, MIT Licence
Functionality customisation	Customisation of metadata export protocols	Integration with entity registries, customisation metadata exports and UI functions	Extendible software, a large pool of extensions is available	Extendible software. Building on the Invenio integration framework
Metadata model customisation	Metadata is DataCite but can be customised	Sets community metadata schemas which can be customised; supports schema.org	Metadata is Qualified Dublin Core and can be customised	Metadata is DataCite but can be customised
Custom storage infrastructure	Decouples storage from indexing and web portals	Decouples web portals from storage	Decouples web portals from storage	Decouples web portals from storage and indexing
Integration with scientific services	SWORD and OAI-PMH	SWORD and OAI-PMH	Deposition API, OAI-PMH, ResourceSynch	Deposition API, OAI-PMH, index synchronisation
Persistent identifiers	DOIs via DataCite	DOIs and Handles	Handle system, optional DOI from DataCite and EZID identifiers (ARK, DOIs)	DOIs via DataCite
Usage statistics	Make Data Count	Make Data Count	COUNTER Code of Practice	COUNTER Code of Practice

## References

1. Rodrigues, E., et al.: Next generation repositories: behaviours and technical recommendations of the COAR next generation repositories working group. Zenodo (2017). <https://doi.org/10.5281/zenodo.1215014>
2. Dempsey, L.: Library collections in the life of the user: two directions. *LIBER Q.: J. Assoc. Eur. Res. Librar.* **26**(4), 338–359 (2016). <https://doi.org/10.18352/lq.10170>
3. Austin, C., Brown, S., Fong, N., Humphrey, C., Leahey, A., Webster, P.: Research data repositories: review of current features, gap analysis, and recommendations for minimum requirements (Version 0) (2015). <https://doi.org/10.29173/iq904>
4. Assante, M., Candela, L., Castelli, D., Tani, A.: Are scientific data repositories coping with research data publishing? *Data Sci. J.* **15**, 6 (2016). <https://doi.org/10.5334/dsj-2016-006>

5. Jean-Claude, B., et al.: Open science, open data, and open scholarship: European policies to make science fit for the twenty-first century. *Front. Big Data* (2) (2019). <https://doi.org/10.3389/fdata.2019.00043>
6. RepOSGate: Open Science Gateways for Institutional Repositories, Michele Artini, Leonardo Candela, Paolo Manghi & Silvia Giannini
7. Jwa, A.S., Poldrack, R.A.: The spectrum of data sharing policies in neuroimaging data repositories. *Hum. Brain Mapp.* **43**(8), 2707–2721 (2022). <https://doi.org/10.1002/hbm.25803>
8. Forero, D.A., Curioso, W.H., Patrinos, G.P.: The importance of adherence to international standards for depositing open data in public repositories. *BMC Res. Notes* **14**, 405 (2021). <https://doi.org/10.1186/s13104-021-05817-z>
9. Liaw, S.-T., et al.: Quality assessment of real-world data repositories across the data life cycle: a literature review. *J. Am. Med. Inform. Assoc.* **28**(7), 1591–1599 (2021). <https://doi.org/10.1093/jamia/ocaa340>
10. Löffler, F., Wesp, V., König-Ries, B., Klan, F.: Dataset search in biodiversity research: do metadata in data repositories reflect scholarly information needs? *PLoS ONE* **16**(3), e0246099 (2021). <https://doi.org/10.1371/journal.pone.0246099>
11. Bashir, S., Gul, S., Bashir, S., Nisa, N.T., Ganaie, S.A.: Evolution of institutional repositories: managing institutional research output to remove the gap of academic elitism. *J. Librarianship Inf. Sci.* **54**(3), 518–531 (2022). <https://doi.org/10.1177/09610006211009592>
12. Barrueco, J.M., Termens, M.: Digital preservation in institutional repositories: a systematic literature review. *Digit. Libr. Perspect.* **38**(2), 161–174 (2022). <https://doi.org/10.1108/DLP-02-2021-0011>
13. Boch, M., et al.: A systematic review of data management platforms. In: Rocha, A., Adeli, H., Dzemyda, G., Moreira, F. (eds.) *WorldCIST 2022*. LNNS, vol. 469, pp. 15–24. Springer, Cham (2022). [https://doi.org/10.1007/978-3-031-04819-7\\_2](https://doi.org/10.1007/978-3-031-04819-7_2)