

Policy Research and Analysis of Foreign General-Purpose Data Knowledge Bases

Authors: Liu Jingjing, Gu Liping, Fan Shaoping

Date: 2016-06-03T00:00:00+00:00

Abstract

[Objective] Conducting policy investigation for establishing policy norms for institutional repositories and data repositories, and performing rights and interests analysis for librarians recommending data repositories. Abstract: [Objective] To conduct policy investigation for establishing policy norms for institutional repositories and data repositories, and to perform rights and interests analysis for librarians recommending data repositories. [Method] Employing a literature review approach, policy investigation and analysis were conducted to systematically review policy components and their substance. [Results] The study identified the rights and obligations of repository managers (establishing review mechanisms, formulating data identification standards, promulgating dissemination and usage regulations); the rights of submitters (free storage, updating metadata, setting embargo periods) and obligations (ensuring reliable data sources, complying with repository policies, avoiding intellectual property disputes); and the rights and obligations of users (free access, following citation requirements). [Limitations] The research lacks investigation into policies for specialized data repositories; a comprehensive policy framework can be developed in the future. [Conclusion] Establishing a comprehensive data repository policy, based on balancing the interests of all parties, is conducive to advancing open sharing of scientific research data.

Full Text

Policy Research and Analysis of General Research Data Repositories Abroad

Liu Jingjing^{1,3}, Gu Liping¹, Fan Shaoping^{2,3}

¹(National Science Library, Chinese Academy of Sciences, Beijing 100190, China)

²(Lanzhou Library, Chinese Academy of Sciences, Lanzhou 730000, China)

³(University of Chinese Academy of Sciences, Beijing 100049, China)

Abstract

[Objective] This study investigates and analyzes policies to establish norms for institutional repositories and research data repositories, and explores the rights and obligations involved when librarians recommend data repositories to researchers. **[Methods]** Through literature review and policy analysis, we systematically examine policy elements and their content. **[Results]** The findings reveal the rights and obligations of repository managers (establishing audit mechanisms, developing data identification standards, issuing dissemination regulations), submitters (rights to free storage, metadata updates, and embargo settings; obligations to ensure data reliability, comply with repository policies, and avoid intellectual property disputes), and users (rights to free access; obligations to follow citation requirements). **[Limitations]** The study lacks investigation into specialized research data repository policies, and future work should establish a comprehensive policy framework. **[Conclusion]** Developing robust data repository policies that balance the interests of all stakeholders will facilitate the advancement of research data sharing.

Keywords: Research data; Research Data Repository; Policy framework; Analysis of rights and obligations

Classification Number: G237.6

1. The Role of Data Repository Policy in Research Data Sharing

Data has become a “First-Class Citizen” in the scientific community, and relying solely on research papers for academic evaluation represents a historical limitation [1]. To effectively promote open sharing of research data, the scientific community has gradually developed standardized approaches for data publication and a system for recognizing data contributions. These include: storing and publishing research data in data repositories; publishing peer-reviewed data descriptors in data journals [2]; and publishing research papers supported by reproducible data in academic journals [3].

The relationships among these three pathways are illustrated in [Figure 1: see original paper] [4]. Research Data Repositories (RDRs) serve not merely as infrastructure for open data sharing but as the cornerstone of a standardized research data evaluation system. The focus of repository research extends beyond information system construction to encompass data quality review and assessment of data contributions, with particular emphasis on their role within the broader policy framework governing research data publication.

Science thrives on mutual criticism, growth, and learning. In the cyclical and dynamic process of scientific research, comprehensible research data provides

crucial support for verifying scientific conclusions [5-6]. Research data represents both the outcome of scientific inquiry and a valuable resource for inspiring further research [7]. To ensure data and associated information remain accessible, understandable, and usable, numerous stakeholders—including government agencies, research funders, educational institutions, data centers, researchers, librarians, and repository administrators—have made significant contributions [8], with related rights analyses gradually unfolding [9] and the foundational environment for research data sharing taking shape [10]. Open sharing requires not only reliable data storage but also rigorous and meticulous provenance management [11].

In response to mandatory storage requirements and encouraging recommendations from national governments, funding agencies, research institutions, and journals regarding data supporting research outcomes, numerous high-quality repositories serving scientific communities have been established internationally, such as Edinburgh DataShare [12], Open Data LMU [13], PANGAEA® [14], Dryad [15], figshare [16], and The Ber(li)n Digital Pantheon Project [17]. To help researchers identify and select appropriate repositories for storing and reusing data, registry and directory systems have emerged, including OAD [18] and re3data.org [19], which integrate, link, and present repository information in innovative ways [20].

China has also achieved progress in research data infrastructure, with projects such as the Research Data Sharing Initiative [21], Earth System Science Data Sharing Platform [22-23], Cold and Arid Regions Science Data Center [24], National Population and Health Science Data Sharing Platform [25], and National Agricultural Science Data Sharing Center [26]. These platforms focus on big data and project-specific data, forming a scientific data sharing system [27], though policy research on general-purpose data repositories remains nascent.

Data repositories evolve with changing scientific methods and environments, exhibiting different characteristics and classifications across disciplines, data types, and application levels [20]. Based on preliminary investigations [28-29], this paper categorizes data repositories into general and specialized types, as shown in .

** Classification of Data Repositories**

Type	Example	Discipline
General Research Data Repository	Dryad [15], figshare [16]	Natural sciences, engineering, medicine, humanities, and social sciences

Type	Example	Discipline
Specialized Research Data Repository	GenBank [30], PubChem [31], NOAA National Climatic Data Center [32], SIMBAD Astronomical Database [33], IQSS Dataverse Network [34], Durham HepData [35], ORNL DAAC [36]	Biomedical-DNA sequences, chemical information, environmental and earth sciences, astronomy, etc.

2. Research Framework

2.1 Research Questions

A data repository functions not merely as an information system but more importantly as an accelerator for research data sharing. Traditional information system policy research focusing on “creation-submission-management-use-preservation” workflow management proves inadequate for repository investigation. Instead, we adopt a “reverse thinking” approach from the perspective of “users” and “non-users” —specifically, why some repositories gain widespread researcher support while others remain underutilized.

This study examines the issue from a stakeholder perspective, encompassing research funders, managers, information service providers, publishers, and researchers. Simplifying around data itself reveals three key roles critical to effective repository operation: data managers (repository management and operations teams), data submitters (research teams or individuals contributing data), and data users (who may be the same or different researchers).

Consequently, our research question becomes: Who (the first role) advocates for whom (the second role) to store data, and stipulates how whom (the third role) may use the data? We frame the concerns of these three roles as our research questions, summarized in , with the corresponding policy observation questions guiding our analysis of international general repository policies.

** Analytical Framework for Data Repository Policies**

Stakeholder	Policy Observation Questions
Data Repository	1. How to guide submitters to effectively deposit high-quality data? 2. How to identify data for management purposes? 3. How to ensure permanent data accessibility after identification? 4. How to regulate appropriate use of repository data?

Stakeholder	Policy Observation Questions
Data Submitter	<ol style="list-style-type: none">1. What rights do submitters have regarding data/metadata storage, modification, and deletion?2. What self-review obligations apply to the data itself?
Data User	<ol style="list-style-type: none">1. Can users access data free of charge?2. What should users observe when using data?

2.2 Research Subjects

This study examines repository policies, focusing primarily on Dryad’s terms [37] with supplementary analysis of figshare [38], to systematically 梳理 the rights and obligations of all parties in data storage, use, and management, thereby summarizing policies for international general-purpose repositories. Dryad, managed by a non-profit organization, provides a generic home for diverse data types, primarily storing datasets associated with peer-reviewed publications in international repositories without format restrictions. Figshare, currently supported by Digital Science, allows users to upload various document types including images, datasets, multimedia, papers, and posters.

2.3 Research Methods

Our research design follows the process from “policy content collection” to “policy recommendation formulation” outlined in *On Foresight: Sharing Future with Forwarding Policy* [39], establishing observational priorities. During analysis, we reference Chapter 2 of *On Personas* [40], applying four principles of pragmatism to determine content selection, and employ rational methods to identify optimal solutions that explain current conditions, thereby distilling policy elements.

3. Policy Analysis: Repository Management

3.1 Audit Mechanisms for Ensuring Service Quality

Repositories review and organize published data to ensure metadata is standardized, accurate, and usable. Most data packages in Dryad undergo peer review [41], while figshare requires that submitted data contain no personal or medical information and comply with the UK Data Protection Act (1998) [42]. To guarantee service quality, repositories state they may conduct several reviews [43]: content review for personal, sensitive, or inappropriate information; copyright review; and review against file format and minimum reporting standards, with notifications to submitters or publishers when content fails to meet criteria. Additionally, repositories require submitter data files to be openable, uncorrupted, virus-free, and free from commercial interests or related disputes.

3.2 Data Identification Standards

Clear data provenance is required, with sources potentially including experimentally generated data, data obtained from other databases, or derived data compiled from others' work—all requiring explicit documentation. The core identification standard uses Digital Object Identifiers (DOI). In Dryad, each data package's DOI typically follows the format “[http://dx.doi.org/10.5061/dryad.\[NNNN\]](http://dx.doi.org/10.5061/dryad.[NNNN])” [41], where “[NNNN]” represents a 4-digit package number, appended with version information such as “/1”, “/2” (indicating file sequence). When new versions are released, update information like “.2”, “.3” is added. Thus, “[http://dx.doi.org/10.5061/dryad.\[NNNN\].2/2.3](http://dx.doi.org/10.5061/dryad.[NNNN].2/2.3)” indicates the third update of the second version of the second file in package [NNNN].

Figshare's DOI format is “[http://dx.doi.org/10.6084/m9.figshare.\[NNNNNNN\]](http://dx.doi.org/10.6084/m9.figshare.[NNNNNNN])” [44], complying with DataCite metadata standards and requiring users to add: Title, Authors, Categories, Tags, and detailed Description [45]. To maintain research integrity and continuity, version control is implemented with update indicators (Retrieved) on data pages. Theoretically, DOIs remain unchanged after data updates [46]. [Figure 2: see original paper] illustrates this using the dataset “All Hands to the Pump: Notes from NCCARF's 2010 International Climate Adaptation Futures Conference” [47].

3.3 Long-term Data Preservation

For long-term preservation, data migration is performed [37]: based on intellectual property agreements with authors, repositories may convert data formats to optimize storage capacity and efficiency for dissemination and reuse, ensuring daily updates and addressing potential conflicts with CC0 waiver provisions in licensing statements. Both Dryad and figshare partner with CLOCKSS [48] to preserve data content copies, migrating formats to latest versions to guarantee indefinite accessibility.

3.4 Dissemination and Use Regulations

Open Access Embargo Periods: Data publication should balance stakeholder interests through reasonable embargo periods, with repositories providing time range options. Dryad specifies embargos of 1-10 years [37].

Licensing Agreements: Repositories may provide explicit stipulations or recommend licenses (e.g., CC0 for metadata, CC-BY for data) [41,49]. International licenses include the Open Data Commons Public Domain Dedication and License (PDDL) [50], which establishes social norms for databases [51].

Third-party Indexing: Repositories allow content to be copied or indexed by third parties. Dryad supports linking with data journal articles and specialized repositories (e.g., GenBank, DataONE, TreeBASE).

Withdrawal Clauses: Post-publication, repositories monitor usage, documenting questions and withdrawal notices from submitters and users [37]. Dryad

may temporarily or permanently remove inappropriate content (sensitive, infringing, illegal, or legally risky material) after consultation with legal counsel. Modified metadata or documentation is updated, with revised data files linked to publisher errata. Some journal publishers collaborate with repositories (e.g., PLoS and Dryad) to jointly address data publication issues.

4. Policy Analysis: Submitter Rights and Obligations

4.1 Submitter Rights

1. **Free Storage Quota:** Submitters receive free storage up to a limit, beyond which additional fees apply. Dryad provides 10GB [37]; figshare offers 1GB [52].
2. **Metadata Update Rights:** Dryad permits submitters to update metadata of published packages and submit new or revised files without additional publication fees (though extra storage fees may apply) [37]. Figshare allows updates to categories, tags, and descriptions without creating new versions; only title, author, or file modifications generate new versions [53].
3. **Embargo Setting Rights:** Submitters may set open access embargos during which data is accessible only to project teams and authorized users, with full or batch public release after expiration [54]. Practice requires distinguishing among embargo period, open access period, and minimum preservation period.

4.2 Submitter Obligations

Submitters must ensure content correctness and legality while complying with repository regulations [37]:

1. **Authorship Verification:** Submitters must be authors or have obtained authorization from content authors, ensuring content accuracy without false or misleading information.
 2. **License Compliance:** Submission implies acceptance of repository licensing terms, permitting open access, promotion, format conversion, metadata modification, and partial content removal. Submitters must ensure compliance with publisher, funder, and institutional guidelines to avoid disputes.
 3. **Legal Compliance:** Submitters must guarantee content does not infringe intellectual property rights, personal privacy, or other national laws and regulations.
-

5. Policy Analysis: User Rights and Obligations

5.1 User Rights

Under open access principles, users may freely utilize research data. Dryad policies grant users rights to download, reprocess, reuse, and share data content, provided they follow repository licensing agreements. Scientific data open sharing traces back to the 2010 Panton Principles, which advocate allowing any user to download, copy, analyze, and reuse data via the internet for any purpose, free from financial, legal, or technical barriers [55]. Most repositories currently adhere to these principles.

5.2 User Obligations

Citation Standards: Dryad requires the following format [41]:

```
<Creator>( <Publication Year>) Data from:<Title>. Dryad Digital  
Repository.<Identifier>
```

[Figure 3: see original paper] demonstrates this citation format, including both the original article reference and the data package identifier. Similarly, figshare provides a “Cite this” link below each dataset [56], as shown in [Figure 4: see original paper].

Repository URL Citation: Beyond DOIs, repository URLs may be cited. The ArrayExpress functional genomics database requires citations to include both data identifiers and the ArrayExpress homepage URL [57].

Recommended Formats: Some repositories recommend DataCite formats. GEO (Gene Expression Omnibus) suggests citing both the original article and the corresponding data record identifier (GSExxx) [58].

Legal Compliance: Users must abide by repository policies and local national laws, refraining from illegal activities or actions that harm other users or disrupt repository functionality.

6. Implications and Limitations

6.1 Practical Significance

Based on the above analysis, we distill key policy elements for data repositories in , which can guide policy planning for new repositories, help librarians explain rights issues to researchers, and supplement institutional repository policies.

** Policy Elements for General Research Data Repositories**

Stakeholder	Policy Elements
Data Repository	(1) Establish audit mechanisms (metadata compliance, spot-check rights, data usability) (2) Develop data identification standards (provenance documentation, DOI format, versioning) (3) Implement data migration for long-term preservation (4) Issue dissemination regulations (embargo periods, licenses, third-party use, withdrawal clauses)
Data Submitter	Rights: (1) Free storage quota; (2) Metadata updates; (3) Embargo setting Obligations: (1) Ensure reliable data sources; (2) Comply with repository policies; (3) Avoid IP disputes
Data User	Rights: Free data access Obligations: Follow repository citation standards

6.2 Research Limitations

This study's systematic policy framework for general repositories is based on limited policy investigations. Generalizing conclusions requires attention to variations in national laws, research ecosystems, and funding models. For instance, Dryad applies CC0 licenses to all stored data and metadata, requiring IP transfer agreements with submitters and charging certain fees. Some international repositories use CC0 for metadata but CC-BY or other licenses for data itself.

Furthermore, repository policies cannot be directly transferred to institutional repositories without adaptation. General repositories serve global researchers, while institutional repositories serve specific institutional faculties. Both Dryad and figshare provide limited free storage, charging for additional capacity. Data differs from publications in requiring ongoing storage expansion and version management, with policies typically allowing metadata modification but requiring applications or additional fees for data changes. Institutional repositories differ not only in financial models but also in content scope and management approaches—important policy details requiring careful consideration.

6.3 Future Research

This study focuses on general repositories; specialized repositories require further investigation. The National Science Library, Chinese Academy of Sciences has compiled a *Policy Compilation for Research Data Repositories* report. Future work will survey Chinese researcher needs and exemplary domestic repository practices to develop best practice compilations providing more detailed selection references.

Author Contributions

Liu Jingjing: Paper revision, information supplementation, case collection and analysis, reference organization.

Gu Liping: Research design, information source identification, policy analysis, final paper revision.

Fan Shaoping: Initial draft, research question and framework development.

Acknowledgments: We thank Mr. Liu Feng from the Computer Network Information Center, Chinese Academy of Sciences and anonymous reviewers for their guidance.

Received: March 23, 2015

Revised: June 15, 2015

References

- [1] Bolikowski L, Houssos N, Manghi P, et al. Data as “First-class Citizens” [J/OL]. D-Lib Magazine, 2015, 21(1-2): DOI: 10.1045/january2015-guest-editorial. http://www.dlib.org/dlib/january15/01guest_editorial.print.html.
- [2] Liu Jingjing, Gu Liping. The Policy Research and Analysis of Data Journals: Taking Scientific Data as an Example [J]. Chinese Journal of Scientific and Technical Periodical, 2015, 26(4): 331-339.
- [3] Wu Rong, Gu Liping, Liu Jingjing. Research on the Data Policy of Academic Journals [J]. Library and Information Service, 2015, 59(7): 99-105.
- [4] Hayashi K, Murayama Y. Trend of Research Data Publishing and the Supported-article Data to Public Access [J]. Science and Technology Trends - Quarterly Review, 2015(1-2): 4-9.
- [5] The Royal Society. Science as an Open Enterprise [EB/OL]. (2012-06-21). [2015-03-14]. <http://royalsociety.org/policy/projects/science-public-enterprise/report/>.
- [6] Hey T, Tansley S, Tolle K, et al. The Fourth Paradigm: Data-Intensive Scientific Discovery [M]. Translated by Pan Jiaofeng, Zhang Xiaolin, et al. Beijing: Science Press, 2012.
- [7] Marchionini G. Research Data Stewardship: Ensuring Data Quality to Enable New Science in iSchools[J]. Translated by Yang Guancan, Lu Kun. Documentation, Information & Knowledge, 2013(4): 4-9.
- [8] Davidson J, Jones S, Molloy L, et al. Emerging Good Practice in Managing Research Data and Research Information Within UK Universities [J]. Procedia Computer Science, 2014, 33: 215-222.
- [9] Gu Liping. Basic Framework of the Analysis on Research Data Rights [J]. Documentation, Information & Knowledge, 2014(1): 34-51.
- [10] Zhang Xiaolin. Open Access, Open Knowledge, and Open Innova-

us.

[39] Gu Liping. On Foresight: Sharing Future with Forwarding Policy [M]. Taipei: Taipei Designer Press, 2013.

[40] Gu Liping. On Perssonas: Web User Information Behavior and Differentiated Services Strategy [M]. Scientific and Technical Documentation Press, 2013: 11-13.

[41] What Kinds of Data does Dryad Accept? [EB/OL]. [2015-03-14]. <http://datadryad.org/pages/faq>.

[42] How do I Publish My Data? [EB/OL]. [2015-03-14]. <https://figshare.zendesk.com/hc/en-us/articles/203712033-How-do-I-publish-my-data->.

[43] Dryad Obligations, Representations&Warranties to Purchasers and Submitters [EB/OL]. [2015-03-14]. <http://datadryad.org/themes/Mirage/docs/TermsOfService-Letter-2013.08.22.pdf>.

[44] figshare-Browse Data[EB/OL]. [2015-03-14]. <http://figshare.com/articles/browse#thumb>.

[45] What Metadata does Figshare Assign to Files? [EB/OL]. [2015-03-14]. <https://figshare.zendesk.com/hc/en-us/articles/203979333-What-metadata-does-figshare-assign-to-files->.

[46] How do I Upload a New Version of My Data? [EB/OL]. [2015-03-14]. <https://figshare.zendesk.com/hc/en-us/articles/203727056-How-do-I-upload-a-new-version-of-my-data->.

[47] All Hands to the Pump: Notes from NCCARF' s 2010 International Climate Adaptation Futures Conference [EB/OL]. [2015-05-31]. http://figshare.com/articles/All_hands_to_the_pump_Notes_from_NCCARF_s_2010_International_Cli

[48] The CLOCKSS Archive [EB/OL]. [2015-03-14]. <http://www.clockss.org/clockss/Home>.

[49] Why does figshare Use CC Licenses? [EB/OL]. [2015-03-14]. <https://figshare.zendesk.com/hc/en-us/articles/201953883-Why-does-figshare-use-CC-licenses->.

[50] Open Data Commons Public Domain Dedication and License (PDDL) [EB/OL]. [2015-03-14]. <http://opendatacommons.org/licenses/pddl/>.

[51] ODC Attribution-Sharealike Community Norms [EB/OL]. [2015-03-14]. <http://opendatacommons.org/norms/odc-by-sa/>.

[52] figshare Pricing [EB/OL]. [2015-03-14]. <http://figshare.com/pricing>.

[53] Does figshare Support Version Control [EB/OL]. [2015-03-14]. <https://figshare.zendesk.com/hc/en-us/articles/203922158-Does-figshare-support-version-control->.

[54] Si Li, Li Yueting. Analysis and Revelation on Research Data Preservation Policy Abroad [J]. Journal of Information Resources Management, 2014, 4(2): 40-50.

[55] Panton Principles-Principles for Open Data in Science [EB/OL]. [2015-03-14]. <http://pantonprinciples.org/>.

[56] Is the Research I Put on figshare Citable / Do I Get a DOI? [EB/OL]. [2015-03-14]. <https://figshare.zendesk.com/hc/en-us/articles/201954103-Is-the-research-I-put-on-figshare-citable-Do-I-get-a-DOI->.

[57] ArrayExpress-Submitting Data to ArrayExpress (general) [EB/OL]. [2015-03-14]. <http://www.ebi.ac.uk/arrayexpress/help/faq.html#cite>.

[58] GEO-Citing and Linking to the GEO Database [EB/OL]. [2015-03-14]. <http://www.ncbi.nlm.nih.gov/geo/info/linking.html>.

Note: Figure translations are in progress. See original paper for figures.

Source: ChinaXiv – Machine translation. Verify with original.