# Analysis and Formalization of Data Curation Process for Open Science

by

**Yasuyuki Minamiyama**

## Dissertation

submitted to the Department of Informatics

in partial fulfillment of the requirements for the degree of

## *Doctor of Philosophy*

SOKENDAI

The Graduate University for Advanced Studies, SOKENDAI

September 2022

# Evaluation Committee

Dr. Hideaki Takeda (Chair)
Professor, National Institute of Informatics/ SOKENDAI

Dr. Asanobu Kitamoto
Professor, National Institute of Informatics/ SOKENDAI

Dr. Hitoshi Okada
Associate Professor, National Institute of Informatics/ SOKENDAI

Dr. Ikki Ohmukai
Associate Professor, The University of Tokyo

Dr. Masao Takaku
Associate Professor, University of Tsukuba

.

# Abstract

With the trend of open science, there are growing expectations for research data sharing and reuse across fields. For research data reuse, data curation is essential to make the data interpretable and reusable. In some leading fields, field-specific tasks and procedures have been developed to implement systematic data curation. However, even in those leading fields, the research data reuse is often closed within the field. For interdisciplinary reuse of research data, research data must be interpretable by researchers from different fields. The problem here is the difference in data curation, which depends on the field. This variation in the data curation by field reduces the interpretability of research data. Without a method to overcome this challenge, open science will not be realized. This study will analyze and formalize the practices of data curation in each field for interpreting the data curation process. Furthermore, we will conduct practical studies for interpreting the data curation process with a formalized method.

In Chapter 2, we provide an overview of the research background and related areas. First, we review the current status of research data sharing and reuse with its social positioning. Then we discuss the significant studies from the data curation field. Finally, we show our approach to address the issues.

In Chapter 3, we analyze and formalize the knowledge of the data curation process to provide an interdisciplinary interpretation across different fields. The granularity of tasks and procedures that constitute data curation building blocks is not formalized, so it is impossible to identify common tasks and procedures across different fields. Here we propose an approach using ontology theory and techniques to interpret data curation tasks and procedures across fields. Our proposed ontology will allow data re-users to interpret the tasks and procedures performed with a formalized method.

The formalized interpretation leads to eliminating potential risks such as improperly misusing metadata. This study contributes to building a knowledge framework for a common understanding of the data curation process in different fields.

In Chapter 4, we deal with the processes in the Data Evaluation category formalized in Chapter 3 through practical implementations. To interpret research data from a field to be independently, the data and their documentation must be reviewed from an interdisciplinary perspective and revised as necessary. In contrast, many data repositories either do not have a policy for evaluating research data and encouraging improvement, or the evaluation process is done only from the perspective of a specific field. This chapter focuses on publishing mechanisms for data papers that divert the journal peer-review and clarifies the mechanism from a process perspective. Furthermore, we conduct a technical implementation for interpreting the mechanism using the reference model for data publishing. This study is a practical interpretation of the relationship and boundaries between the existing quality assurance process and the peer-review process in the Data Evaluation category.

In Chapter 5, we deal with the processes in the Appraisal category formalized in Chapter 3, focusing on the processes related to setting conditions of use. These processes have different variations due to the differences in legal restrictions and disciplinary norms by jurisdiction. Also, the conditions of use granted by data providers are more diverse. As a result, it is difficult for data re-users to interpret the results of the process accurately. This chapter investigates the actual processing status related to setting conditions of use in different fields and clarifies the correspondence of information tied to each formalized process. Furthermore, we conduct a practical implementation for stepwise interpretation of these processing results. This study reframes the processes in the Appraisal category and complements Chapter 4 practices each other.

In Chapter 6, we discuss the results of this study and prospects. Through this study, we provided our framework for understanding data curation activities as processes from an interdisciplinary perspective. This framework allows the data curation process to be interpreted in a decoupled way of the original research context. Data re-users will be able to formally assess the increased interpretability by verifying that the processes included in the framework have been properly executed. Furthermore, we have demonstrated some practical implementations through our problem-solving

approach as a stepwise formalization to the level of interpretation. We believe that as we move forward with these formalized efforts, we can improve the interpretability of research data and thereby contribute to the realization of open science.

# Contents

# List of Figures

# List of Tables

# 1

# Introduction

## 1.1   Motivation

In recent years, with the trend of open science, there have been many efforts to share and reuse research data on the Web (Kowalczyk and Shankar, 2011). The main purpose of researchers sharing research data is to improve research efficiency, increase verifiability, and generate new knowledge by research data reuse (Piwowar, 2011; Tenopir et al., 2011; OECD, 2015b). Research data reuse is an essential act for researchers to achieve open science (Fecher et al., 2015).

Research data reuse occurs when the data provider processes the research data to make it interpretable and reusable (Peer et al., 2014), and the data re-user uses the processed research data. The set of activities that make research data interpretable and reusable is called data curation (Sun and Khoo, 2017). The sequence of data curation processes includes various tasks such as cleaning, documenting, standardizing, formatting, and associating metadata with relevant research data and codes (CASRAI, 2019a). The high-quality metadata given by these tasks and mutual understanding of

the tasks makes published research data interpretable.

The practice of data curation has been developed mainly in fields such as life sciences (Venkatesan et al., 2019), earth sciences (Gray et al., 2002), and social sciences (Johnson, 2008). Through its historical efforts, tasks and procedures have been developed in these fields to implement systematic data curation (Ball, 2012). With the increasing reliable and interpretable research data, the research style of reusing other's research data is becoming the norm (Hemphill et al., 2022). However, even in those leading fields, research data reuse is often closed within the field (Yoon, 2016). For further development of open science, it is necessary to expand this norm across fields.

For interdisciplinary reuse of research data, research data must be interpretable by researchers from different fields (UNESCO, 2021). The problem here is the difference in data curation, which depends on the field. This variation in the data curation by field reduces the interpretability of research data. Without a method to overcome this challenge, open science will not be realized.

In order to interpret the variation in the data curation by field in common, it is necessary to interpret the data curation process from an interdisciplinary perspective. This study will analyze and formalize the practices of data curation in each field for interpreting the data curation process. Furthermore, we will conduct practical studies for interpreting the data curation process with a formalized method.

## 1.2   Objectives

This study aims to interpret the data curation process in various fields from an interdisciplinary perspective. To achieve this purpose, we set the following two objectives:

***Objective 1: Analysis and formalization of knowledge representing interdisciplinary data curation process***

***Objective 2: Practical studies to interpret the interdisciplinary data curation process***

Objective 1 is to formalize knowledge regarding the data curation process to build a framework for unifying interpretations. The process targeted for formalization will be reviewed for interdisciplinary understanding by data curators in each field.

Objective 2 is to interpret the data curation process with the framework through

Figure 1.1: The structure of this thesis.

practical implementations. It also includes exploring the possibility of providing appropriate technical support.

## 1.3   Thesis outline

Figure 1.1 shows the structure of this thesis.

**Chapter 2 Research background and related area**

In Chapter 2, we provide an overview of the research background and related areas. First, we review the current status of research data sharing and reuse with its social positioning. Then we discuss the significant studies from the data curation field. Finally, we show our approach to address the issues.

**Chapter 3 Formalizing the knowledge of the data curation process across different fields**

In Chapter 3, we analyze and formalize the knowledge of the data curation process to provide an interdisciplinary interpretation across different fields. The granularity of tasks and procedures that constitute data curation building blocks is not formalized, so it is impossible to identify common tasks and procedures across different fields. Here we propose an approach using ontology theory and techniques to interpret data curation tasks and procedures across fields. Our proposed ontology will allow data re-users to interpret the tasks and procedures performed with a formalized method. The formalized interpretation leads to eliminating potential risks such as improperly misusing metadata. This study contributes to building a knowledge framework for a common understanding of the data curation process in different fields.

**Chapter 4 Practical interpretation of the interdisciplinary data evaluation processes using the reference model for data publishing**

In Chapter 4, we deal with the processes in the Data Evaluation category formalized in Chapter 3 through practical implementations. To interpret research data from a field to be independently, the data and their documentation must be reviewed from an interdisciplinary perspective and revised as necessary. In contrast, many data repositories either do not have a policy for evaluating research data and encouraging improvement, or the evaluation process is done only from the perspective of a specific field. This chapter focuses on publishing mechanisms for data papers that divert the journal peer-review and clarifies the mechanism from a process perspective. Furthermore, we conduct a technical implementation for interpreting the mechanism using the reference model for data publishing. This study is a practical interpretation of the relationship and boundaries between the existing quality assurance process and the peer-review process in the Data Evaluation category.

**Chapter 5 Reframing the interdisciplinary appraisal processes by analyzing conditions of use of research data**

In Chapter 5, we deal with the processes in the Appraisal category formalized in Chapter 3, focusing on the processes related to setting conditions of use. These

processes have different variations due to the differences in legal restrictions and disciplinary norms by jurisdiction. Also, the conditions of use granted by data providers are more diverse. As a result, it is difficult for data re-users to interpret the results of the process accurately. This chapter investigates the actual processing status related to setting conditions of use in different fields and clarifies the correspondence of information tied to each formalized process. Furthermore, we conduct a practical implementation for stepwise interpretation of these processing results. This study reframes the processes in the Appraisal category and complements Chapter 4 practices each other.

**Chapter 6 Conclusion**

In Chapter 6, we discuss the results of this study and prospects.

# 2

# Research background and related area

In Chapter 2, we provide an overview of the research background and related areas. First, we review the current status of research data sharing and reuse, which serves as the research background, and its social positioning. Then we discuss the significant studies from the data curation field. Finally, we show our approach to address the issues.

## 2.1  Definition

In this Section, we define "Researcher," "Data provider," and "Data user" in this thesis.

"Researcher" is a person who conducts research using research data. Both a data provider and a data user are included.

"Data provider" is either a person or an institution with research data to provide a data user.

"Data user" is also a person or an institution that uses research data. It is sometimes

described as a "data re-user" to emphasize the use in different fields.

## 2.2 Research background

### 2.2.1 Research data

The term research data is defined as "Data that are used as primary sources to support technical or scientific enquiry, research, scholarship, or artistic activity, and that are used as evidence in the research process and/or are commonly accepted in the research community as necessary to validate research findings and results" (CASRAI, 2019b). It includes not only observation data and experimental data, which are generally easy to imagine, but also digitized data of historical documents, data of documentary photographs and images, text data of reprints, and so on. The term research data is used diverse ways in different fields: such as intermediate data, samples, and specimens may also be described as research data.

Among the various types of research data, research data to validate research findings and results is essential to ensure reproducibility. Research data validating research findings and results are also called evidence data; the evidence data should be stored in the same manner as academic papers as it constitutes part of the research. In recent years, some academic publishers have established data availability policies and made registering evidence data in data repositories a condition for starting the peer-review (CHORUS, 2021).

Research data treated as evidence in the research process are generally produced through the life of research projects. In fields such as life sciences (Koso, 2013) and earth sciences/astronomy (Murayama and Hayashi, 2014), where data acquisition costs are high, there has been the promotion of systematic collection and data sharing using large-scale data repositories. Such project-based data are not only used as primary sources for research but also as archives of research activities.

### 2.2.2 Research data sharing

Research data sharing is broadly defined as "providing access for use and reuse of data" (Tenopir et al., 2011). There are two possible methods of research data sharing: private data sharing among researchers and public data sharing through journal supplements and/or repositories (Borgman, 2015). Private data sharing among researchers are reported to be more common than public data sharing through journal supplements and/or repositories (Faniel and Jacobsen, 2010; Wallis et al., 2013). But having said that, public data sharing has an advantage over private data sharing in terms of long-term access and scalability (OECD, 2017). For this reason, more expectations have been placed on public data sharing.

Research data sharing is also perceived positively by the research community. In a recent international survey, 86.7% to 89.6% said they are willing to share data with broader research groups, although there are differences by field (Tenopir et al., 2020). Relatedly, there has been many discussion about research data sharing policies. One of the most well-known community norms for research data sharing is the FAIR principles (Wilkinson et al., 2016), which stands for Findable, Accessible, Interoperable, and Reusable. The FAIR Principles are expanding internationally, as exemplified by the EU's Horizon 2020 programme, which explicitly states that data management will be aligned with the FAIR principles (Commission, 2016). That being said, the FAIR principles are somewhat abstract and subject to diverse interpretations. Therefore, efforts to measure the FAIR-compliance degree are underway internationally, including the development of a FAIR maturity model (Bahim et al., 2020) and services to evaluate the degree of FAIR (Wilkinson et al., 2019; Clarke et al., 2019).

### 2.2.3 Research data reuse

Research data reuse occurs when the data provider processes the research data to make it interpretable and reusable (Peer et al., 2014), and the data re-user uses the processed research data. When research data reuse, researchers evaluate information such as relevance to their own research, trustworthiness, and reproducibility of research data (Faniel and Jacobsen, 2010). However, the information needed for evaluation is often tacit knowledge; researchers contact data providers as needed to discuss how the data will be created, cleaned, processed, analyzed, and reported (Zimmerman, 2008; Wallis

et al., 2013). Prior research also suggests that some contextual information is essential to most data re-users, including information about the physical, technological, and social environment in which the data was collected (Baker and Yarmey, 2009; Chin and Lansing, 2004).

Despite the high level of consensus regarding research data sharing, there has not been much progress in research data reuse. The previous research data sharing survey introduced in Section 2.2.2 also asked about the experience of research data reuse generated by others. According to this survey, the percentage of researchers who regularly reuse research data varies widely from 13.6% to 50.4% (Tenopir et al., 2020).

There have been several empirical studies on research data reuse in each field: engineering (Howard et al., 2010), astronomy (Sands et al., 2012), cancer epidemiology (Rolland and Lee, 2013), archaeology (Faniel et al., 2013), social sciences (Faniel et al., 2016; Yoon, 2014) in addition to the examples already mentioned. Through these studies, the difficulty of documenting context information is highlighted as challenge when research data reuse (Yoon, 2016).

### 2.2.4   Increasing social significance for research data sharing and reuse

Traditionally, research data on the basis of research papers have been shared to the extent necessary for verification, either as figures and tables or supplementary information. Efforts to share broader research data, such as project-based data have been limited to leading fields mentioned above. However, with the expansion of international collaborative research and the trend of open science, there are growing expectations for research data sharing and reuse across fields (Kowalczyk and Shankar, 2011). Starting with the G8 Science Ministers Statement in 2013 (G8, 2013), a full-scale activities for open research data have begun. In a recent example, UNESCO defines the ideal way of open research data as "available in a timely and user-friendly, human- and machine-readable and actionable format, in accordance with principles of good data governance and stewardship, notably the FAIR (Findable, Accessible, Interoperable, and Reusable) principles, supported by regular curation and maintenance" (UNESCO, 2021). Even in Japan, the Cabinet Office published the report "Promoting Open Science in Japan -Opening up a new era for the advancement of science-" (Cabinet Office,

2015) in 2015. This report made the keyword "open science" widely known in Japan. Moreover, the 5th Science and Technology Basic Plan includes promoting open science as a fundamental national stance (Cabinet Office, 2016). This national policy has been continued by the recent 6th Science, Technology, and Innovation Basic Plan (Cabinet Office, 2021). The research data sharing and reuse across fields are not only advancing research in the academic community but are also beginning to have strong social significance.

## 2.3 Related area

### 2.3.1 Data curation

Research data sharing and reuse are positioned as part of research activities. Since research data are collected, organized, and stored according to the interests of the researchers or research group, research data are deeply embedded in the research context. A knowledge gap exists in interpreting research data between the research group and third parties outside the community. To expand the possibilities research data reuse, it is necessary to make the data interpretable in a detached way from the research context.

Efforts to ensure the long-term continuous access and reusing digital information have been practiced mainly in the field of data curation. Data curation is the active and on-going management of data through its lifecycle of interest and usefulness to scholarship, science, and education; curation activities and policies enable data discovery and retrieval, maintain data quality and add value, and provide for re-use over time (of Illinois Urbana-Champaign, 2017).

### 2.3.2 Data curation activities, actions, and processes

The term data curation has two origins: as a subset of digital curation in information science and preprocessing in other scientific databases. In the former context, digital curation as a discipline was established around the 1990s and has its origins in archive/records management and library science (Higgins, 2011). Data curation activities include selection and evaluation by creators and archivists and developing

digital repositories as a technical foundation (Lee and Tibbo, 2007). In the latter, data curation is an older concept describing the process of selecting, normalizing, annotating, and integrating data from journals, reports, or other databases (Buneman et al., 2006). In fields such as astronomy, genomics, ecology, and economics that use data-driven methods, data curation has long been a specialized skill (Gray et al., 2002).

Despite these differences in origin, the two have been treated as almost identical concepts (Ball, 2010). According to the early definition (Lord and Macdonald, 2003), curation is defined as follows:

Curation: The activity of, managing and promoting the use of data from its point of creation, to ensure it is fit for contemporary purpose, and available for discovery and re-use. For dynamic datasets this may mean continuous enrichment or updating to keep it fit for purpose. Higher levels of curation will also involve maintaining links with annotation and with other published materials.

Thus, the latter process is incorporated into one of the activities as a "Higher levels of curation. " In a recent study, the components for expressing data curation activities are organized into four categories: work, action, curator, and repository (Johnston et al., 2016). The term 'action' is used as a concept that includes processes directed toward data preprocessing. It can be said that the terms 'action' and 'process' are interpreted in an integrated manner with different origins.

### 2.3.3   Theoretical model

Data curation activities are modeled in various fields; two representative theoretical models that support the field are the Reference Model for an Open Archival Information System (OAIS reference model) and the DCC curation lifecycle model. These two models provided a specialized language for the field, complementing each other.

**OAIS reference model**

The OAIS reference model was developed for use in facilitating a broad, discipline independent, consensus on the requirements for an archive or repository to provide long-term, preservation of digital information (Consultative Committee for Space

Data Systems, nd). This model describes the technical workflow for digital preservation and "establishes a common framework of terms and concepts" (Consultative Committee for Space Data Systems, 2012).

The OAIS reference model adopts a mechanism for storing multiple pieces of information placed around the content itself in a single package. This package information is called the Information Package, and consists of two major components: Content Information and Preservation Descriptive Information. OAIS approach is considered valid for preventing loss of meaning by identifying and preserving the information contained in both components.

Originally, the OAIS reference model was developed based on the space physics field practices. However, its usefulness was found early on. After a draft of the OAIS reference model became available, the CEDARS project was conducted on its applicability to the digital curation field (Jones, 2009). After the initial validation is complete, the OAIS reference model was widely recommended, adopted, and critiqued as a theoretical construct in the digital curation field. OAIS reference model is now an ISO standard (ISO, 2012), with version 2 being the latest as of May 2022.

**DCC curation lifecycle model**

The DCC curation lifecycle model was developed as a training tool to help curators understand the processes involved in successful curation and develop curation and preservation methodologies for their organizations (Higgins, 2008). This lifecycle model is unique because it provides widely used active verbs to describe the set of activities and the sequence that needs to be performed (Higgins, 2018). DCC curation lifecycle model was developed by reviewing previous lifecycle models; its functions, roles, and responsibilities appear to be strongly influenced by the OAIS reference model (Higgins, 2008; Lord and Macdonald, 2003). This conceptual model played a central role in driving data curation activities across multiple fields (Tibbo, 2012). In 2020, ver. 2 is being developed, which reviews issues such as multidisciplinary practices, the emergence of big data, and AI to address changes in social conditions after the development of ver. 1 (Choudhury et al., 2020).

## 2.4   Gaps and Challenges

Data curation's definition and the lifecycle are gradually expanded in response to the digital information revolution of the last two decades. Despite sharing these foundations, data curation activities in each field have developed independently. The problem here is the difference in data curation activities, especially at the level of actions or processes. For example, manual data cleaning and related tasks are often tacit knowledge and not documented in data curation records (Claerbout, 2010). Even if they were recorded, the granularity of the recorded information varies widely among the fields (Mayernik et al., 2013). Moreover, even if the granularity of recorded information is partially the same, identification is often difficult due to different representations of tasks and procedures (Borgman, 2007). This variation in the data curation actions or processes by field reduces the interpretability of research data activities in different fields. Without a method to overcome this gap, interdisciplinary reuse of research data will not be promoted.

## 2.5   Approach

To overcome the differences in data curation activities in different fields, it is necessary to unify interpretations regarding data curation actions or processes. This study aims to provide a framework for understanding data curation activities as processes from an interdisciplinary perspective. Data re-users will be able to formally assess the increased interpretability by verifying that the processes included in the framework have been properly executed. To build an interdisciplinary data curation process framework, this study will adopt a stepwise formalization approach to the level of interpretation.

For objective 1, we adopt a knowledge engineering approach. Methodologies for clarifying and systematically expressing certain knowledge have been studied mainly in the knowledge engineering field. The theory and practice of the knowledge engineering field are also likely to help formalize the data curation process framework.

For objective 2, we adopt a problem-solving type approach. The formalized framework needs to be elaborated based on practice. By focusing on the differences from an interdisciplinary perspective, we can analyze the data curation process in more detail. Also, a practical implementation may lead to the discovery of areas for

improvement and the interpretation of the entire process. By interpreting with a formalized data curation process, it will be possible to provide appropriate technical support.

# 3

# Formalizing the knowledge of data curation process across different fields

In Chapter 3, we analyze and formalize the knowledge of the data curation process to provide an interdisciplinary interpretation across different fields. The granularity of tasks and procedures that constitute data curation building blocks is not formalized, so it is impossible to identify common tasks and procedures across different fields. Here we propose an approach using ontology theory and techniques to interpret data curation tasks and procedures across fields. Our proposed ontology will allow data re-users to interpret the tasks and procedures performed with a formalized method. The formalized interpretation leads to eliminating potential risks such as improperly misusing metadata. This study contributes to building a knowledge framework for a common understanding of the data curation process in different fields.

## 3.1 Introduction

As outlined in Chapters 1 and 2, the differences in data curation tasks and procedures across different fields are a significant obstacle to realizing research data reuse. In particular, the different granularity of the process results in a significant reduction in interpretability. Without having a method for interpreting the data curation process across fields, it will be difficult to put it into practice on open science.

In order to interpret the tasks and procedures performed in different fields at the same granularity, it is necessary to identify the tasks and procedures in an interdisciplinary method as a knowledge. Methodologies for clarifying and systematically expressing certain knowledge have been studied mainly in the knowledge engineering field. Among them, ontology theory has been established and widely supported for constructing a conceptual system of knowledge (Uschold and Gruninger, 1996). Ontology theory and techniques has a possibility for interdisciplinary understanding for structural knowledge sharing of the data curation tasks and procedures.

This study investigates the practices of data curation conducted in each field. Then, to formalize the data curation tasks and procedures in different fields, we propose an approach using ontology theory and techniques that are considered helpful for structural knowledge representation. This study contributes to building a knowledge framework for a common understanding of the data curation process in different fields.

## 3.2 Literature Review

The data curation is commonly described with a research data lifecycle model (Kowalczyk, 2018). In a research data lifecycle model, the decisions involved in a set of data curation are divided into abstracted steps (Wallis et al., 2008). By performing data curation according to a lifecycle model, the data provider can perform each data curation task and procedure with accuracy and the data re-user can understand in detail the methodology and workflow used (Ball, 2012).

Two frameworks, knowledge creation and knowledge transfer, are presented as perspectives to better understand the data curation that takes place at each stage of the life cycle model (Humphrey, 2006). Regardless of the theoretical framework, the actual model is a mixture of both. Table 3.1 shows an example of the fields and steps involved

in a representative research data lifecycle (ICPSR, 2021; Faundeen et al., 2014; Oostdijk et al., 2013; DataONE, 2013; Chao et al., 2014; Johnston et al., 2016; Higgins, 2008; Griffin et al., 2018; UK Data Archive, nda).

Table 3.1: List of data curation activities by field.

| Name of Institutions / Communities | CLARIN-NL | Data Curation Network | DataONE | Digital Curation Centre | DPCVocab | EMBL Australia Bioinformatics Resource | ICPSR | UK Data Archive | U.S. Geological Survey |
|---|---|---|---|---|---|---|---|---|---|
| Fields | Humanities / Linguistics | Multiple | Earth Sciences | Multiple | Earth sciences/Life sciences | Life sciences | Social sciences | Social sciences | Earth Sciences |
| Steps | A: Identification and assessment B: Development of a curation plan C: Curation D: Validation E: Archiving | Ingest Appraise / Accept Curate Access Preserve | Plan Collect Assure Describe Preserve Discover Integrate Analyze | Conceptualise Create or receive Appraise / Select Dispose Ingest Preservation action Store Access, use and reuse Transform | Ingest Representation Provenance management Systems management Data storage Policies Preservation Public access provision | collecting integrating processing analyzing storing sharing publishing finding | Proposal development and data management plan Project start-up Data collection and file creation Data analysis Preparing data for sharing Depositing data | Transfer of data Assigning processing standard Data processing Documentation processing Metadata creation Additional user information Publishing data Delivering data Preserving data | Plan Acquire Process Analyze Preserve Publish/ Share |

The "Steps" row contains the steps defined by each organization, starting from the top. The steps defined by each field differ in terms of granularity. It is not easy to standardize decisions at each step throughout the life cycle of research data (Borgman, 2007).

Each step in the life cycle consists of multiple tasks and procedures necessary to make decisions. The tasks and procedures included in each field are more diverse than the steps themselves, and there is no comprehensive list of tasks and procedures performed in data curation across fields. In one of the few efforts to formalize definitions of tasks and procedures across fields, the Data Curation Network has drafted a glossary of terms to be used in a survey of important data curation activities in the U.S. (Johnston et al., 2016). This glossary is based on the existing glossary provided by the Digital Curation Centre (DCC), Society of American Archivists (SAA), CASRAI, RDA Data Foundation and Terminology Group, Digital Preservation Coalition (DPC), RDC (Research Data Canada), ICPSR, and practices in U.S. university libraries. Such efforts can be evaluated as potentially helpful in capturing the data curation tasks and procedures at the level of activities and supporting knowledge sharing. However, there still some issues: Some definitions do not distinguish between an 'action' performed by data curators and a 'process' performed by appropriate tools or repositories. An ambiguous expression may lead to misinterpretation in different fields where data curation activity is performed differently.

## 3.3   Research questions

This study aims to interpret the data curation activities in an interdisciplinary perspective and formalize the data curation activities in different fields. To achieve this objective, we set the following research questions (RQ):

RQ1: What are the commonalities and differences in the data curation tasks and procedures being performed in different fields?

RQ2: How to formalize the data curation activity's structure?

## 3.4 Preliminary analysis and survey

To clarify the commonalities and differences in the data curation tasks and procedures being performed in different fields, we conducted a preliminary analysis and survey in the following steps.

### 3.4.1 Analysis of existing data curation vocabulary

To clarify the commonalities and differences in data curation tasks and procedures, we need a working framework that can be used as a yardstick for different fields. As observed in Section 3.2, the Data Curation Network defines 47 vocabularies for the most important data curation activities derived from multiple lexical analyses. These vocabularies have been used in various fields of investigation (Hudson-Vitale et al., 2017) and are highly comprehensive; we have chosen to use the Data Curation Network vocabulary as our working framework for these reasons. Firstly, we analyzed the vocabularies using the IPO (Input - Process - Output) model to interpret from a process perspective. Table 3.2 shows a list of the 47 vocabulary sets subjected to analysis and the control structure expressed at the definition level.

Table 3.2: Results of Input – Process - Output analysis of data curation activity vocabularies.

| No | Activity | Input information | Process | Output information | Control structures |
|----|----------|-------------------|---------|--------------------|--------------------|
| 1 | Authentication | Data depositor identity information | Authenticate the identity of data depositors | Data depositor's identity authentication results | Sequential |
| 2 | Chain of custody | Data files | Generate data file provenance information | Data file provenance information | Sequential |

*Continued from previous page*

| No | Activity | Input information | Process | Output information | Control structures |
|---|---|---|---|---|---|
| 3 | Deposit agreement | Deposit agreement application information | Verify that deposited agreement file is fit for data repository's policies and conditions | Verification results of deposited agreement file | Sequential |
| 4 | Documentation | Information describing any necessary information to use and understand the data | Generate all information describing any necessary information to use and understand the data | Data document file | Sequential |
| 5 | File Validation | Data files | Generate and verify checksums for data files | Checksum verification result of the data files | Sequential |
| | | | Verify the data file format | File Format Verification Results | |
| 6 | Metadata | Information about a data set that is structured for purposes of search and retrieval | Generate necessary information about a data set that is structured for purposes of search and retrieval | Metadata file for purposes of search and retrieval | Sequential |
| 7 | Rights management | Data document file | Verify that retention and copyright rights inherent in data files are consistent with policies and conditions for access and reuse | Verification results on data file ownership and copyright | Sequential |

*Continued from previous page*

| No | Activity | Input information | Process | Output information | Control structures |
|----|----------|-------------------|---------|--------------------|--------------------|
| 8 | Risk management | Data files/Data document file | Verify that external constraints contained in data files are consistent with policies and conditions | Verification results of external constraints contained in the data files | Sequential |
| 9 | Selection | Verification results of deposit agreement/file format/data file ownership and copyright/external constraints contained in the data files | Verify that the results of the various verifications conform to the collection policy of the repository | Results of acceptance/rejection decision | Sequential |
| 10 | Arrangement and description | Data files | Re-organize data files according to standards and policies set by the repository | Data files (re-organized) | Sequential |
| 11 | Code review | Computer code | Verify the computer code | Verification results of the computer code | Sequential |

*Continued from previous page*

| No | Activity | Input information | Process | Output information | Control structures |
|----|----------|------------------|---------|--------------------|--------------------|
| 12 | Contextualize | Data document file/Metadata file for purposes of search and retrieval | Generate link information related to data files | Link information related to data files | Sequential |
| 13 | Conversion (Analog) | Analog data | Convert information into machine-readable format | Data files (converted into machine-readable format) | Sequential |
| 14 | Curation log | Execution results of the data curation process and executor information | Record changes made to the data and executor information during the data curation process | Information that records the execution results of the data curation process and executor information | Sequential |
| 15 | Data cleaning | Data files | Detect and fix (or remove) defects and errors in data files | Data files (cleaned) | Sequential |
| 16 | Deidentification | Data files | "Redact or remove personally identifiable or protected information (e.g., sensitive geographic locations) contained in data files" | Data files (de-identificated) | Sequential |

*Continued from previous page*

| No | Activity | Input information | Process | Output information | Control structures |
|----|----------|-------------------|---------|---------------------|---------------------|
| 17 | File format transformations | Data files | "Transform files into open, non-proprietary file formats" | Data files (trans-formatted) | Sequential |
| 18 | Transcoding | Data files | Encode audio/video files in ways that optimize reuse and long-term preservation actions | Data files (encoded) | Sequential |
| 19 | File inventory or manifest | Data files | "Verify the number of data files, file types (extensions), and file sizes periodically" | Verification results of data files | Sequential |
| 20 | File renaming | Data files | Rename data files | Data files (renamed) | Sequential |
| 21 | Indexing | Data document file/Metadata file for purposes of search and retrieval | Crosswalk to descriptive and administrative metadata compliant with a standard format for repository interoperability | Metadata files that conform to the repository's standard format | Sequential |
| 22 | Interoperability | Data files | Format the data using a disciplinary standard | Data files (formatted) | Sequential |

*Continued from previous page*

| No | Activity | Input information | Process | Output information | Control structures |
|----|----------|-------------------|---------|--------------------|--------------------|
| 23 | Peer-review | Data files/Data document file/Computer code | Validation of data files/data document file/computer code according to discipline-specific criteria by peers | Validation results of data files/data document file/computer code by peers | Sequential |
| 24 | Persistent Identifier | Data files/Metadata files that conform to the repository's standard format | Generate persistent identifier for data files | Persistent identifier for data files | Sequential |
|    |          |                   | Set up redirection when necessary | Redirect URL for data files | |
| 25 | Quality assurance | Data files/Data document file/Computer code | Validate data files/data document file/computer code according to the standards set by the repository | Validation results of data files/data document file/computer code | Sequential |
| 26 | Restructure | Data files | Organize and/or reformat poorly structured data files | Data files (restructured) | Sequential |
| 27 | Software registry | Data document file/Metadata file for purposes of search and retrieval | Maintain copies of modern and obsolete versions of software (and any relevant code libraries) | Copies of modern and obsolete versions of software (and any relevant code libraries) | Occasional |

*Continued from previous page*

| No | Activity | Input information | Process | Output information | Control structures |
|---|---|---|---|---|---|
| 28 | Contact information | Data document file/Metadata file for purposes of search and retrieval | Generate contact information for the data depositor and/or contact person | Contact information for the data depositor and/or contact person | Occasional |
| | | | Update contact information for the data depositor and/or contact person | Latest contact information for the data depositor and/or contact person | |
| 29 | Data citation | Metadata files that conform to the repository's standard format | Display of a recommended bibliographic citation | Recommended bibliographic citation text | Sequential |
| 30 | Data visualization | Data files/Data document file | Generate visualized data | Visualized data | Sequential |
| 31 | Discovery Services | Information on applying for connection to the discovery services/Metadata files that conform to the repository's standard format | Connect external discovery services | Discovery Service connection results | Sequential |

*Continued from previous page*

| No | Activity | Input information | Process | Output information | Control structures |
|----|----------|-------------------|---------|--------------------|--------------------|
| 32 | File download | Identifying information of authorized third parties/Metadata files that conform to the repository's standard format | Generate access URLs to data files by authorized third parties | Access URLs to data files by authorized third parties | Sequential |
| 33 | Full-text indexing | Data files | Generate text inherent in data file in search-engine-optimized formats | Full text information of the data files | Sequential |
| 34 | Metadata brokerage | Information on harvesting requests for metadata search and discovery services/Metadata files that conform to the repository's standard format | Set harvesting requests for metadata search and discovery services | Results of harvesting settings for metadata search and discovery services | Sequential |

| No | Activity | Input information | Process | Output information | Control structures |
|---|---|---|---|---|---|
| 35 | Restricted access | Access permission information/Access URLs to data files by authorized third parties | Set access permissions for data files based on access permission information | Access URLs to data files by authorized third parties restricted by access authority information | Sequential |
| 36 | Embargo | Embargo period information/Access URLs to data files by authorized third parties | Set an appropriate embargo period | Access URLs to data files with the embargo period set | Sequential |
| 37 | Terms of use | Metadata files that conform to the repository's standard format | Display information about the requirements or conditions for use provided to the end user of the data files | Information on the requirements or conditions for use of data files | Sequential |
| 38 | Use analytics | Data files/Data document file/Metadata files that conform to the repository's standard format | "Generate information on the frequency of data views, requests, and downloads" | Various usage information about data files | Occasional |

*Continued from previous page*

| No | Activity | Input information | Process | Output information | Control structures |
|----|----------|-------------------|---------|--------------------|--------------------|
|    |          |                   | Generate reuse metrics information such as data citations and impact measures for the data over time | | |
| 39 | Cease data curation | Information on data file storage and disposal plans | Plan for any contingencies that will ultimately terminate access to the data | Data Storage and Disposal Policy | Occasional |
| 40 | Migration | Data files | Transform obsolete file formats to new formats | Data files (migrated) | Occasional |
| 41 | Emulation | Copies of current versions of software (and any relevant code libraries) | Store and/or provide software to use the data files available in legacy systems | Software for emulation | Occasional |
| 42 | Secure storage | Data files | Back up data files on a regular basis | Backup data files | Occasional |
| 43 | File audit | Data files | Verify the digital integrity of data files | Verification results of digital integrity of data files | Occasional |

*Continued from previous page*

| No | Activity | Input information | Process | Output information | Control structures |
|---|---|---|---|---|---|
| 44 | Repository certification | A set of information about repository certification | Verify the technical and administrative capabilities of the repository by a trusted third-party accreditation body | Trusted third-party review results for repositories | Occasional |
| 45 | Succession planning | Information about the repository's long-term management plan | Develop a succession plan for the repository | Succession plan for the repository | Occasional |
| 46 | Technology monitoring and Refresh | Technical information about repository | Validate the performance of the repository against the latest technical requirements | Verification results of technical information | Occasional |
| 47 | Versioning | Data files | Generate version information for data files | Version information for data files | Occasional |

*End of table*

Next, we classified the control structure of the vocabulary into two categories based on the pairs of input and output information extracted from each vocabulary: Sequential processing, in which the output information of activity becomes the input information of a different activity (35 vocabularies); and occasional processing, in which activities are carried out independently from the time series (12 vocabularies). This classification is consistent with existing models (Higgins, 2008) so we judged it to be appropriate as a working framework. However, the following three points should be noted:

**1) Lack of vocabulary corresponding to the output information**    Some of the "generation" activities corresponding to the output information are not defined. For example, several activities have "data files" as input information, such as "Chain of Custody" or "File Validation," but the vocabulary for activities that output data files is not defined.

**2) Lack of a vocabulary with different hierarchies**    There are parallel and sequential processes that require multiple input information for some output information. However, some activities that aggregate multiple input information are not exist. For example, activities that have data files as input information ("Arrangement and Description," "Conversion," "Data Cleaning," "Data Visualization," "De-identification," "File Format Transformation," "File Renaming," and "Interoperability") are a series of activities that aggregate these activities to create an individual processed data file. However, "File Download" targets the processed data file that aggregates a series of these activities.

**3) No staffing/software information is included**    Each vocabulary does not include staffing information, so it is difficult to know the roles required to perform these activities. Also, some vocabularies are assumed to process by a repository software, which may make a difference depending on the software implemented.

## 3.4.2  Field survey

We conducted a field survey for several organizations conducting data curation activities in Japan. The purpose of the survey was to evaluate the actual status of data curation in each field using the working framework created in 3.4.1. As a supplement, we also intended to verify the validity of the working framework. First, we conducted interviews with the data curators at each organization. Table 3.3 shows an overview of the surveyed repositories.

Table 3.3: List of surveyed repositories.

| Organization name | Repository name | Name abbreviation | Repository type | Field | Repository Description |
|---|---|---|---|---|---|
| The Center for Global Environmental Research, Earth System Division, National Institute for Environmental Studies | Global Environmental Database | GED | Institutional | Global environmental issues | The Center for Global Environmental Research (CGER) at the National Institute for Environmental Studies (NIES) has created a Global Environmental Database (GED), which comprises data and research results collected and compiled from natural and social sciences. The GED serves as a fundamental database related to global environmental problems with an emphasis on global warming and climate change. |
| Center for Statistics and Information, Rikkyo University | Rikkyo University Data Archive | RUDA | Institutional | Social sciences | Rikkyo University Data Archive "RUDA" aims to collect, organize, and store social survey data which are valuable public assets, and make them widely available for research purposes such as academic secondary analysis and educational use in classes. |

*Continued from previous page*

| Organization name | Repository name | Name abbreviation | Repository type | Field | Repository Description |
|---|---|---|---|---|---|
| Japan Agency for Marine-Earth Science and Technology | Data and Sample Research System for Whole Cruise Information | DARWIN | Institutional | Marine-earth science | On the "Data and Sample Research System for Whole Cruise Information (DARWIN)" the Japan Agency for Marine-Earth Sciences (JAMSTEC) disseminates information for data, rock samples, and sediment core samples obtained by its research vessels and submersibles, and links to related databases. |
| Japan Science and Technology Agency National Bioscience Database Center | Life Science Database Archive | NBDC archive | Institutional | Life science | The Life Science Database Archive maintains and stores the datasets generated by life scientists in Japan in a long-term and stable state as national public goods. The Archive makes it easier for many people to search datasets by metadata (description of datasets) in a unified format, and to access and download the datasets with clear terms of use (see here for detailed descriptions). |

*Continued from previous page*

| Organization name | Repository name | Name abbreviation | Repository type | Field | Repository Description |
|---|---|---|---|---|---|
| National Museum of Japanese History | Knowledgebase of Historical Resources in Institutes | khirin | Institutional | Japanese history | "khirin (https://khirin-ld.rekihaku.ac.jp)" is the information infrastructure system that has been developed by the National Museum of Japanese History. "khirin" is an attempt to provide access to historical materials held by universities and museums on their networks as well as to offer data in a stable and sustainable manner in collaboration with the Japan Search. |
| National Institute for Materials Science | Materials Data Repository | MDR | Institutional | Materials science | MDR : Materials Data Repository is a data repository that hosts materials research data and publications. Discover various data and publications using metadata tailored for materials. MDR is operated by the National Institute for Materials Science (NIMS), Japan. |

*Continued from previous page*

| Organization name | Repository name | Name abbreviation | Repository type | Field | Repository Description |
|---|---|---|---|---|---|
| National Museum of Ethnology | Digital Picture Library for Area Studies | DiPLAS | Project | Ethnology | The purpose of this project is to support the representatives of Grant-in-Aid for Scientific Research projects conducting research in various regions of the world (including Japan), and to contribute to the research advancement by promoting the digitization and creating photographic materials database. |
| The Research Organization of Information and Systems, National Institute of Polar Research; Tohoku University; Nagoya University; Kyoto University; Kyushu University | Inter-university Upper atmosphere Global Observation NETwork | IUGON-ET | Project | Upper atmospheric physics | "We have three action plan in the second term (FY2015-) as follows: To provide the infrastructure and opportunity of the upper atmospheric research for users, in particular, in emerging countries. To provide our products and now-how for other fields and nurture human resources who can develop future database and utilize it. To promote the use of various data in a wide range of fields and support the advanced integration science. |

*End of table*

In selecting interviewees, we collected as various fields of practice as possible. On this basis, we limited our interviewees who can provide the following verification method: some form of documentation and/or the data curator's review. As a result, we conducted these interviews with people committing these repositories; four institutional repositories, i.e., Global Environmental Database (GED), Data and Sample Research System for Whole Cruise Information (DARWIN), Knowledgebase of Historical Resources in Institutes (khirin), and Materials Data Repository (MDR)) and two project-based repositories, i.e., Digital Picture Library for Area Studies (DiPLAS) and Inter-university Upper atmosphere Global Observation NETwork (IUGONET) from August to November 2020. We conducted additional interviews with those committing two institutional repositories, i.e., Rikkyo University Data Archive (RUDA) and Life Science Database Archive (NBDC archive) in August 2021. Each repository adopts various data curation models based on the nature and characteristics of the research data in each field. By comparing the models through an abstracted process, it is possible to extract commonalities and differences in structure. Each interview survey took about 1.5 to 2 hours. We used a topic guide to share the specific phase of data curation activities with the interviewee. In the topic guide, we set nine questions referred to the previous study categories (Johnston, 2017). The interview results were assigned to our working framework under the authors' responsibility and checked by each interviewer. The topic guide template used for the interviews is shown in Appendix 1.

Next, we read and referred each organization's data curation process manuals and related documents for the rationale for the activities to ensure consistency with the interview results. We mapped the specific description of the activities and the data curators' information onto the working framework for those activities for which we were able to identify a description of the rationale for the activities. The description of the rationale for the activities is shown in Appendix 2.

We can learn from the interview results concerning commonalities and differences among data curation activities in different organizations and fields. Firstly, we observed how much data curation activities can be interpreted in the same way. Figure 3.1 shows the implementation rates of data curation processes in eight repositories. We calculated the implementation rates using the following procedure: First, we classified each data curation process mapped to the working framework into three levels: 1. Implemented, 2. Partially implemented, and 3. Not implemented. The processes classified as "2.

■ all (8)   ■ multiple (2 to 7)   ■ individual (1)   ■ none (0)

Figure 3.1: Implementation rates of data curation processes in eight repositories.

Partially implemented" were mainly found when the vocabulary included multiple processes such as "generating and verifying checksums of data files" and "verifying file formats," as in "File Validation." Next, we aggregated the implementation number of organization by each activity. We counted "2. Partially implemented" as one organization. Finally, we classified the implementation number by four categories (all / multiple / individual / none) from the perspective of interpretability. As a result, we found that about 87.2% of the processes are interpretable across multiple fields. Among them, about a quarter of the processes were found to be fully interpretable across all fields.

Secondly, we observed the variety of staffing. Table 3.4 shows an overview related to the staffing of each repository.

Table 3.4: List of roles and number of appearances in eight repositories.

| Repository name (abbreviated) | Roles | Number of appearances |
|---|---|---|
| khirin | Researcher | 4 |

*Continued on next page*

*Continued from previous page*

| Repository name (abbreviated) | Roles | Number of appearances |
|---|---|---|
| | Related committee | 2 |
| | Center for Integrated Studies of Cultural and Research Resources | 27 |
| | Photographer | 2 |
| | System administrator | 1 |
| | Department of Rekihaku museum | 6 |
| | Department of internal database | 10 |
| | External organization | 1 |
| DiPLAS | Researcher | 2 |
| | Technical staff | 10 |
| | System administrator | 15 |
| | Data provider | 1 |
| | Project staff | 8 |
| | Digitization support staff | 1 |
| | Operation support staff | 1 |
| | Graduate students | 1 |
| | Review board | 1 |
| Materials Data Repository | Researcher | 6 |
| | Data system group | 14 |
| | Data service team | 13 |
| | System administration division | 1 |
| DARWIN | Researcher | 9 |
| | Data Management group | 42 |
| | Technician | 9 |
| | Navigation planning department | 2 |
| GED | Data provider | 14 |
| | Data curator | 29 |
| | Technical support staff | 1 |

*Continued from previous page*

| Repository name (abbreviated) | Roles | Number of appearances |
|---|---|---|
| | Web application developer | 1 |
| RUDA | RUDA manager | 33 |
| | Research assistant | 10 |
| | Researcher | 5 |
| | System administrator | 1 |
| | Related committee | 2 |
| IUGONET | IUGONET manager | 23 |
| | Researcher | 16 |
| NBDC archive | Contact information staff | 9 |
| | Researcher | 14 |
| | Data curator | 17 |
| | System operator | 6 |
| | Repository manager | 1 |

*End of table*

The roles defined by each repository are different, and there is no noticeable trend in the number of appearances. Each repository's data curation activities are carried out in different ways. For example, there are three staffing patterns in the "Data Cleaning" process: the data holders themselves, the data curator(s), and the 2 or 3 parties working together. It should be noted that some of these processes are covered by support systems or tools.

### 3.4.3   Section summary

First, we conducted a vocabulary analysis of data curation activities as a preliminary study to interpret commonalities and differences in data curation activities. As a result, we obtained a working framework to compare the data curation activities. Next, we conducted a field survey using the above working framework and found that about 87.2% of the processes are interpretable across multiple fields. In contrast, there were no similarities in staffing by fields. To accurately compare data curation activities in

different fields, the working framework should be improved by including the staffing as well as the three issues pointed out in 3.4.1.

## 3.5    Development of Data Curation Process Ontology

### 3.5.1    Formalizing commonalities and differences in data curation activity's structure

This section explores methods for formalizing commonalities and differences in data curation activities. Based on the preliminary analysis and the survey results, there needs an expressive method of the structure including the relationships among Input-Output objects, processes, hierarchical relationship among activities, and staffing in order to accurately represent the data curation activity's structure in different fields. Since these relationships are complicated, it is not impossible to represent the relationships in a simple tabular form. Some model is needed to describe these relationships adequately.

In this study, we adopt ontology as a model representation. We have developed an ontology that collects and structures knowledge to represent the data curation activity's structure. Ontology is one of the methods for constructing conceptual systems used in the knowledge engineering field. The ontology theory provides a framework for knowledge sharing by clearly defining concepts and describing the logical relationships between concepts. Developing an ontology makes it possible to manage processes in which people and information systems are mixed.

### 3.5.2    Development process

In order to develop an appropriate ontology, it is recommended to follow some ontology developing procedure. Developing an ontology is not an easy task since explicating and formalizing the conceptual system behind the target system need a very complex and abstract thinking and reasoning. To ease the task, several procedures to develop an ontology are proposed. As for the ontology development procedure, we followed the seven steps proposed by Noy McGuiness (Noy and McGuinness, 2001). In the actual work, we made several iterations between Step 4 and Step 6 to maintain consistency

with the hierarchical relationship. The OWL description of the ontology is shown in Appendix 3.

## Step 1: Determine the domain and scope of the ontology

In this step, we determine the domain and scope of the ontology to design an ontology. The decisions to be made include those for the domain to be covered by the ontology, the intended use of the ontology, the problem that the ontology to be developed can solve, and the maintainer of the ontology.

In our ontology, we aim to represent commonalities and differences in data curation activity's structure. The domain to be covered by this ontology is that of data curation. Providing the structured data curation activity in a machine-readable format can support the knowledge-sharing process between humans and information systems in a scalable manner. It would be desirable to maintain the ontology through the collaboration of data curators in each field and ontologists who deal with knowledge sharing in information systems.

## Step 2: Consider reusing existing ontologies

In this step, we consider the reusing existing ontologies. The description of this ontology is based on the vocabulary of the PROV ontology endorsed by W3C (Lebo et al., 2013). The PROV ontology is an ontology that provides a set of classes, properties, and restrictions that can be used to represent and exchange provenance information generated by different systems and different contexts. Since this study aims to interpret data curation activity's structure in different fields, including different systems, we determined the PROV ontology to be compatible for representing the data curation activity's structure.

As a basic structure of the PROV ontology, the information is represented by three classes and their relationships: Activity, Entity, and Agent. In the case of data curation activities, the data curation process can be represented as the "Activity" class, the input information and output information as the "Entity" class, and the staffing as the "Agent" class.

We mainly used the relationships defined in the PROV ontology to describe the relationships among Activities, Entities, and Agents. For some of the relationships, we used the FOAF ontology (http://xmlns.com/foaf/spec/) as a compliment.

**Step 3: Enumerate important terms in the ontology**

In this step, we enumerate important terms in describing the data curation activity' s structure. First, we added the Data Curation Network vocabulary used in our preliminary analysis. Also, we added four additional vocabularies to organize the input-output information pairs. The vocabularies we added are "SubmitData," "ActualDataProcessing," "MetadataProcessing," and "CreatingLandingPage." We extracted input information, output information, and the staffing role from the vocabulary in an abstract form to express the relationship between the data curation activity' s structure. The criteria for the extraction are described in detail in Step 4.
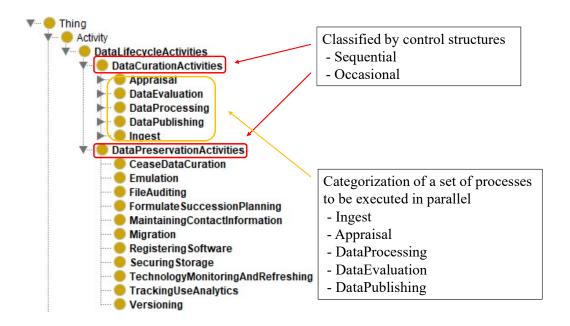


Figure 3.2: Data Curation Process Ontology structure.

**Step 4: Define the classes and the class hierarchy**

In this step, we define the classes and hierarchical relations of the ontology. Figure 3.2 shows the overall picture of this ontology's classes and hierarchical relations. Before determining the logical hierarchical relationship between the classes, we performed a categorical division of the activities; as shown in the preliminary analysis in Section 3.4.1, the extracted processes are a mixture of sequential and occasional processes. To separate the two types of activities with different control structures, we divided the classes into 'Data Curation Activities' for the sequential processes and 'Data Preservation Activities' for the occasional processes.

Next, we examined the logical structure of the 'Data Curation Activities'. Figure 3.3 shows the list of classes associated with each category. We set the following five categories under 'Data Curation Activities': "Ingest," "Appraisal," DataProcessing," "DataEvaluation," and "DataPublishing." We have already known that some sets of the data curation activity are performed in parallel from the preliminary analysis in Section 3.4.1. When handling this ontology, categorizing the process sets to be performed parallel helps interpretation. We set 22 processes under the five categories. In addition, four of the 22 processes have subclasses.

**Step 5: Define the properties of classes-slots**

In this step, we define the properties of the class-slots. Table 3.5 shows the list of properties used in this ontology.

This study adopted eight properties from the PROV ontology and one from the FOAF ontology. In describing the relationships in this ontology, we kept the description to the minimum necessary. In particular, the relationship between Activity and Entity is limited to "used" and "generated." In the reality of the data curation activity's structure, the relationship between Activity and Entity is far more diverse. For example, "CodeReview (Activity)" has the relationship of reviewing "sourceCode (Entity)."

But having said that, describing the elaborate relationship intends to complicate the properties' semantics. Since the complexity of semantics may affect the data curation activity's structure in different fields, we adopted the above policy as the first step in this ontology.
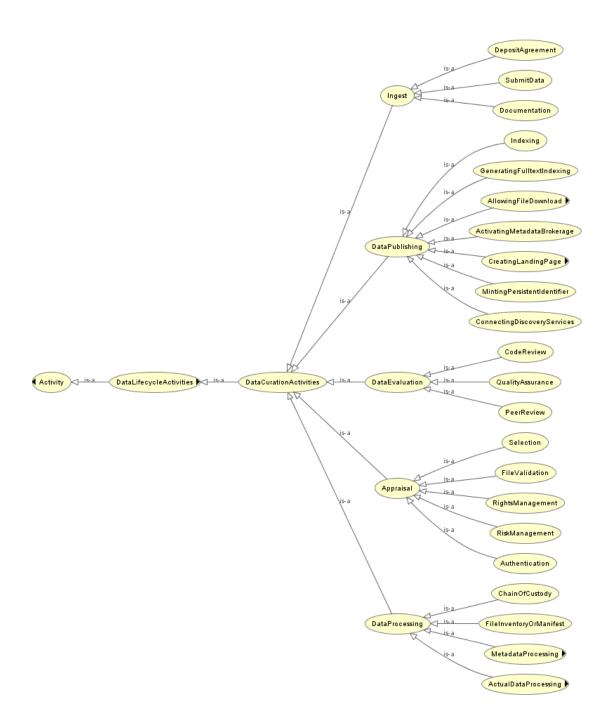
Figure 3.3: List of classes by category for sequential data curation processes.

Table 3.5: List of properties used in Data Curation Process Ontology.

| prefix | property |
|---|---|
| prov | used |
| | generated |
| | wasAssociatedWith |
| | influenced |
| | wasInformedBy |
| | wasAttributedTo |
| | wasEndedBy |
| | wasStartedBy |
| foaf | primary topic |

**Step 6: Define the facets of the slots**

In this step, we define the value type, allowed values, the number of the values (cardinality), and other features of the values as the facets that can be set for each slot. Since facets' values can vary depending on the type of research data handled, it is necessary to accumulate data based on actual output information. Here, we have set tentative values for constraint types that align with the actual situation obtained from preliminary survey results in Section 3.4.1.

**Step 7: Create instances**

In this step, we create an instance corresponding to the class of this ontology. Since this ontology abstracts the commonalities and differences of data curation activity's structure, it does not deal with the description of instances, which are individual phenomena. The description of the actual data curation activity's structure is treated in Section 3.6.

### 3.5.3 Basic structure of Data Curation Process Ontology

In this section, we explain the basic structure of the ontology created in Section 3.5.2. Figure 3.4 shows the basic structure of the configured class and an example of its

output in OWL format.

In this figure, "ActualDataProcessing (Activity)" is represented as generating "curated data (Entity)" by using "submittedData (Entity)." Also, "curator A (Agent)" is represented as a performer of "ActualDataProcessing." The corresponding OWL description is shown in red.

In organizing the hierarchical relationship between activities, we defined the structure so that the relationship with Entity can be expressed only by "used" and "generated" as described Figure 3.4 "ActualDataProcessing" is a vocabulary introduced to maintain this structure, with subclasses "ArrangementAndDescription," "Conversion," "DataCleaning," "DataVisualization," "Deidentification," "FileFormatTransformation," "FileRenaming," "Interoperability," and "Restructure". Similarly, "MetadataProcessing" has "Contextualization" and "MetadataGeneration", "CreatingLandingPage" has "DisplayingDataCitation" and "SettingTermsOfUse" as subclasses. This ontology has 184 classes and 313 subclasses generated as of version 1.

## 3.6 Case study

This section discusses a case study using the Data Curation Process Ontology. We used this ontology to describe the actual data curation activity's structures. Figure 3.5 through Figure 3.12 show the data curation activity's structure covered in Section 3.4.

```
:ActualDataProcessing rdf:type owl:Class ;
           rdfs:subClassOf :DataProcessing ,
                    [ rdf:type owl:Restriction ;
                      owl:onProperty <http://www.w3.org/ns/prov#generated> ;
                      owl:someValuesFrom :visualizedData
                    ] ,
                    [ rdf:type owl:Restriction ;
                      owl:onProperty <http://www.w3.org/ns/prov#influenced> ;
                      owl:someValuesFrom :accessRestriction
                    ] ,
                    [ rdf:type owl:Restriction ;
                      owl:onProperty <http://www.w3.org/ns/prov#used> ;
                      owl:someValuesFrom :dataDocument
                    ] ,
                    [ rdf:type owl:Restriction ;
                      owl:onProperty <http://www.w3.org/ns/prov#used> ;
                      owl:someValuesFrom :dataProcessingPolicy
                    ] ,
                    [ rdf:type owl:Restriction ;
                      owl:onProperty <http://www.w3.org/ns/prov#generated> ;
                      owl:minQualifiedCardinality "1"^^xsd:nonNegativeInteger ;
                      owl:onClass :curatedData
                    ] ,
                    [ rdf:type owl:Restriction ;
                      owl:onProperty <http://www.w3.org/ns/prov#used> ;
                      owl:minQualifiedCardinality "1"^^xsd:nonNegativeInteger ;
                      owl:onClass :submittedData
                    ] ,
                    [ rdf:type owl:Restriction ;
                      owl:onProperty
<http://www.w3.org/ns/prov#wasAssociatedWith> ;
                      owl:minQualifiedCardinality "1"^^xsd:nonNegativeInteger ;
                      owl:onClass :dataCurator
                    ] .
```
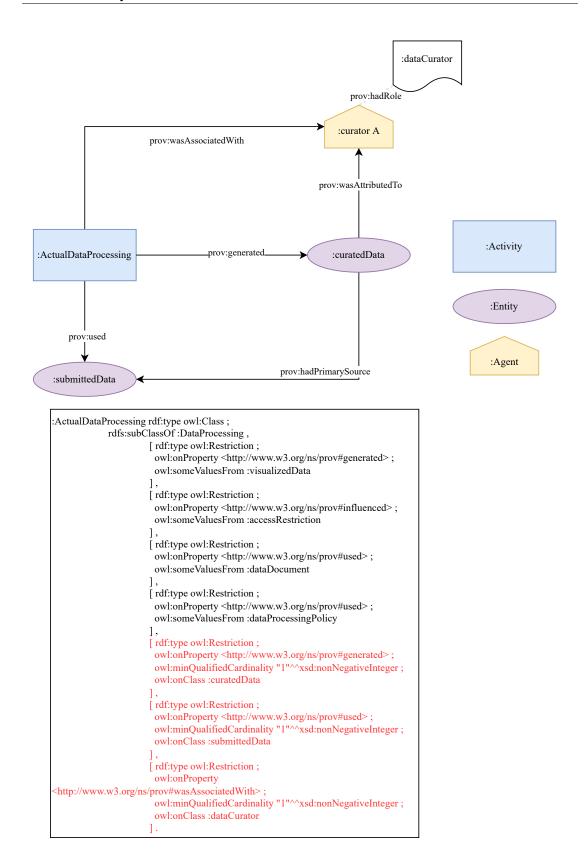
Figure 3.4: Example of the basic structure of "ActualDataProcessing" class and an output example of "ActualDataProcessing" class in OWL format.

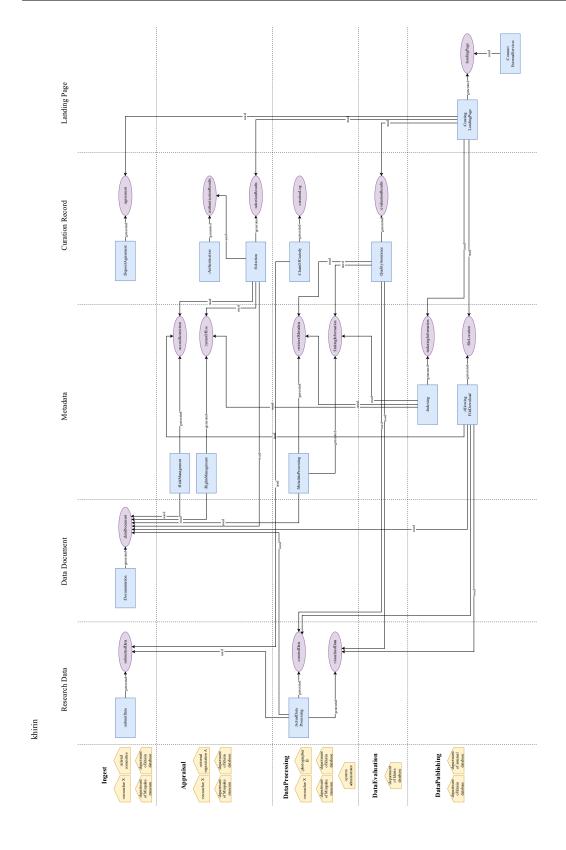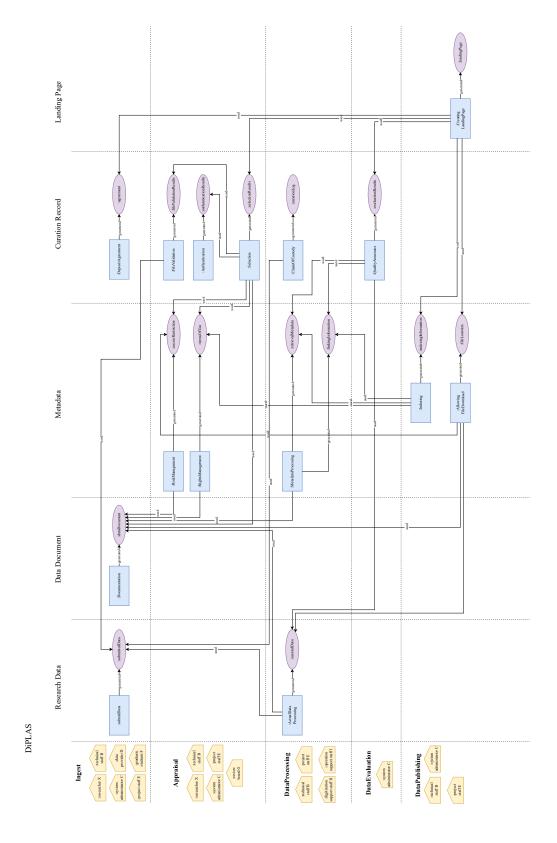Figure 3.5: khirin data curation activity's structure.

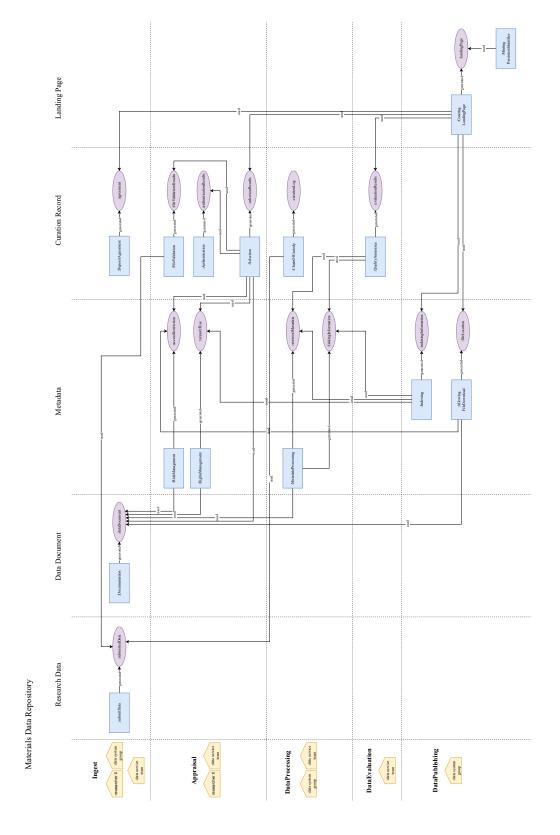Figure 3.6: DiPLAS data curation activity's structure.
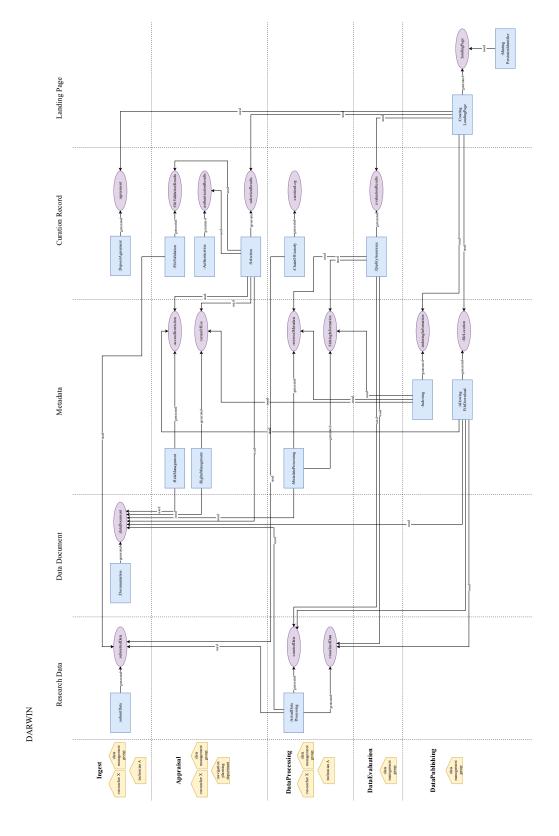
Figure 3.7: MDR data curation activity's structure.

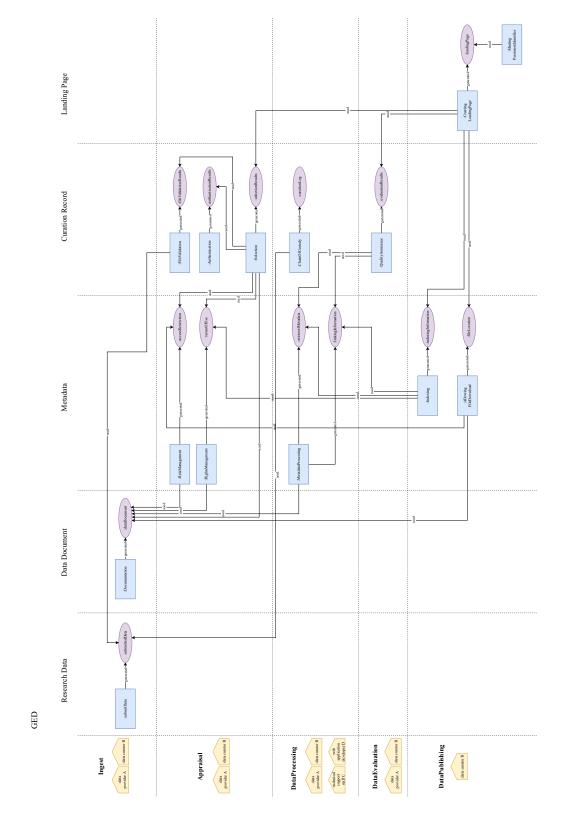Figure 3.8: DARWIN data curation activity's structure.

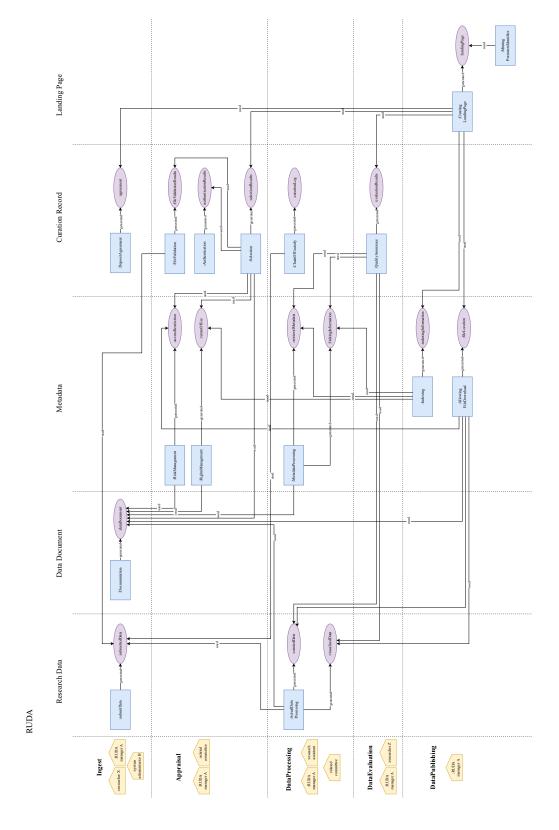Figure 3.9: GED data curation activity's structure.

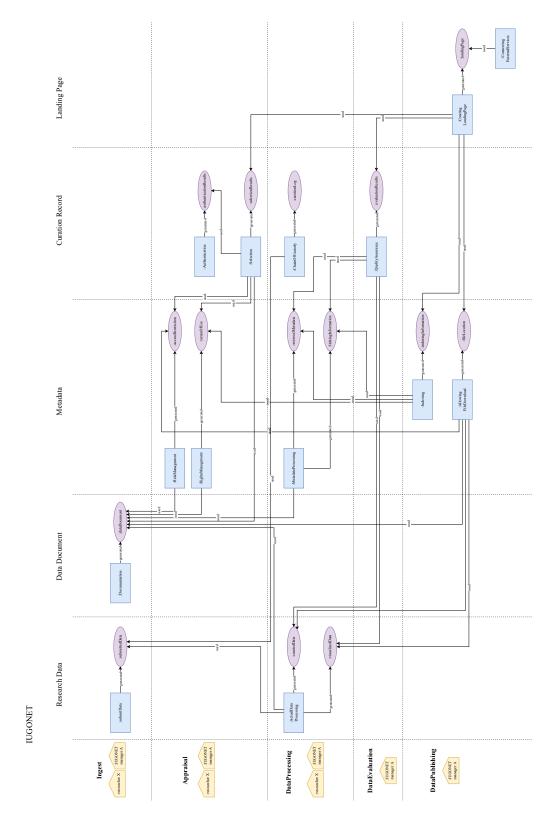Figure 3.10: RUDA data curation activity's structure.

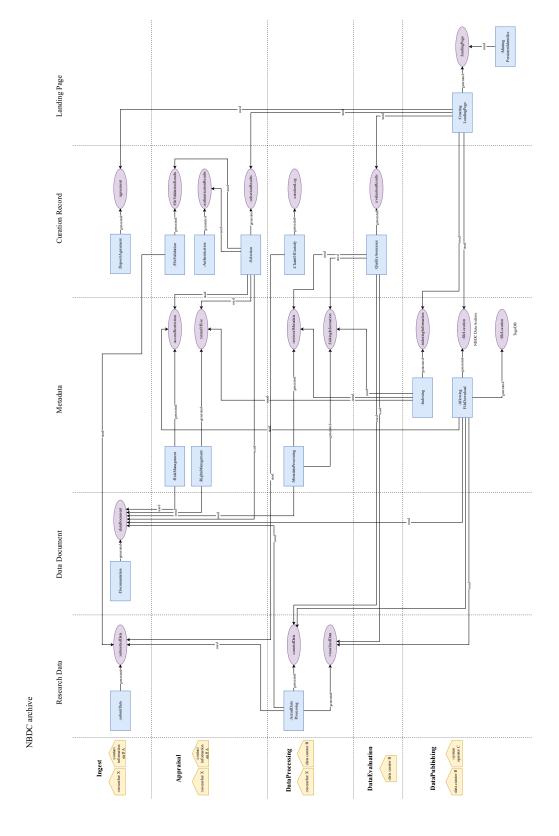Figure 3.11: IUGONET data curation activity's structure.

Figure 3.12: NBDC archive data curation activity's structure.

The essential information of the figure is the graph structure consisting of activity, entity, and agent. For ease of understanding, we arranged the graph with the following rules. The columns in each Figure show five typical Entities: "Research Data," "Data Document," "Metadata," "Curation Record," and "Landing Page." The rows in each Figure show the categories of "Ingest," "Appraisal," "Data Processing," "Data Evaluation," and "Data Publishing." A corresponding data curation process and the generated entity are placed at the intersection of the categories and the entities. The generated entity is connected to another data curation process in which the entity is used by a "used" line. Also, we described agent information on the horizontal axis in each Figure. Agent should be associated with each activity in PROV ontology. Since there are many agent-activity linkages, we described agent information in the simplified form; the agent linked to the activity is described at the left-most column on the same row.
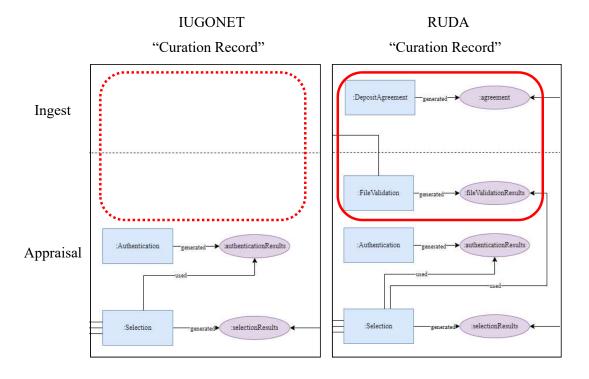


Figure 3.13: Comparison of data curation activity's structures in different fields.

As such, this format can describe the data curation activity's structures in multiple fields in a single model. First of all, this format enables users from different fields to trace the provenance information of each entity. More consistent with the intent

of this chapter, it also contributes to comparing the commonalities and differences between the data curation activity's structure in different fields by describing them by this ontology. Figure 3.13 shows an example of the "Curation Record (Entity)" comparison between IUGONET data curation activity's structure in Figure 3.11 (left) and RUDA in Figure 3.10 (right). The comparison shows that there are no "DepositAgreement" in the Ingest category and no "FileValidation" in the Appraisal category in the IUGONET data curation processes. The reason why these process has not been implemented in IUGONET is that the IUGONET is a metadata distribution service that relies on the data provider for data access. There is no need to verify the data or obtain permission for publication. Therefore, the "Authentication" process is positioned as more important duty of the data curator in terms of comparison with other fields.

## 3.7 Application: Specification of ontology-based data curation process support functions

In this section, we discuss the specification of a data curation process support function using this ontology. Table 3.6 shows the mapping to the functions possessed by the repository software WEKO3 (https://rcos.nii.ac.jp/en/service/weko3/), which is widely used in Japan.

Table 3.6: Functional mapping with WEKO3.

| Category | Data Curation Processes | Function name (WEKO3) | Shareable process flag | Remarks |
|---|---|---|---|---|
| Ingest | SubmitData | Item registration | * | |
| | Documentation | (No function) | * | |
| | Deposit-Agreement | (No function) | * | |
| Appraisal | Authentication | Log-in | * | |
| | FileValidation | Item registration | * | Except for file format validation |

*Continued from previous page*

| Category | Data Curation Processes | Function name (WEKO3) | Shareable process flag | Remarks |
|---|---|---|---|---|
| | RightsManagement | (No function) | * | |
| | RiskManagement | (No function) | * | |
| | Selection | (Partly) Workflow | * | Except for selection criteria support |
| DataProcessing | ActualDataProcessing | (for journal article) Cover page creation | * | Required processes vary by field |
| | ChainOfCustody | Workflow | * | |
| | FileInventoryOrManifest | (No function) | * | |
| | MetadataProcessing | Item registration/Item linking | * | |
| DataEvaluation | CodeReview | (No function) | | |
| | PeerReview | (No function) | | |
| | QualityAssurance | Item approval | * | |
| DataPublishing | ActivatingMetadataBrokerage | OAI-PMH harvesting / ResourceSync | * | |
| | CreatingLandingPage | LandingPage displaying | * | |
| | GeneratingFulltextIndexing | Full-text indexing | | |
| | Indexing | Index creation | * | |
| | AllowingFileDownload | Download URL creation | * | |
| | ConnectingDiscoveryServices | OAI-PMH harvesting / ResourceSync / Google Scholar metadata / schema.org | * | |

*Continued from previous page*

| Category | Data Curation Processes | Function name (WEKO3) | Shareable process flag | Remarks |
|---|---|---|---|---|
| | MintingPersistentIdentifier | DOI registration / CNRI handle | * | |

*End of table*

Table 3.6 shows the WEKO3 functions corresponding to the Data Curation Process Ontology and the flags indicating the interpretable processes by multiple organizations analyzed in Section 3.4.2 as the shareable processes. WEKO3 is a data publishing platform for researchers to publish research data and related materials. WEKO3 supports basic data registration routes such as "SubmitData" and "FileValidation," and also supports a wide range of metadata registration, editing, and publishing functions such as "MetadataProcessing," "ChainOfCustody," "QualityAssurance," and "DataPublishing." Whereas, WEKO3 does not support some shareable processes related to data itself in each field such as "Documentation," "RightsManagement," "RiskManagement," "Selection," and "ActualDataProcessing." It is recommended to provide the support functions for these processes to support data curation in different fields by WEKO3 comprehensively.

Thus the formalized processes enable analysis system functions to identify the shareable data curation processes in different fields. Moreover, the ontology supports integrating with external tools to handle missing processes. Data curators often use multiple data curation systems, then the ontology can help them to allocate support functions among the systems.

## 3.8 Results and remaining issues in chapter 3

As the first step to interpret the tasks and procedures performed in different fields at the same granularity, we investigated the practices of data curation conducted in each field. As a result, we found that about 87.2% of the processes are interpretable across multiple fields. Also, we realized that there needs a suitable model to describe the structure such as the relationships among Input-Output objects, processes, and staffing to accurately represent the data curation activity's structure in different fields.

Based on the preliminary analysis and survey results, we developed the Data Curation Process Ontology to formalize the data curation activity᾽s structure in different fields. To verify the usefulness and validity of this ontology, we described and compared the several actual data curation activity᾽s structures. Users can visualize the data curation activity᾽s structure and the provenance information of the research data by using the format presented in Section 3.6. It is the important contribution of this study to compare the activity᾽s structure of eight diverse repositories in a single model. Also, we showed that the ontology can use the specification of data curation process support functions by systems. This ontology can also help to allocate data curation process among multiple systems and improve the data curation process during data integration in different fields. With the formalized method enabled by the ontology, we can expect to expand the data curation process in various fields from an interdisciplinary perspective. Thus this study contributed to building a knowledge framework for a common understanding of the data curation process in different fields.

A possible direction for utilizing ontology is to evaluate the quality of research data. The ontology provides a well-formalized structure for the data curation activities and may function as a framework for the process management of research data. To achieve process management of research data, it is necessary to accumulate case study practices based on a formalized process definition at an application level. Technical implementation in terms of formalized processes will help validate and solidify the framework.

# 4

# Practical interpretation of the interdisciplinary data evaluation processes using the reference model for data publishing

In Chapter 4, we deal with the processes in the Data Evaluation category formalized in Chapter 3 through practical implementations. To interpret research data from a field to be independently, the data and their documentation must be reviewed from an interdisciplinary perspective and revised as necessary. In contrast, many data repositories either do not have a policy for evaluating research data and encouraging improvement, or the evaluation process is done only from the perspective of a specific field. This chapter focuses on publishing mechanisms for data papers that divert the journal peer-review and clarifies the mechanism from a process perspective. Furthermore, we conduct a technical implementation for interpreting the mechanism

using the reference model for data publishing. This study is a practical interpretation of the relationship and boundaries between the existing quality assurance process and the peer-review process in the Data Evaluation category.

## 4.1   Introduction

In Chapter 4, we deal with the processes in the Data Evaluation category formalized in Chapter 3 through practical implementations. The Data Evaluation category includes Code review, Quality Assurance, and Peer-review processes. Generally, these processes review the quality of the research data and associated documentation and take the necessary actions to make the subject independently interpretable (Peer et al., 2014; Consultative Committee for Space Data Systems, 2012). The target objects are data files, data documents, metadata, codes, etc., and may include data papers (Lawrence et al., 2011). The data evaluation is essential in maintaining control over scholarly communication.

The data evaluation processes have been put into practice in some leading fields. These experiences have shown several gaps and challenges in operationalizing the evaluation process. First, many data repositories either do not have a policy for evaluating research data and encouraging improvement, or the evaluation process is done only from the perspective of a specific field (Peer et al., 2014). Dataset quality seems to be a mixed bag due to the diversity of the data repository's policy (Ruggles, 2017). Also, describing metadata is time-consuming and often not recognized as a research contribution (Edwards et al., 2011). As a result, it lacks positive motivation for researchers. This situation is one of the primary reasons for the lack of progress in depositing research data even in the leading fields.

Mechanisms that include incentives for researchers regarding data curation are needed to overcome challenges in the data evaluation process. In recent years, peer-review of data using the media type of data papers has come into the spotlight (Hrynaszkiewicz and Shintani, 2014; Assante et al., 2016). Traditionally, peer-review has been the standard mechanism across all fields as a framework for accrediting and evaluating research papers (Spier, 2002). Peer-review of data papers could be a more effective mechanism for evaluating research data.

This chapter focuses on publishing mechanisms for data papers that divert

the journal peer-review and clarifies the mechanism from a process perspective. Furthermore, we conduct a technical implementation for interpreting the mechanism using the reference model for data publishing. This study is a practical interpretation of the relationship and boundaries between the existing quality assurance process and the peer-review process in the Data Evaluation category.

## 4.2   Literature review

### 4.2.1   Peer-review of data

Peer-review of data is not yet a clearly defined concept; its meaning varies depending on the peer-reviewed scenarios (Mayernik et al., 2015). There are three main scenarios in which data may be subject to peer-review: 1) when data is published in a data repository, 2) when data is published as an appendix to a paper, or 3) when data is published as a data paper (Lawrence et al., 2011; Mayernik et al., 2015). These scenarios can be positioned as a framework for data publishing (Candela et al., 2015; Pampel et al., 2012), with variations depending on the division of roles among researchers, repository managers, and publishers involved in publishing data (Peer et al., 2014; Lawrence et al., 2011).

1) When data is published in a data repository, the data curator reviews the data in terms of file interpretability and accessibility, and long-term preservation (Peer et al., 2014). This process includes dataset dimension checks, validity checks, confidentiality checks, metadata checks and enhancement, and format transformation checks (UK Data Archive, ndb). As already mentioned, many data repositories either do not have a policy for evaluating research data and encouraging improvement, or the evaluation process is done only from the perspective of a specific field (Peer et al., 2014).

2) When data is published as an appendix to a paper, the peer-reviewer reviews to substantiate the statements contained in the paper (Grootveld and van Egmond, 2012). The items checked to depend on the journal's peer-review policy. In most cases, data documents and metadata are not subject to peer-review, making it difficult to interpret the data independently.

3) When data is published as a data paper, the peer-reviewer reviews from the

perspectives of both 1) and 2). The data paper describes mainly ancillary information, including where the research data is stored. Unlike ordinary research papers, a data paper rarely describes some findings or interpretations derived from research data (Candela et al., 2015). The peer-reviewers will review the data documents described in the data paper and confirm that the data have been processed as described. In this case, the challenge is establishing a system to realize the peer-review of data.

## 4.2.2 Data paper in data publishing

The following section will then provide an overview of the mechanisms that make the peer-review process of data papers available. Reflecting the recent open science trends, many data papers are published in an open access form (NISTEP, 2015). Dataset tied to the data paper can be easily utilized by a wide range of research fields now. Having said that, the concept of data publication itself is not new. There have been many attempts to publish research data in reports or chronologies (Beagrie, 2008). While the open access initiative has given new meaning to the data publication (Klump et al., 2006), data journals that mainly deal with data papers have emerged as a unique media format (Candela et al., 2015). Data journals handle research data using the same business model as a typical research paper (NISTEP, 2015). This means that the publication functions of registration, certification, awareness, archiving, and rewarding are applied to research data (Van de Sompel et al., 2004). According to this perspective, the following benefits can be expected.

**a) Enriching description regarding research data**

The authors can describe more detailed information in the data paper compared with existing research papers. Table 4.1 shows the general elements of the data paper's description (Candela et al., 2015). Note that the precision and emphasis of the required descriptions vary from journal to journal. These descriptions include information that is rarely described in general research papers. For example, "Format," "Microattribution," "Provenance," and "Reuse" are likely to be excluded from general research papers.

Increased information could simply improve reproducibility and ensure research transparency. Also, project-isolated and not yet cited research data can be stored and

traced step-by-step on a common platform via a data paper. It is expected to preserve and reuse datasets as research contribution, 80% of which are said to be lost in 20 years (Vines et al., 2014).

Table 4.1: Data paper description style. Source: Candela, L et al., 2015

| No | Element | Description |
|----|---------|-------------|
| 1 | Availability | Name of the data provider, DOI/URI, etc. |
| 2 | Competing interests | Individual/Organizational relationship related to the dataset |
| 3 | Coverage | Spatial/Temporal |
| 4 | Format | File format, Encode, Language, etc. |
| 5 | License | – |
| 6 | Microattribution | Contribution type related to the dataset development |
| 7 | Project | – |
| 8 | Provenance | Including acquisition methods, equipment, etc. |
| 9 | Quality | Information on data limits, outliers, etc. |
| 10 | Reuse | Potential Value |

**b) Ensuring the quality of research data**

The conceptual process of data publishing is shown in figure 4.1 (Kratz and Strasser, 2014). Authors store the dataset in a repository and submit a detailed description of the dataset as a data paper. After the prescribed procedure, the publisher checks the submitted dataset and data paper and assigns an identifier to both.

As mentioned above, data papers submitted to data journals go through the same third-party peer-review process as existing research papers; ensuring quality through peer-review improves interpretability. Note that the meaning of "quality" of a data paper differs from that of an existing research paper. In a data paper, a peer-reviewer does not evaluate "highly novel" but rather "accurately and richly described according to a certain descriptor."

**c) Expansion of distribution channels**

Since data publishing maintains the form of paper, data papers can be published on the same platform as existing research papers. Also, data papers will be given the same metadata as existing research papers. This metadata can improve research
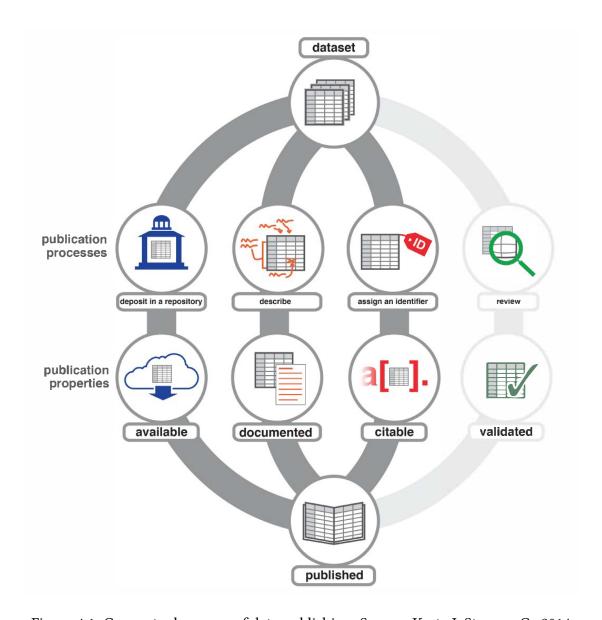
Figure 4.1: Conceptual process of data publishing. Source: Kratz J, Strasser C., 2014

data discoverability through searching data papers; the search route via data papers avoids complex interdisciplinary use of metadata problem across different fields (Willis et al., 2012). Moreover, from the data user's perspective, they can search the research data centralized with a familiar search platform. Publishing data paper leads to an expansion of the distribution channel.

**d) Provide incentives for data curation**

A high-quality description of research data will be recognized as a research contribution to evaluating the data paper's description by peers. This direction may motivate researchers to create and preserve higher-quality data. Also, research data have traditionally not been subject to citation (Robinson-García et al., 2015) so the data creators have not been fully recognized. In the data publishing norm, the data creators will be the author of the data paper. The data creators will evaluate their work as a research community's contribution (Lawrence et al., 2011).

## 4.2.3   A reference model for data publishing

Traditional data deposition flow is accomplished in a manner tied to a repository, making it difficult to find a clear and consistent human-readable workflow representation for the repository (Austin et al., 2016). As one of the few examples of how this issue is being addressed, the RDA/WDS Publishing Data Workflows WG has analyzed 25 journals/repositories/guidelines and is developing standard workflows.

Figure 4.2 shows a framework for data publishing discussed in the RDA/WDS WG (Austin et al., 2016). No. 2-1 shows the traditional publication workflow model. In this model, both ensuring access and creating an explanation of the research data are left to the researchers. No. 2-2 shows the data publishing workflow model with "data publication." In this model, researchers should describe additional information such as methods, detailed descriptions, and related computer codes. No. 2-3 provides a detailed flow of the data publication and article publication process as defined in No. 2-2. "Data repository submission" shows the checking mechanism for mutual use and reuse of actual data. The "data article submission" shows the process of reviewing data papers. Data papers are published in scholarly journals using the same peer-review process as traditional research papers. The data paper contains link information about
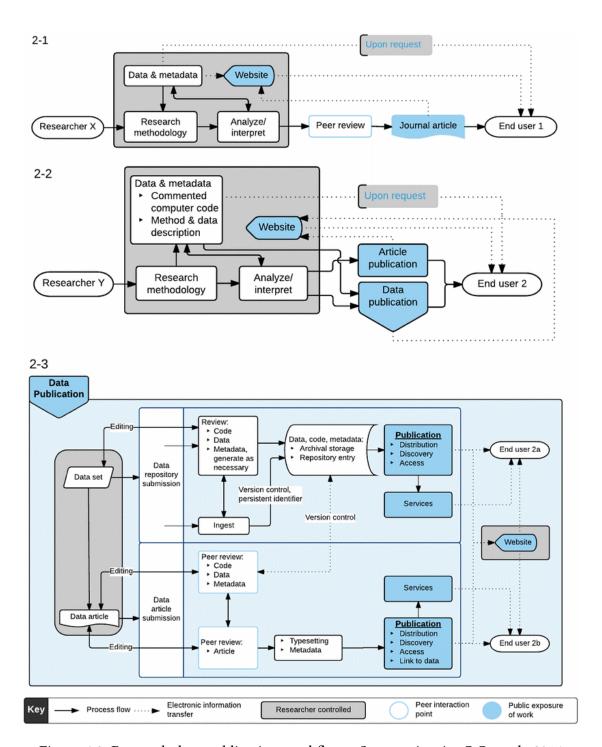
Figure 4.2: Research data publication workflows. Source: Austin, C.C. et al., 2016

the actual data and will be published on the journal platform in the same manner as the research paper.

"Data repository submission" and "Data article submission" can be evaluated as flow representations of the quality assurance process and the peer-review process formalized in Chapter 3. We have positioned this model as an effort consistent with a practical interpretation of the formalized evaluation process.

## 4.3 Case study

As detailed in Figure 4.2, the data publishing flow for data paper combines a quality assurance process and a peer-review process. Separating the two evaluation process will help make the system lighter and more flexible. However, there are few accumulated examples of application designs that follow this reference model. Therefore, we planned to implement some data curation support systems with each process as an application boundary. In this section, we introduce a case study in polar science field for referring to this reference model.

### 4.3.1 Context

The NIPR has acquired a wide variety of polar scientific data through its historical Antarctic and Arctic Research program. The acquired data are shared internationally, contributing significantly to the development of the field. The NIPR has published two reports: JARE Data Reports and NIPR Arctic Data Reports for international data sharing (Matsuzato, 1992).

Since the 2000s, NIPR started to operate some data repositories: the "NIPR Science Database" in 2007 and the "Arctic Data Archive System" in 2012. These data repositories realize the distribution of research data regardless of media format, but they have similar operational issues as mentioned in Section 4.1. Especially, the quality assurance of research data to be published under the responsibility of the NIPR was a critical issue for NIPR data curators. To address this issue, the data journal was proposed as a mechanism for both data distribution and quality assessment. The NIPR launched the Polar Data Journal (https://pdr.repo.nii.ac.jp/), the first data journal in a Japanese academic institution in January 2017. 32 data papers have been published as of March

2022.

## 4.3.2 Policy design

Since data journals are the same media format as existing journals, there are no significant differences in policy design considerations. Polar Data Journal is published as a free-access online journal, inheriting the characteristics of a general data journal. The license granted to the data paper is either the Creative Commons Attribution 4.0 International License (CC-BY 4.0) or the Creative Commons Attribution - Non-Commercial - No Derivative Works 4.0 International License (CC-BY-NC-ND 4.0). Also, DOIs (digital object identifiers) are assigned to data papers to enhance reusability.

## 4.3.3 Establishment of data publishing flow

According to the RDA/WDS reference model, research data and data papers need to be managed in parallel under version control. To achieve parallel flow management with multiple stakeholders, we represented each stakeholder' s data curation process in a flow diagram. The parties involved in the workflow are as follows: the author who submit data papers, the referee who scrutinizes the papers, the editor and officer who perform the editing, and the data repository manager. Figure 4.3 shows an overview of the editorial process for the Polar Data Journal. The adopted system is described in detail in Section 4.3.4.

The description of each process is as follows:

(1) Authors submit a manuscript and associated files to EM.

(2) The information entered at the time of submission is sent to the editorial office.

(3) The editorial office confirmed no missing information and that the link to access the actual data is valid.

(4) The actual data in a data repository is copied to the office server; the officer calculates the checksum of the actual data.
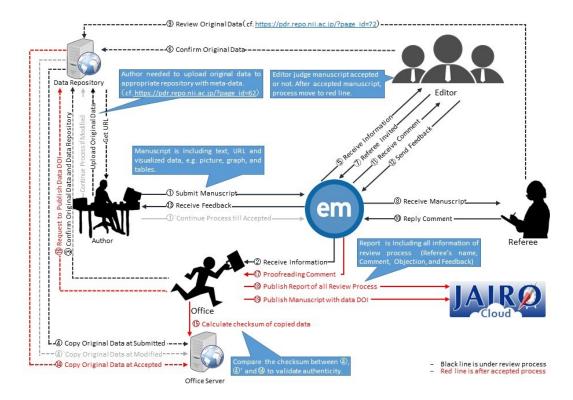
Figure 4.3: The editorial process for the Polar Data Journal.

(5) When the above process (2), (3), and (4) have been completed without any problems, the manuscript is sent to an editor.

(6) The editor checks the manuscript and the actual data.

(7) The editor selects two referees who are experts in the relevant fields.

(8) When they accept the review request, the selected referees become official referees and receive the manuscript and the link to the actual data.

(9) Referees scrutinize the content of the manuscript and prepare a review comment on the logicality of the manuscript, consistency with linked actual data, and the quality of the actual data.

(10) The referees reply to the editor with the acceptance or rejection of the manuscript with the referees' comments.

(11) The editor confirms the referees' comments and decides whether to accept or reject the manuscript as an editorial board member.

(12) The editor sends the decision results to the authors with the proofreading comments.

(13) Authors receive the decision results with the proofreading comments. If there are instructions for revision, authors revise the manuscript based on the comments and resubmit, and the set of processes is repeated.

(14) If accepted, the officer copies the actual data to the officer server again.

(15) The officer re-calculates the checksum of the actual data.

(16) After the checksum matches are confirmed before and after peer-review, the officer requests the data repository to issue a DOI for the actual data through the authors. The DOIs of the actual data will be included in the accepted data papers.

(17) The officer compiles all comments made by the editor and referees during the peer-review process, as well as corrections and rebuttals by the authors.

(18) Comments compiled by the officer will be published with the accepted data paper as a peer-reviewed report.

(19) The data paper will be published on the JAIRO Cloud.

Since authors can modify actual data, there remains the possibility that only the actual data may be changed during/after the peer-review period. If a version change of the actual data occurs during peer-review process, there will be a discrepancy between the description of the data paper and the actual state of the stored data. To avoid this situation, the Polar Data Journal requires authors to store their data statically from submission to acceptance. If authors wish to make additions or changes to the actual data submitted to the Polar Data Journal, they need to submit information about the differences. The mechanism for detecting irregular version changes is described in detail in Section 4.4.

### 4.3.4   System requirements

Next, we discussed system requirements to support each process. Since the data publishing system is expected to work with data repositories in multiple fields, the system must be designed to be flexible. We selected a mainly open-source system that allows flexible development. Conversely, since the peer-review flow is the same process as for existing research papers, we selected systems emphasizing the sustainability of development. As a result, we set the following three system requirements.

**Peer-review process system**

For the peer-review process system, it handles functions such as version control of data articles, assignment of peer-reviewers, management of peer-review schedules and comments, and notification of acceptance/rejection decisions that arise in the course of the peer-review process in addition to data papers and metadata associated with the papers. We adopted Editorial Manager, which is the proven online editorial workflow management system.

**Quality assurance process system**

For the quality assurance process system, it handles functions such as the management of research data and associated metadata, as well as functions such as issuing DOIs to research data and temporary data viewing during the peer-review period. We set the following three criteria for the quality assurance system: 1) Assignment of a persistent identifier (e.g., DOI), 2) Free access, and 3) Open license.

Unlike the previous peer-review process system, we did not adopt our own quality assurance process system and defined only the criteria on the system side. This is due to the following two reasons: a) research data management methods vary depending on the institution and/or field to which the authors belong, and b) that multiple data repositories have already been established within the field. Also, the previous three points must be clearly stated as operational policies.

The editorial office will check the link information provided by the authors and pre-screen the data repository. If the data repository does not meet the above criteria, the authors will be guided to resubmit the research data to the appropriate data repository. Note that this link need not be entirely public during the submission process; it only needs to be accessible to referees.

**Data publishing process system**

For the data publishing process system, it handles functions such as issuing DOIs to data papers that have been selected for acceptance, publishing peer-review reports, and backing up eligible research data. These requirements correspond to the processes in the Data Publishing category formalized in Chapter 3. We adopted the JAIRO Cloud service (Hayashi et al., 2021), developed by the National Institute of Informatics (NII).

## 4.4 Discussion

In this section, we discuss some issues that emerged during applying the model between multiple systems environments, including the possibility of inappropriate data manipulation in peer-reviews and how to prevent it.

**Verification measure during peer-review process**

We discuss inappropriate data manipulation that may occur during the peer-review period. Suppose the actual data in the data repository is altered or modified somehow, and the content differs from the actual data at the time of peer-review process. In that case, the reliability of the content described in the data paper cannot be certified. To ensure reliability, we added three steps in the original model: (4), (14), and (15) shown in Figure 4.4. The editorial office calculates a checksum of the actual data in each phase and compares both. By comparing the checksums, it is possible to verify that no modification or loss has occurred to the actual data in the data repository during the peer-review period.

**Verification measure after peer-review process**

We also discuss inappropriate data manipulation that may occur after the peer-review period. Any manipulations performed during peer-review must also be verifiable by a third party as necessary. Third parties cannot determine whether the data is reliable or not without accurate records; it is crucial for reproducibility to record the data curation process. We added two steps in the original model: (17) and (18) shown in Figure 4.4. The editorial office aggregates the history of each preceding process: recording the data curation process, including the start and end dates and times of each review process, the duration of the review, comments by the editor and referees, corrections and rebuttals by the authors, and the results of the checksum calculations. By publishing the peer-review report, anyone can check whether the peer-review process of the data paper was conducted correctly, which is expected to improve reliability.

## 4.5   Results and remaining issues in chapter 4

This chapter focused on publishing mechanisms for data papers that divert the journal peer-review. First, we evaluated the reference model for data publishing from a process perspective, positioning the model as a combination of existing quality assurance process and peer-review process. Furthermore, we conduct a technical implementation by interpreting RDA/WDS model as a case study. In applying the model, we clarified possible inappropriate activity between multiple systems environments in the data

evaluation process. We developed two verification measures to deal with them. Through this study, we interpreted practically the relationship and boundaries between the existing quality assurance process and the peer-review process. This effort can also be positioned as an application-level interpretation of the processes in the Data Evaluation category formalized in Chapter 3.

We point out a limitation of this study regarding the processes in the Appraisal category formalized in Chapter 3. The tasks and procedures addressed in the data evaluation process are closely related to the risk management process and rights management process performed in the Appraisal category. Since the model does not explicitly mention the risk management process and rights management process, it can be assumed that it is under the researcher's control. However, both processes are assumed to vary greatly depending on the nature of the handled research data. In the case of sensitive data, data curators may require involvement in the early stage. In order to apply this model to different fields with more diverse characteristics, it will be necessary to clarify the relationship between the processes in the Appraisal category and the Data Evaluation category.

# 5

# Reframing the interdisciplinary appraisal processes by analyzing conditions of use of research data

In Chapter 5, we deal with the processes in the Appraisal category formalized in Chapter 3, focusing on the processes related to setting conditions of use. These processes have different variations due to the differences in legal restrictions and disciplinary norms by jurisdiction. Also, the conditions of use granted by data providers are more diverse. As a result, it is difficult for data re-users to interpret the results of the process accurately. This chapter investigates the actual processing status related to setting conditions of use in different fields and clarifies the correspondence of information tied to each formalized process. Furthermore, we conduct a practical implementation for stepwise interpretation of these processing results. This study reframes the processes in the Appraisal category and complements Chapter 4 practices each other.

## 5.1  Introduction

In Chapter 5, we deal with the processes in the Appraisal category formalized in Chapter 3, focusing on the processes related to setting conditions of use. These processes have different variations due to the differences in legal restrictions and disciplinary norms by jurisdiction (Ball, 2014). Also, the conditions of use granted by data providers are more diverse (Kindling et al., 2017). As a result, it is difficult for data re-users to interpret the results of the process accurately. Research data reuse based on ambiguous interpretation of processing results often leads to unintended reuse for the data providers. The conditions of use must be classified when research data are published in different fields (OECD, 2015a).

To accurately interpret the processing results involved in setting the conditions of use, it is necessary for data re-users to understand the actual status of processing from an interdisciplinary perspective. In recent years, there has been a move toward making data generated from public funds openly available in principle (ARDC, 2021; CODATA, 2019; Tsoukala et al., 2016). However, non-public funded data are not subject to policy implications. Standardization efforts based on the actual situation for conditions of use for data publishing has not been substantially explored.

This chapter investigates the actual processing status related to setting conditions of use in different fields and clarifies the correspondence of information tied to each formalized process. Furthermore, we conduct a technical investigation for stepwise interpretation of these processing results. This study is positioned as an effort to complement chapter 4 practices each other and support the building of an interdisciplinary infrastructure from a rule perspective.

## 5.2  Literature review

Following factors have been identified as occurring risks in research data sharing: lack of trust for proper use by data re-users, overly restrictive policies and unclear guidelines on research data sharing, and confusion over the ownership of data (Van Panhuis et al., 2014). The data curation activity involved in addressing these factors can be divided into two main processes: One is Risk Management, the process of reviewing data for known risks such as confidentiality issues inherent to human subjects data,

sensitive information. This process also includes taking actions to reject or facilitate remediation when necessary. Another one is Rights Management, the process of tracking and managing ownership and copyright inherent to a data set as well as monitoring conditions and policies for access and reuse. As a result of the execution of these two processes, conditions of use are set.

Among the setting conditions of use, unclear copyrights and licenses are often a high barrier among the list of factors that hinder data reuse (Mayernik et al., 2020). A large international survey found that unsure about copyright and licensing (37%) are the second most crucial factor as the obstacles to publishing research data after organizing the data in a useful manner (46%) (Stuart et al., 2018). As a complementary report to the difficulty of interpreting conditions of use, the most widely used type of condition of use registered in re3data.org is "Other" at 57.2% and "Copyright" at 38.6%, with setting its own conditions of use or copyright notices being the most used situation (Kindling et al., 2017). The use of standardized tools, such as Creative Commons (CC) licenses, is limited to 21.8% at most. In some cases, CC licenses are granted for "non-copyrighted" works. The conditions of use for research data required by data providers are very diverse; hence, the interpretive costs for reuse are significant.

Initiatives for setting appropriate conditions of use have already been set up. The Digital Curation Centre mentioned early on the importance of promoting licensing as a way of maximizing the economic and social impact of data publishing (Ball, 2014). A survey conducted in 2017-2018 reported that 20 initiatives had been set to address ethical or legal issues (Grabus and Greenberg, 2019). Among these initiatives, the RDA/CODATA Legal Interoperability Interest Group (IG) published the article "Legal interoperability of research data: principles and implementation guidelines" (RDA-CODATA, 2016). These guidelines are primarily for data produced in or funded by the public sector and focus on legal interoperability to address misunderstanding and lack of knowledge and guidance on the legal issues generally related to research data.

In some fields, data that are not funded by public funds are widely used in research. In the case of non-public-funded data, the data are held by an individual or a company, and the conditions of use are determined by the data provider, except in cases provided by laws or regulations. A number of existing licensing tools can be applied to these data.

Table 5.1 shows a comparison of the CC and Open Data Commons (ODC) licenses and the "Government of Japan Standard Terms of Use" used for Japanese government websites.

Table 5.1: License comparison chart.

| License | Permissions | | | Requirements | | | Prohibitions |
|---|---|---|---|---|---|---|---|
| | Reproduction | Distribution | Derivative Works | Notice | Attribution | Share Alike | Non-Commercial |
| CC0 | X | X | X | | | | |
| CC-BY | X | X | X | X | X | | |
| CC-BY-SA | X | X | X | X | X | X | |
| CC-BY-NC | X | X | X | X | X | | X |
| CC-BY-ND | X | X | | X | X | | |
| CC-BY-NC-SA | X | X | X | X | X | X | X |
| CC-BY-NC-ND | X | X | | X | X | | X |
| ODC-PDDL | X | X | X | | | | |
| ODC-BY | X | X | X | X | X | | |
| ODC-ODbL | X | X | X | X | X | X | |
| Government of Japan Standard Terms of Use | X | X | X | X | X | | |

These licenses are designed to provide flexible conditions of use, considering the copyright law. However, the current laws and regulations that may apply to research data are more complex and do not fully reflect the actual conditions. In some cases, CC licenses are granted to data that are not subject to protection under the Japanese copyright law, such as numerical data, which may cause confusion about reuse. While there is no doubt that CC licenses are still a useful solution in many cases of data publishing, some challenges exist for handling non-copyrighted data. Some proposals have recently been pushed forward (e.g., UK Scholarly communication license (Baldwin and Pinfield, 2018) and Microsoft Open Use of Data Agreement (https://github.com/microsoft/Open-Use-of-Data-Agreement)). Meanwhile, the requirements for legal decisions on privacy and other matters slightly differ from country to country. For example, fair use in the U.S. have not been introduced in Japan, and alternatives are still under consideration. The process of setting conditions of use

needs to be considered in the local context.

## 5.3 Survey

This study aims to investigate the actual processing status related to setting conditions of use in different fields and clarifies the correspondence of information tied to each formalized process. To achieve this purpose, we set the following research questions (RQs) herein.

RQ1: What are the limitations that arise in using the research data?

RQ2: What do the data providers desire or request when they publish research data?

RQ3: What support can be effective in promoting reuse in aspects of the data providers' desire or data re-users' request?

The conditions of use of research data are a complex combination of external constraints and the data provider's requests. We start clarifying the two boundaries based on the use-case of actual conditions of use in various fields.

### 5.3.1 Interview survey

First, we conducted the interview survey with data repository practitioners to organize the conditions of use that could be granted according to the types of research data. Since the data provider's requests were assumed to be more diverse, we aimed to identify external constraints across fields as a first step. We also aimed to obtain clues for the subsequent questionnaire survey.

In selecting interviewees, we targeted data repositories with data submission policies and researchers knowledgeable about data policies. The reason was that data policies generally include items that stipulate the conditions of use, and we assumed that we could collect information on the external constraints behind the enactment. As a result, the survey included five experts, including data repository managers and researchers from space science, environmental sciences, social sciences, materials science, and humanities. Note that the interview study was limited to Japan, since we

think that external constraints must be judged in the local context.

The survey period was from December 12th, 2017 to February 1st, 2018 and was conducted for approximately one hour each. In conducting the survey, we stated the research request document before the interview and obtained permission. We used a topic guide to share our interview outline with the interviewee. We set three major sections within the topic guide: "Sharing and publishing of research data," "Regarding granting licenses to research data," and "Licensing of research data and promotion of legal interoperability." The topic guide used for the actual questions is shown in Appendix 4. Table 5.2 shows the summary of the results.

We coded in terms of data characteristics, data holders, requests for users, penal regulations, and rights protections that may affect Risk Management Process. We also provided columns to describe issues and aspirations related to conditions of use.

Typical data holders are the researchers from whom the data were obtained, their affiliated institutions, research funders, and third-party data suppliers. In various cases, it is unclear who can claim rights because of the passage of time or the circumstances of the funding agency. Penalties for violations have not been established or strictly enforced. As the scope of data sharing becomes more expansive, it seems to be regarded as less effective. The demand for rights protection varies by research field. In a research field that deals with both constrained and unconstrained data, there seems to be that the less constrained datasets tend to be more commonly used. Challenges in data repository operations included national security, sensitive data, and fostering a culture to achieve data protection. They also indicated a demand for an interdisciplinary understanding of these issues.

### 5.3.2   Questionnaire survey

Next, we conducted a web questionnaire survey based on the data obtained in Section 5.3.1. This survey aims to take an exploratory look at the data provider's request. We provided following ten questions without mandatory items:

(1)  Which of the following terms best describes your research field?

(2)  Have you ever obtained or published any data, including the cases in which user

Table 5.2: Summary of the interview survey results.

| Date | 2017/12/12 | 2017/12/18 | 2017/12/20 | 2018/1/30 | 2018/2/1 |
|---|---|---|---|---|---|
| Fields | Space Physics | Environmental Sciences | Social Sciences | Materials Science | Humanities |
| Characteristics | Mostly numerical data | Image data, numerical data, etc.. | Survey data | Measurement data, calculation data, Material Informatics data, software | Image data, bibliographic data |
| Data holder (Representative) | The person(s) who is/are acquired the data | The funder(s) | The data provider | Institution of the data holder (with exception) | uncertain |
| Requests for users | (None) | Request for user name and purpose of use for searching metadata ; Primary data are negotiated on an individual basis | Only available by the researchers who belongs in the social sciences discipline ; Submission of a research proposal is mandatory ; Submit usage reports every year / Inclusion in the acknowledgments for using the data | Provide a provenance information of the data as well as literature citations | The metadata should be CC0 ; The conditions of use of the data should be clearly stated by data holder |
| Penal regulations | There are no publishing constraints from a scientific point of view | Under consideration | If the data is passed on to a third party without permission, the use of the data may be suspended | (None) | Considering the use of rightsstatements in case of publishing data |
| Rights protections | No rights protection is provided for the data to be published | There are two levels of access restrictions on data in the repository, depending on the contents | One-year and indefinite licenses are available according to the wishes of the data holder | Data marked as private will have restricted access | Non-public data will be considered for protection on an individual basis by contract. Among the public data, those with copyright properties will be subject to government of Japan standard terms of use and CC licenses |
| Issues | There are few explanation for national security related data (especially the description of disclosure period) | There are no institution to consult about research data rights in Japan | In addition to an organization that supports data management, data management personnel and personnel who can handle the technical aspects of metadata are needed ; The criteria for determining sensitive data change over time, so it cannot provide past data as it is | A point of contact is needed to receive inquiries about published data | Licensing standards for publishing non-copyrighted or obscure data ; How to develop a culture of data protection, who will bear the cost, and how to spread it The future discussion points are whether or not to do so |
| Aspirations related to conditions of use | Development of data utilization laws | Enhancement of the university's intellectual property department function | (None) | (None) | (None) |

registration and fees are required?

(3) Are you familiar with the following license tools?

(4) Have you ever used any of the following license tools to publish your data?

(5) If you would like to publish your data, would you like to require the following to your users (those who use that data to publish results)?

(6) If the license is complied with, would you be willing to publish the data?

(7) If you are using public data for your part of the research, please choose the method of presentation that you think is appropriate.

(8) Do you have any requests or concerns about using your published data for commercial activities, patents, press, literature, art, etc.?

(9) Please select the initiatives that you believe are desirable to the use and publishing of data.

(10) Free description (any problems or requests regarding the use or publishing of data).

In designing the survey, we created questions without the external constraints suggested by the interview survey. The survey period was set from February 13, 2018 to March 20, 2018. We distributed our questionnaire form via some mailing lists and websites using Questant's questionnaire system. The survey targeted mainly researchers, data managers, and librarians. The final number of responses is 413, of which 409 are valid responses. It should be noted that two limitations of this survey are as follows: (a) This survey was not a random selection. (b) The respondents' research fields were biased from Social Science (17.4%) to Astronomy (0.2%). The aggregated results of the 409 valid responses are presented in the order of the questions presented before. The data from the survey are publicly available (Ikeuchi and Minamiyama, 2020). The "n" in the charts indicates the number of respondents.

**(1) Property of respondents**

Table 5.3 shows the research fields of the respondents. Social sciences (17.4%), earth sciences (12.5%), and humanities (10.3%) were well represented among respondents, while mathematics and astronomy were not (both 0.2%). Other responses included library and information science, nursing, nutrition, and so on. In some cases, affiliation information was recorded for library staff and private companies rather than research fields. Fifty-eight (14.2%) respondents selected, "I am not currently engaged in any research activities."

Table 5.3: Research fields of respondents (n = 409).

| Research Field | Number | Ratio |
|---|---|---|
| Social Sciences | 71 | 17.4% |
| Earth Sciences | 51 | 12.5% |
| Humanities | 42 | 10.3% |
| Medicine | 35 | 8.6% |
| Engineering | 32 | 7.8% |
| Computer Science | 20 | 4.9% |
| Biological Science | 19 | 4.6% |
| Agricultural Science | 18 | 4.4% |
| Psychology | 16 | 3.9% |
| Physics | 8 | 2.0% |
| Chemistry | 2 | 0.5% |
| Mathematics | 1 | 0.2% |
| Astronomy | 1 | 0.2% |
| Other | 35 | 8.6% |
| I am not currently engaged in any research activities | 58 | 14.2% |
| Total | 409 | 100.0% |

**(2) Experience in obtaining published data and publishing data by themselves**

In this question, we asked for experience in obtaining published data and publishing data by themselves from the nine sources. The respondent's choices are "Obtain," "Publish," and "None" and set as follows: "Obtain" and "Publish" are multiple selections, and "None" cannot be selected when "Obtain" or "Publish" is selected. Table 5.4 shows the aggregate results.

Highly selected answers as regards "where to obtain" are institutional repositories/data archives (62.3%), government repositories/data archives (48.4%), and personal/research lab websites and blogs (47.9%). Highly selected answers about "where to publish" are institutional repositories/data archives (25.7%), personal/research

Table 5.4: Experience in obtaining published data and publishing data by themselves (n = 409).

| Sources | Obtain | Publish | None | No Answer |
|---|---|---|---|---|
| Institutional repositories/data archives | 62.3% | 25.7% | 29.1% | 1.5% |
| Government repositories/data archives | 48.4% | 1.7% | 46.0% | 4.6% |
| Personal/research lab websites or blogs | 47.9% | 23.5% | 41.8% | 2.2% |
| Supplementary materials (in research paper) | 36.7% | 9.3% | 54.3% | 6.8% |
| Academic SNS services (e.g. Mendeley, ResearchGate) | 32.0% | 11.5% | 58.2% | 6.6% |
| Data repositories/archives in specific field | 28.6% | 8.3% | 64.3% | 5.4% |
| Code sharing services (e.g. GitHub) | 24.4% | 8.1% | 69.2% | 5.1% |
| Repositories/data archives by Commercial company | 18.1% | 1.5% | 73.6% | 7.1% |
| Other data publishing services (e.g. figshare, zenodo) | 12.7% | 3.7% | 79.2% | 6.8% |

lab websites and blogs (23.5%), and academic SNS services (11.5%). Compared to the experience of obtaining data, the proportion of respondents with experience in publishing data is lower.

Table 5.5 presents the results of obtaining public data and having experience in releasing data. Respondents who selected "Yes" for one or more of the items in Table 5.4 are tabulated as having "Yes" experience in obtaining and publishing. Consequently, 84.1% of the respondents had experience in obtaining data, and 46.5% of the respondents had experience in publishing data. One respondent did not respond at all.

Table 5.5: Experience in obtaining published data and publishing data.

| | Yes | No/No response | Total |
|---|---|---|---|
| Obtain | 344 (84.1%) | 65 (15.9%) | 409 (100%) |
| Publish | 190 (46.5%) | 219 (53.5%) | 409 (100%) |

**(3) Awareness of existing licenses**

We asked for the awareness of three licenses, which are well known in Japan to identify the extent to which existing licenses are recognized. To eliminate answers based on fuzzy memories, we also set a link to the license or a page explaining the license in this question form. Figure 5.1 shows the aggregate results.

The highest recognition is for CC license, but less than half (46.9%) of the respondents are aware of it. ODC (19.3%) and Government Standard Terms of Use (15.9%) follow, and both are less than two in ten. The survey respondents would be expected to
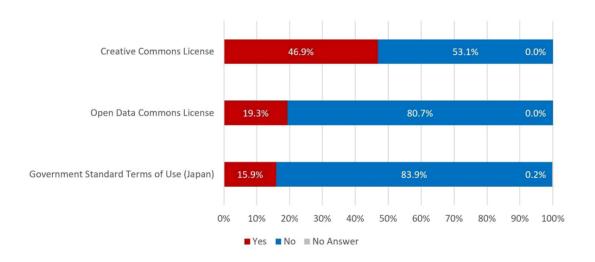
Figure 5.1: Awareness of existing licenses (n = 409).

have some level of interest in licensing their research data, but awareness of existing licenses was not high.

### (4) Usage of existing licenses

To ascertain the use of the licenses listed in the previous question (3), we asked for respondents who are aware of each license about their experience in using each one. Figure 5.2 shows the aggregate results.

Fifty-nine respondents (30.7%) had used the CC license, which was the highest proportion, as was the case with (3). Only four (5.1%) and six (9.2%) respondents had experience using ODC and Government Standard Terms of Use.

### (5) Desired condition of use when respondents publish their research data

We asked respondents to select their desired conditions of use from a list to quantify the extent of the requests they would make. We listed ten items to the possible requests; five items obtained during the interview survey and additional five items included in the CC license elements. Figure 5.3 depicts the aggregate results, with the following order: the sum of "Yes" and "It depends on cases" is the highest.

The aggregate results revealed a diverse reality; all categories were used to some degree. However, a certain degree of difference between preferred and non-preferred
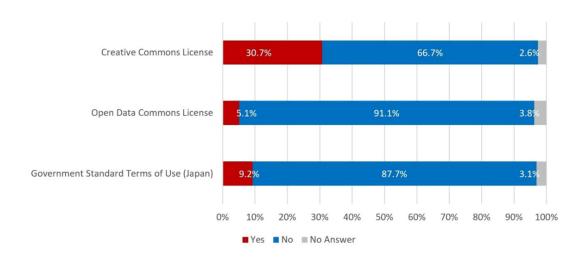
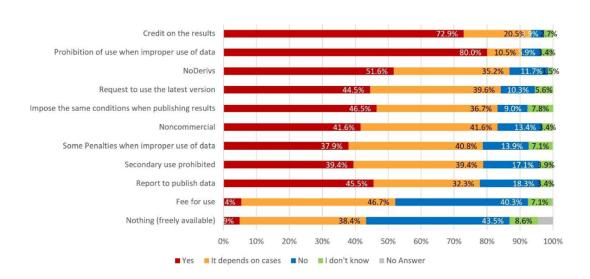Figure 5.2: Usage of existing licenses (n = 409).



Figure 5.3: The desired conditions of use when respondents publish their research data (n = 409).

Figure 5.4: License compliance and willingness to publish data (n = 409).

categories can be observed. The highest percentage of "Yes" and "It depends on cases" is for "Credit on the results" (93.4%). The "Yes" percentage was higher than the credit indication (80.0%) for the "Prohibition of use when improper use of data." The total of "It depends on cases" was 90.5%. The items for which the total of "Yes" and "It depends on cases" exceeded 80% were "Request to use the latest version" (84.1%), "Impose the same conditions when publishing results" (83.1%), and "Noncommercial" (83.1%).

On the contrary, 43.5% of the respondents selected "No" when they asked for "Nothing (freely available)." In other words, just over 40% of respondents wanted to set some kinds of conditions to release their data. In addition, 40.3% of the respondents answered "No" to the question of "Fee for use," indicating that a certain number of respondents do not want to be compensated for data publishing. Note that there is an error of 0.1% between the Figure 5.3 and the main text due to rounding off numbers after the decimal point.

## (6) License compliance and willingness to publish data

We asked whether they would be willing to publish their own data if the conditions of use listed in (5) were complied with. This question was asked to all respondents, including those whose data had already been exposed. Figure 5.4 depicts the aggregate results.

Consequently, 64.1% of the respondents said that they were "Agree," and 24.4% said they were "Somewhat agree" (total: 88.5%), exceeding "Somewhat disagree" (4.6%) and "Disagree" (2.4%).

**(7) An appropriate method of displaying the use of published data**

We asked the respondents about the display method they thought was appropriate for using published data by a third party. This question allows multiple choices. Table 5.6 presents the number and percentage of respondents who chose each option. Note that three respondents, who did not select any of the options, were excluded from the tabulation.

Table 5.6: Appropriate methods of displaying the use of published data (n = 406).

| Choices | Numbers | Rates |
|---|---|---|
| Cite the source of the data in a paper (include it in the bibliography) | 367 | 90.4% |
| Include source of the data information in the main text | 224 | 55.2% |
| Include source of the data information in the acknowledgment | 99 | 24.4% |
| Add the data provider name as a co-author | 47 | 11.6% |
| It is not necessary to describe the data in a paper | 0 | 0.0% |

The highest selection rate was "Cite the source of the data in a paper (include it in the bibliography)" (90.4%). Most of the respondents judged that it would be appropriate to cite the data and the paper if the data were used. 55.2% of the respondents selected "Include source of the data information in the main text" as their next choice. None of the respondents selected, "It is not necessary to describe the data in a paper."

**(8) Requests and concerns about data reuse**

We asked the respondents, "Do you have any requests or concerns if the data you've published will be used for commercial activities, patents, press, literature, art, etc.?" in an additional comment space. This question aims to identify any other requests or concerns not raised in the literature or interview survey. As of 197 respondents' major concerns were as follows: citation or indication of authorship (99 respondents), concern about misuse or inappropriate use (35 respondents), and concern about commercial use (14 respondents).

**(9) Desired approach to data use and publishing**

We asked the respondents about their preferred approach to data use and publishing. This question allows multiple choices. The choices were made with reference to the

interview survey results. Table 5.7 presents the number and rate of respondents who chose each option. Note that five respondents who did not select any of the options were excluded.

Table 5.7: Desired approach for data use and publishing (n = 404).

| Choices | Numbers | Rates |
|---|---|---|
| Establishment of standard data licenses (conditions of use) | 312 | 77.2% |
| Development of appropriate guidelines for data licensing | 285 | 70.5% |
| Establishment of a data licensing consultation, support, and management department (organization) | 168 | 41.6% |
| Enabling a license to be specified in the data retrieval system | 155 | 38.4% |
| Development of data rights legislation | 148 | 36.6% |
| Establishment of a governing body for data licensing (external organization) | 95 | 23.5% |
| Nothing in particular | 21 | 5.2% |

The highest rate was "Establishment of standard data licenses (conditions of use)" (77.2%), followed by "Establishment of guidelines for data licenses" (70.5%). Moreover, "establishment of a data licensing consultation, support, and management department (organization)" (41.6%) was selected higher than "establishment of a data licensing management organization (external organization)" (23.5%) as a contact point of data licensing issues.

### (10) Free comments

A total of 84 respondents described the situation in the free comments. Regarding data publishing in general, various issues were pointed out, including inadequate systems, infrastructure, and technical difficulties and concerns about data publishing.

## 5.4   Analysis and design

This section discusses the correspondence between formalized processes and conditions of use based on two surveys in Section 5.3. The two possible reasons for not publishing research data are as follows: 1) external constraints, such as legal restrictions or disciplinary norms, and 2) data provider's request. We clearly separate the two and discuss the Risk Management process and the Rights Management process in relation to each. We also discuss the feasibility of these interpretations by the data re-user.

### 5.4.1   External constraints

In this section, we organize the external constraints regarding when to publish research data from the survey results in Section 5.3.1. Since the Risk Management process has been conducted on a field-by-field basis, individual risk assessments have been conducted for each keyword, such as personal information and intellectual property rights. We re-summarized these keywords into five categories from an interdisciplinary perspective. This categorization allows for a comprehensive understanding of what is considered in the Risk Management Process. Table 5.8 shows its category, definition, constraint subject matter, and some examples. Note that the examples described are not exhaustive.

Table 5.8: List of external constraints.

| Category | Definition | Subject | Example |
|---|---|---|---|
| Discipline agreement and international treaties | Practices and standards in a specific discipline or research community that limit the data publishing. In some cases this is stated as an international treaty, but in others it is not always explicitly stated. | Disciplines and Norms | Convention on International Trade in Endangered Species of Wild Fauna and Flora (CITES) |
| | | | Convention on the Means of Prohibiting and Preventing the Illicit Import, Export and Transfer of Ownership of Cultural Property 1970 |

*Continued on next page*

*Continued from previous page*

| Category | Definition | Subject | Example |
|---|---|---|---|
| | | | Convention Concerning the Protection of the World Cultural and Natural Heritage |
| | | | The Convention on the Protection and Promotion of the Diversity of Cultural Expressions |
| | | | The Nagoya Protocol on Access and Benefit-sharing |
| | | | Recommendation on the Safeguarding of Traditional Culture and Folklore |
| | | | Bereaved family′s request |
| Personal Information | It stipulates the handling of data that can identify individuals. It includes guidelines that define individual policies on anonymization and information disclosure. | Societies | The Personal Information Protection Commission, Government of Japan. "Laws and guidelines" (only in Japanese) |

*Continued on next page*

*Continued from previous page*

| Category | Definition | Subject | Example |
|---|---|---|---|
| | | | Japan External Trade Organization(JETRO). "About General Data Protection Regulation (GDPR)" (only in Japanese) |
| | | | Ministry of Health, Labor and Welfare (Japan). "About research guidelines" (only in Japanese) |
| Diplomatic / National security | Research data pertaining to national security. Data related to the development of weapons of mass destruction, etc. (as defined in the Foreign Exchange and Foreign Trade Act) and defense secrets (the Self-Defense Forces). law), important data that may affect national life (e.g., domestic energy (e.g., location of resources, blueprints for critical equipment, etc.). | State | Japan Society for Intellectual Production. "Security Trade Control Guidelines for Researchers in universities and other institutions of higher education. Revised 2nd ed" |

*Continued from previous page*

| Category | Definition | Subject | Example |
|---|---|---|---|
| Agreements, contracts, Intellectual Property rights | An agreement with a research partner, contractor, etc. that restricts the data publishing in joint research or contract research. | Companies, etc. | Ministry of Economy, Trade and Industry (Japan). "Operation guidelines for data management in contract research and development" (only in Japanese) |
| | | | Ministry of Economy, Trade and Industry (Japan). "Contract Guidelines on Utilization of AI and Data. Data Section" |
| Data Policy | Where the research funder has a policy on limited data sharing for the research to be funded, or where a strategic business decision restricts the data publishing relating to pending industrial property rights or research data where the commercialization of the research results is envisaged. | Institutions | National Institute for Environmental Studies (Japan). "NIES Data Policy" (only in Japanese) |

*Continued from previous page*

| Category | Definition | Subject | Example |
|---|---|---|---|
| | | | Teikyo University (Japan). "Intellectual Property policy in Teikyo University" (only in Japanese) |
| | | | Japan Agency for Medical Research and Development (Japan). "Data sharing policy for realization of genomic medicine" (only in Japanese) |

*End of table*

**1) Discipline agreement and international treaties**

Data publishing is sometimes restricted by the discipline agreement of the field or research community, such as cases in which the data publishing causes harm to the research subject or cases in which the subsequent research activities themselves are severely affected. Although protection policies are established as international treaties in many cases, these policies are not always clearly defined, known, nor applied.

**2) Personal Information**

Data publishing is sometimes restricted to protect personal information by social demand. Its cases include disclosure, transfer, and anonymization restriction by relevant local laws, cross-border rules such as GDPR, and specific globalized guidelines.

**3) Diplomatic/national security**

Data publishing is restricted if the data is related to national security or international relations. These data are strictly operated in the global context, including the conditions

of use.

**4) Agreements, contracts, and intellectual property rights**

Data publishing is sometimes restricted by agreements or contracts. In general, the parties' agreement determines jurisdiction over contracts and intellectual property rights; sometimes the conditions and the embargo period for data publishing are not uniformly determined because it is not a direct matter of concern.

**5) Data Policy**

Data publishing is restricted when there is a conflict with the policies of the data provider's organization or funding agency. There are many possible reasons for data policy restrictions; for example, a funding agency has a policy limiting data publishing for the research to be funded when a patent has been applied or the commercialization of research results is expected. These cases are restricted as an individual strategic decision and similar to previous data disclosure through contracts. The difference is that it is a management decision based on an "Open and Closed Strategy", which is a strategy to handle data by separating what should be released and what should be protected based on the characteristics of the data.

## 5.4.2   Data provider's request

This section discusses the data provider's requests subject to the Rights Management process. In Section 5.3.2, we observed quantitative differences among categories preferred and not preferred by the data providers. Therefore, we qualitatively analyzed the ten categories used in the questionnaire survey to understand the data provider's request diversity. Table 5.9 shows each condition of use analyzed from the data publishing perspective. This categorization will allow us to better understand the data provider's requests and give some direction to the Rights Management Process.

Note that items observed in the questionnaire survey that belong to the "Requests" category are not included in Table 5.9. "Requests" for the public are not legal contracts, but mainly moral matters; hence, no-obligation, prohibition, nor permission easily occur. The violation does not immediately imply termination of use, but data

Table 5.9: List of each condition of use analyzed from the perspective of expected users, duties, and constraints.

| Condition of use | Expected user | Type of duty | Target of constraint | Compatible with CC licenses | Suggested categories |
|---|---|---|---|---|---|
| 1) Waiver | Public | - | - | CC0 | Preferable |
| 2) Credit on the results (CC license term: Attribution) | Public | Obligation | Redistribution | BY | |
| 3) Impose the same conditions when publishing results (CC license term: ShareAlike) | Public | Obligation | Redistribution and combination | SA (only for redistribution) | Available |
| 4) Noncommercial | Public | Prohibition | Redistribution and data processing | NC (only for redistribution) | |
| 5) NoDelivs | Public | Prohibition | Redistribution and data processing | ND | |
| 6) Improper use of data | Public/Specific | Prohibition | Data processing | – | |
| 7) Reporting | Public/Specific | Obligation | Continue to use | – | |
| 8) Secondary use prohibited | Specific | Prohibition | Redistribution | – | Not Preferable |
| 9) Request to use the latest version | Public (latest version only) | Obligation | Redistribution | – | |
| 10) Fee for use | Specific | Obligation | Redistribution | – | |

providers ask data re-users to comply as much as possible with the data provider's
request for the appropriate use of their data. A further study is needed.

**Basis of categorization**

This section discusses three categories shown in Table 5.9. From the data publishing
perspective, an unspecified number of users must be given access, even if under some
limited conditions. We categorized the targets assumed by each condition of use as
"public," "specific," or both. If the conditions of use are only "specific," we
categorized these conditions of use as "Not Preferable" for data publishing.

We also analyze what obligations are imposed in the conditions of use and
summarize the types and targets of these obligations. Consequently, we found two
cases in which restrictions were placed on the "Redistribution of data" and on the
"Data processing or related processes." Restrictions on data processing or related
processes are not desirable for research data reuse. We categorized the rest conditions

of use as "Available" if they have restrictions on data processing or related processes. Conditions of use with no data constraints or only redistribution constraints were categorized as "Preferable." We also discuss the individual conditions of use in more detail below.

**Preferable**

**1) Waiver**

This declaration waives copyright and all other related rights, including any indication of the data provider's name; it can be evaluated to be definitely intended for the public. The cost to the data re-user is minimized because it is consistent with only the legal requirements. However, there seems not to include any incentives for data publishing in this declaration. Moreover, changing the conditions of use in a manner that makes them more stringent is extremely difficult, even when there is a request to prevent unwanted use due to changes in circumstances, such as increased property values of data. Although this declaration is ideal from the viewpoint of research data reuse, the number of data actually published may be limited.

**2) Credit on the results (CC license term: Attribution)**

This condition of use requires displaying the data provider's name and data URL information. The cost to data re-users can be assessed as a negligible level because it remains a "minimal constraint" found in the openness debate (Open Knowledge, nd). We can evaluate this condition of use aiming for the public.

Relatedly, we observed from the questionnaire survey that the data citation expectation is particularly high. Crediting names could be a certain incentive for data providers; although the past practice has mentioned using the data in the acknowledgments or the text, data citation is meant to treat research data as an independent academic artifact.

Whereas, as data-specific concerns, it may be too costly and impractical for the data re-user to deal with many different data sources as the source data for data-intensive science (e.g., machine learning). Describing all credits in the presence of multiple dataset takes time and effort. One of the commenters stated that this should be resolved as a problem with citation and notation methods. Although it is out of the scope of

this study, more discussions are underway at the Data Citation Synthesis Group in
FORCE11 (FORCE11, 2014) and elsewhere.

### 3) Impose the same conditions when publishing results (CC license term: Share Alike)

This condition of use requires the same condition of use under redistribution and
combination of multiple pieces of data. Unlike "Credit on the results," it may
prevent them from combining the data with another dataset that have an incompatible
license. Although it remains the "minimum constraint" in the openness debate for
redistribution, it should be used with caution in research data reuse.

### 4) Noncommercial

This condition of use requires the "non-commercial" use of data. The habit of
prohibiting commercial use is deeply rooted in the academic community, and it seems
unavoidable given the significance of academia's freedom from the society. There
is another point of contention with this condition of use; the criteria for judgment
fluctuate depending on people because "commercial use" is not clearly defined. For
example, some people will have different interpretations of commercial use when
selling visualizations derived from the data. The ambiguity of commercial use was also
pointed out in discussions on copyrighted materials (Creative Commons, 2009). It
should be used with a more clarifying scope of commercial use.

### Available

### 5) NoDerivs (No Derivatives)

This condition of use prohibits data publishing after any modification. Although
opening to the public is not restricted, the cost to the data re-user is relatively high due
to the permitting process. In general, the data will be published for new knowledge
through processing or combination. Furthermore, the effective case in which the
prohibition of modification largely depends on the type of data and the manner of use
(e.g., image data that are practically a work of art). From the viewpoint of research data
reuse, it should be used in a limited manner.

**6) Improper use of data**

This condition of use prohibits "Improper use" of data in research data reuse. This condition of use has the problem that the meaning of "Improper use" is ambiguous. The inappropriate use in a legal context is limited in scope by the laws of each country. The inappropriate use of field conventions may exist, but sometimes, tacit knowledge remains for data re-users in different fields. As a result, this condition of use is difficult to be the basis for any restriction. Unclear conditions of use lead to contraction of usage; it should be used with more specific provisions.

**7) Reporting**

This condition of use requires post-reuse reporting, suggesting that the objective is to know detailed usage practices rather than mechanical access statistics. Although the data can be reused by both public and specific situations, the condition of use is stricter because of the proactive obligation. This usage condition will work well when the data re-user can be identified; in situations where the data re-user cannot be identified, it will remain at the same level of efficacy as a "request."

**Not Preferable**

**8) Secondary use prohibited**

This condition of use prohibits the secondary use of data. It clearly prohibits data redistribution, translation, or adaption and is intended for one-on-one use of its original form. The background for setting this condition of use is a mix of concerns about misuse/responsibility for quality and a desire to understand the data re-users accurately. Although the data processing or combination is not restricted, they cannot be re-distributed and have to be excluded from the definition of data publishing.

**9) Request to use the latest version**

This condition of use is used to limit the research data reuse to the latest one. Data in the past cannot be reused under this condition of use; thus, a large amount of data will be replaced when the latest data are published, resulting in marked costs of data usage. Also, it is impossible to know in advance when the condition will be violated, and the

manner to notify the version update is very limited. The constraints for research data reuse are too strict to be included in the "Not Preferable" category.

**10) Fee for use**

This condition of use requires some fee for data reuse. The requirement of a user fee before data reuse is considered to be out of the scope of a condition of use that assumes that the data will be open to the public. The survey results also suggest that approximately half of the respondents are still uncomfortable with the act of monetizing data. However, given the sustainability of the data repository, monetization may be a major challenge in the future. The fee could be obtained in various ways, including shareware on a request basis, charging through a freemium model, download speed limits, and whether or not ads are displayed. In cases where the data provider itself requires a user fee, under what conditions the fee will be incurred must be clarified.

### 5.4.3   Feasibility of interpretation by the data re-user

The data re-user is required to understand the external constraints behind the granted conditions of use. From the data provider's perspective, compliance with the granted conditions of use leads to safe data publishing. However, as observed in Section 5.4.1, data re-user needs interdisciplinary and specialized knowledge to interpret these constraints or requests. In fact, we can see many concern on the topic of citation, misuse/inappropriate use, and commercial use in Section 5.3.2. The feasibility of interpretation by the data re-user remains an issue.

### 5.4.4   Section summary

This section discussed the correspondence between formalized processes and conditions of use based on two surveys in Section 5.3. In terms of the reasons preventing the data publishing, conditions of use can be divided into two categories: those resulting from external constraints and those resulting from the the data provider's request.

The Risk Management process corresponds to setting conditions of use resulting from external constraints. We re-summarized external constraints into five categories

from an interdisciplinary perspective. It highlighted constraints that should be judged locally and constraints that can be judged globally. This categorization helps to provide a comprehensive understanding of the matters covered by the Risk Management Process.

The Rights Management process corresponds to the setting of usage conditions resulting from the data provider's request. We observed quantitative differences among categories preferred and not preferred by the data providers. Therefore, we qualitatively analyzed the ten categories used in the questionnaire survey to understand the data provider's request diversity. This categorization will allow us to better understand the data provider's requests and give some direction to the Rights Management Process.

We also discussed the feasibility of these interpretations by the data re-user. Data re-users need to interpret interdisciplinary and specialized knowledge that differs from each field's practice. This suggests that even if a data provider in one field accurately describes the conditions of use, the description may not be easy to understand for data re-users in different fields. Publishing research data may not progress without a mechanism to bridge this knowledge gap between the data providers and the data re-users.

## 5.5  Practical implementation

This section discusses practical implementation in promoting research data reuse in aspects of the data provider's requests or data re-users' intention. There is a knowledge gap between data providers and data re-users in understanding the conditions of use, so we need to bridge this gap by implementing some mechanisms. To address this issue, we formalized the data publishing flow with licensing scenarios, including an expert consultation process. Furthermore, we also developed "Guidelines for specifying conditions of use in research data publishing" (Minamiyama et al., 2020) for a common understanding of this flow. A tentative English translation of the document is shown in Appendix 5.

Figure 5.5: Data Publishing flow with licensing scenarios.

## 5.5.1 Formalization of data publishing flow with licensing scenarios

From the analysis in Section 5.4, we mapped the Risk Management process and Rights Management process to the appropriate setting of conditions of use. Based on this understanding, we formalized a data publishing flow with licensing scenarios shown in Figure 5.5.

The flow consists of five steps. First, the data provider identifies the data to be published in Step 1. Next, the data provider confirms the external constraints from an interdisciplinary perspective in Step 2. In Step 3, for the external constraints identified in Step 2, the data provider set the necessary conditions of use for releasing constraint. Step 2 and step 3 are positioned as Risk Management processes. In Step 4, the data provider selects an appropriate data repository for the data judged to be open to the public. Steps 3 and 4 clearly state that expert consultation will be held because expert knowledge with an interdisciplinary perspective may be required for judgment. This process can also be understood as a part of the Quality Assurance process discussed in Chapter 4. Finally, the data provider chooses appropriate conditions of use in Step 5. Step 5 is positioned as Rights Management process. The details for each step are shown below.

**Step 1: Appraisal and selection of data to publish**

In this step, the data provider identifies various data used in the research, which can be curated and made available to the public. There are various types of data publishing motivations; mandated by publishers, funders, institutional policies, and researchers' intention. Also, some data types have already established methods of publication. The data provider is required to objectively understand their data as a prerequisite for the subsequent steps in this step.

**Step 2: Confirmation for legal restrictions/regulations/remarks**

In this step, the data provider considers whether or not the data identified in Step 1 falls under the five categories of external constraints shown in Section 5.4.1: 1) Disciplinary customs, including international treaty, 2) Personal information, 3) Diplomatic/national security, 4) Agreements, contracts, and intellectual property rights, and 5) Data policy.

The data provider extracts possible legal restrictions/regulations/remarks by category. If the target research data has not any category of concern, Step 3 can be skipped.

**Step 3: Release constraint**

In this step, the data provider sets the conditions of use for releasing constraints identified in Step 2. It includes limiting the number of users and setting an embargo period. The terms or periods set out here will be written as special conditions.

Step 3 also provides a route to consultation with experts. Even in cases where legal or disciplinary restrictions are imposed, restrictions may be lifted with appropriate data processing, such as anonymization or data release restrictions. Also, specific decisions on releases often require interdisciplinary expertise. Therefore, the involvement of experts should be envisaged from this step as a part of the Quality Assurance process discussed in Chapter 4.

There is another meaning to the involvement of experts in this step. There is no legal provision for the termination of the protection period for data as there is for copyrighted works. For example, even if the term of the collaboration agreement has expired, the data are not open to the public after any length of time unless the term is clearly defined. To prevent these unnecessary restrictions, this flow provides

explicit steps for lifting the restrictions and encourages data providers to keep them to a minimum.

**Step 4: Select a data repository**

In this step, the data provider selects the repository where they want to publish the data. Well-known repositories/archives are likely to be the first candidate. If there are no representative repositories in each field, then the data provider's institutional repository or a generic repository such as Zenodo (https://zenodo.org/), figshare (https://figshare.com/) would be the following candidates. What should be checked in this step is whether the selected repository does not fall under the external constraints identified in Step 2. To complement Step 3, we describe a route to consult an expert in this step.

**Step 5: Choose appropriate conditions of use**

In this step, the data provider selects the appropriate conditions of use for data re-users and completes the data use covenants that set out conditions. The requests are more diverse than those for copyrighted works, and the situation has not yet been systematically organized. We propose data use covenants to be consistent with those protected by copyright law.

Note that there remains some concern regarding the complexity of conditions of use. Data use covenants can standardize the description of conditions of use to some extent, but they are not always simple. There are high demands for standardization and simple explanations. But having said that, the simple recommendation of open licenses avoiding the copyright problem will not promote research data use. If the data provider has further requests compared with copyrighted works, the data provider should choose a route to determine validity in consultation with experts.

## 5.5.2 Developing the guidelines for a common understanding of formalized flow

This section introduces the development of guidelines for a common understanding of formalized flow. In order to bridge the knowledge gap between data providers and data

re-users, not only must the process be formalized, but the knowledge used in the process must also be shared. We developed the "Guidelines for specifying conditions of use in research data publishing" (Minamiyama et al., 2020) to help researchers and stakeholders understand and make appropriate publication decisions.

These guidelines provide necessary information and examples to consider when publishing research data with an interdisciplinary perspective. A critical feature of these guidelines are the inclusion of local implementation perspectives from experts. It should be possible to fill the knowledge gap by focusing the description on the local realities faced by data users. Therefore, some limitations exist in the versatility of the description. The main examples are listed below:

**Scope of "research data"**    In Step 1, the data provider must identify their "research data" to publish. Although the scope of "research data" differs depending on the field of expertise, we defined "research data" as data that can be managed by digital means and released as research results and does not include physical objects such as samples, specimens, and recording media. In addition, although research articles and software can be treated as research data, this flow does not change or override the established methods for publishing in each content area. For example, CC licenses for paper publishing and GPL or other software licenses for software publishing are not an assumed target of the data handled by these guidelines. If a researcher has received research funding, they should follow the rule of the treatment of research data defined by the funding agency. These guidelines do not apply to such data; hence, their rules should be applied.

**"well-known" repositories**    In Step 4, we recommend registration in well-known repositories. However, practical differences arise in this phase; "well-known" repositories differ by field and country. Although there are some famous registry sites such as re3data.org (https://www.re3data.org/) and FAIRsharing (https://fairsharing.org/), only a limited number of registrations are available in Japan. In light of this background situation, the guidelines provide both these registry sites and the list of recommended domestic data repositories in cooperation with the Japan Data Repository Network subcommittee under the RDUF (https://japanlinkcenter.org/rduf/).

**Awareness of existing licensed tools**    In Step 5, the data provider selects the appropriate conditions of use. The guidelines prioritize practical use and present the preferable requirements consistent with existing licensed tools. Therefore, the additional categories of conditions of use revealed in the questionnaire survey analysis are not included. For example, we identify the "Preferable" conditions of use in Section 5.4. Still, we added "Impose the same conditions when publishing results," "Noncommercial," and "Noderivs" to the preferable requirements in the guidelines as a practical consideration. This decision corresponds to some contents that are difficult to distinguish from copyrighted works when giving conditions for data usage. A fundamental solution would be for research data to be managed under clear conditions of use from when it is created.

## 5.6    Results in chapter 5

In Chapter 5, we investigated the actual processing status related to setting conditions of use in different fields and clarified the correspondence of information tied to each formalized process. First, we conducted two surveys to investigate the actual processing status related to setting conditions of use in different fields. The conditions of use of research data are a complex combination of external constraints and the data provider's requests. We started clarifying the two boundaries based on the use-case of actual conditions of use in various fields. In the interview survey, we found the Risk Management process has been conducted on a field-by-field basis. Individual risk assessments have been conducted for each keyword, such as personal information and intellectual property rights. In the questionnaire survey, we observed differences among categories preferred and not preferred by the data providers quantitatively.

Through the analysis of the results of the two surveys, we discussed the correspondence between formalized processes and conditions of use. We clearly separated the external constraints and the data provider's request, and we discussed the Risk Management process and the Rights Management process in relation to each. Regarding the Risk Management process, we re-summarized external constraints into five categories from an interdisciplinary perspective. It highlighted constraints that should be judged locally and constraints that can be judged globally. Regarding the Rights Management process, we qualitatively analyzed the ten categories used in the questionnaire survey

to understand the data provider's request diversity. This categorization will allow us to better understand the data provider's requests and give some direction to the Rights Management Process. We also discussed the feasibility of these interpretations by the data re-user. Data re-users need to interpret interdisciplinary and specialized knowledge that differs from each field's practice.

Finally, we conducted a technical investigation for stepwise interpretation of these processing results. There is a knowledge gap between data providers and data re-users in understanding the conditions of use, so we need to bridge this gap by implementing some mechanisms. As a practical implementation, we developed a data publishing flow with licensing scenarios and guidelines based on our understanding of the formalized process. During the flow development process, we clarified the relationship to the Quality Assurance process discussed in Chapter 4. This study reframes the processes in the Appraisal category and complements Chapter 4 practices each other.

# 6

# Conclusion

In this chapter, we discuss the results of this study and prospects.

## 6.1   Conclusion

This study aimed to interpret the data curation process in various fields from an interdisciplinary perspective. In chapter 1, we set the following two objectives:

***Objective 1: Analysis and formalization of knowledge representing interdisciplinary data curation process***
***Objective 2: Practical studies to interpret the interdisciplinary data curation process***

For Objective 1, we addressed this issue in chapter 3. As the first step to interpret the tasks and procedures performed in different fields at the same granularity, we investigated the practices of data curation conducted in each field. As a result, we found that about 87.2% of the processes are interpretable across multiple fields. Also,

we realized that there needs a suitable model to describe the structure such as the relationships among Input-Output objects, processes, and staffing to accurately represent the data curation activity's structure in different fields.

Based on the preliminary analysis and survey results, we developed the Data Curation Process Ontology to formalize the data curation activity's structure in different fields. To verify the usefulness and validity of this ontology, we described and compared the several actual data curation activity's structures. Users can visualize the data curation activity's structure and the provenance information of the research data by using the format presented in Section 3.6. It is the important contribution of this study to compare the activity's structure of eight diverse repositories in a single model. Also, we showed that the ontology can use the specification of data curation process support functions by systems. This ontology can also help to allocate data curation process among multiple systems and improve the data curation process during data integration in different fields. With the formalized method enabled by the ontology, we can expect to expand the data curation process in various fields from an interdisciplinary perspective. Thus this study contributed to building a knowledge framework for a common understanding of the data curation process in different fields.

For Objective 2, we addressed this issue in chapter 4 and chapter 5. In chapter 4, we focused on publishing mechanisms for data papers that divert the journal peer-review. First, we evaluated the reference model for data publishing from a process perspective, positioning the model as a combination of existing quality assurance process and peer-review process. Furthermore, we conduct a technical investigation by implementing RDA/WDS model as a case study. In applying the model, we clarified possible inappropriate activity between multiple systems environments in the data evaluation process. We developed two verification measures to deal with them. Through this study, we examined practically the relationship and boundaries between the existing quality assurance process and the peer-review process in the Data Evaluation category. This effort can be positioned as an application-level representation of the processes in Data Evaluation category formalized in Chapter 3.

We point out a limitation of this study regarding the processes in the Appraisal category formalized in Chapter 3. The tasks and procedures addressed in the data evaluation process are closely related to the risk management process and rights management process performed in the Appraisal category. Since the model does not

explicitly mention the risk management process and rights management process, it can be assumed that it is under the researcher's control. However, both processes are assumed to vary greatly depending on the nature of the handled research data. In the case of sensitive data, data curators may require involvement in the early stage. In order to apply this model to different fields with more diverse characteristics, it will be necessary to clarify the relationship between the processes in the Appraisal category and the Data Evaluation category.

In Chapter 5, we investigated the actual processing status related to setting conditions of use in different fields and clarified the correspondence of information tied to each formalized process. First, we conducted two surveys to investigate the actual processing status related to setting conditions of use in different fields. The conditions of use of research data are a complex combination of external constraints and the data provider's requests. We started clarifying the two boundaries based on the use-case of actual conditions of use in various fields. In the interview survey, we found the Risk Management process has been conducted on a field-by-field basis. Individual risk assessments have been conducted for each keyword, such as personal information and intellectual property rights. In the questionnaire survey, we observed differences among categories preferred and not preferred by the data providers quantitatively.

Through the analysis of the results of the two surveys, we discussed the correspondence between formalized processes and conditions of use. We clearly separated the external constraints and the data provider's request, and we discussed the Risk Management process and the Rights Management process in relation to each. Regarding the Risk Management process, we re-summarized external constraints into five categories from an interdisciplinary perspective. It highlighted constraints that should be judged locally and constraints that can be judged globally. Regarding the Rights Management process, we qualitatively analyzed the ten categories used in the questionnaire survey to understand the data provider's request diversity. This categorization will allow us to better understand the data provider's requests and give some direction to the Rights Management Process. We also discussed the feasibility of these interpretations by the data re-user. Data re-users need to interpret interdisciplinary and specialized knowledge that differs from each field's practice.

Finally, we conducted a technical investigation for stepwise interpretation of these processing results. There is a knowledge gap between data providers and data

re-users in understanding the conditions of use, so we need to bridge this gap by implementing some mechanisms. As a practical implementation, we developed a data publishing flow with licensing scenarios and guidelines based on our understanding of the formalized process. During the flow development process, we clarified the relationship to the Quality Assurance process discussed in Chapter 4. This study reframes the processes in the Appraisal category and complements Chapter 4 practices each other.

Through this study, we provided our framework for understanding data curation activities as processes from an interdisciplinary perspective. This framework allows the data curation process to be interpreted in a decoupled way of the original research context. Data re-users will be able to formally assess the increased interpretability by verifying that the processes included in the framework have been properly executed. Furthermore, we have demonstrated some practical implementations through our problem-solving approach as a stepwise formalization to the level of interpretation. We believe that as we move forward with these formalized efforts, we can improve the interpretability of research data and thereby contribute to the realization of open science.

## 6.2   Future works

One aspect of our approach is to reveal the bounded context of each field by intentionally providing a different perspective than the original contextual understanding. Even in the same data format, there are parts that can be understood in common and parts that cannot, and these unintelligible parts are the boundaries between fields. We hope to clarify the boundaries of expertise in the field by applying our frameworks. Also, in parallel with efforts in the leading fields, we would like to explore the possibility of using our frameworks outside the designated community, such as in different fields and industries. Through these efforts, we hope our framework will evolve beyond an interdisciplinary perspective to a transdisciplinary one in the future.

# Bibliography

ARDC (2021). Sensitive data. https://ardc.edu.au/resources/working-with-data/sensitive-data/, (accessed 2022-5-11).

Assante, M., L. Candela, D. Castelli, and A. Tani (2016). Are scientific data repositories coping with research data publishing? *Data Science Journal 15*.

Austin, C. C., T. Bloom, S. Dallmeier-Tiessen, V. K. Khodiyar, F. Murphy, A. Nurnberger, L. Raymond, M. Stockhause, J. Tedds, M. Vardigan, and A. Whyte (2016). Key components of data publishing: using current best practices to develop a reference model for data publishing. *International Journal on Digital Libraries 18*(2), 77–92.

Bahim, C., C. Casorrán-Amilburu, M. Dekkers, E. Herczog, N. Loozen, K. Repanas, K. Russell, and S. Stall (2020). The fair data maturity model: An approach to harmonise fair assessments. *Data Science Journal 19*.

Baker, K. S. and L. Yarmey (2009). Data stewardship: Environmental data curation and a web-of-repositories. *International Journal of Digital Curation 4*(2), 12–27.

Baldwin, J. and S. Pinfield (2018). The uk scholarly communication licence: Attempting to cut through the gordian knot of the complexities of funder mandates, publisher embargoes and researcher caution in achieving open access. *Publications 6*(3).

Ball, A. (2010). Review of the state of the art of the digital curation of research data (version 1.2). https://purehost.bath.ac.uk/ws/portalfiles/portal/293012/erim1rep091103ab12.pdf, (accessed 2022-5-11).

Ball, A. (2012). *Review of Data Management Lifecycle Models.* UK United Kingdom: University of Bath. https://researchportal.bath.ac.uk/en/publications/review-of-data-management-lifecycle-models, (accessed 2022-5-11).

Ball, A. (2014). How to license research data. http://www.dcc.ac.uk/resources/how-guides, (accessed 2022-5-11).

Beagrie, N. (2008). Digital curation for science, digital libraries, and individuals. *International Journal of Digital Curation 1*, 3–16.

Borgman, C. L. (2007). *Scholarship in the Digital Age: Information, Infrastructure, and the Internet.* MIT Press.

Borgman, C. L. (2015). *Big Data, Little Data, No Data: Scholarship in the Networked World.* MIT Press.

Buneman, P., A. Chapman, and J. Cheney (2006). *A Provenance Model for Manually Curated Data.* Springer Berlin.

Cabinet Office, G. o. J. (2016). 5th science and technology basic plan. https://www8.cao.go.jp/cstp/kihonkeikaku/index5.html, (accessed 2022-5-11).

Cabinet Office, G. o. J. (2021). 6th science, technology, and innovation basic plan. https://www8.cao.go.jp/cstp/kihonkeikaku/index6.html, (accessed 2022-5-11).

Cabinet Office, Expert Panel on Open Science, G. o. J. (2015). Promoting open science in japan opening up a new era for the advancement of science. https://www8.cao.go.jp/cstp/sonota/openscience/, (accessed 2022-5-11).

Candela, L., D. Castelli, P. Manghi, and A. Tani (2015). Data journals: A survey. *Journal of the Association for Information Science and Technology 66*(9), 1747–1762.

CASRAI (2019a). Data curation. https://codata.org/rdm-terminology/data-curation/, (accessed 2022-8-18).

CASRAI (2019b). Research data. https://codata.org/rdm-terminology/research-data/, (accessed 2022-8-18).

Chao, T. C., M. H. Cragin, and C. L. Palmer (2014). Data practices and curation vocabulary (dpcvocab): An empirically derived framework of scientific data practices and curatorial processes. *Journal of the Association for Information Science and Technology 66*(3), 616–633.

Chin, G. and C. S. Lansing (2004). Capturing and supporting contexts for scientific data sharing via the biological sciences collaboratory. In *Proceedings of the 2004 ACM Conference on Computer Supported Cooperative Work*, CSCW '04, New York, NY, USA, pp. 409–418. Association for Computing Machinery.

CHORUS (2021). Publisher data availability policies index. https://www.chorusaccess.org/resources/chorus-for-publishers/publisher-data-availability-policies-index/, (accessed 2022-5-11).

Choudhury, S., C. Huang, and C. L. Palmer (2020). Updating the dcc curation lifecycle model. *International Journal of Digital Curation 15*(1), 12.

Claerbout, J. (2010). Reproducible computational research: A history of hurdles, mostly overcome. http://sepwww.stanford.edu/sep/jon/reproducible.html, (accessed 2022-5-11).

Clarke, D. J., L. Wang, A. Jones, M. L. Wojciechowicz, D. Torre, K. M. Jagodnik, S. L. Jenkins, P. McQuilton, Z. Flamholz, M. C. Silverstein, B. M. Schilder, K. Robasky, C. Castillo, R. Idaszak, S. C. Ahalt, J. Williams, S. Schurer, D. J. Cooper, R. de Miranda Azevedo, J. A. Klenk, M. A. Haendel, J. Nedzel, P. Avillach, M. E. Shimoyama, R. M. Harris, M. Gamble, R. Poten, A. L. Charbonneau, J. Larkin, C. T. Brown, V. R. Bonazzi, M. J. Dumontier, S.-A. Sansone, and A. Ma'ayan (2019). Fairshake: Toolkit to evaluate the fairness of research digital resources. *Cell Systems 9*(5), 417–421.

CODATA (2019). The beijing declaration on research data. https://zenodo.org/record/3552330, (accessed 2022-5-11).

Commission, E. (2016). Guidelines on fair data management in horizon 2020 ver.3.0. https://ec.europa.eu/research/participants/data/ref/h2020/grants_manual/hi/oa_pilot/h2020-hi-oa-data-mgt_en.pdf, (accessed 2022-5-11).

Consultative Committee for Space Data Systems, C. (2012). Rerence model for an open archival information system (oais): Magenta book. https://public.ccsds.org/Pubs/650x0m2.pdf, (accessed 2022-5-11).

Consultative Committee for Space Data Systems, C. (n.d.). Oais reference model (iso 14721) (website). http://www.oais.info/, (accessed 2022-5-11).

Creative Commons, U. (2009). Defining "noncommercial" : A study of how the online population understands "noncommercial use" . https://mirrors.creativecommons.org/defining-noncommercial/Defining_Noncommercial_fullreport.pdf, (accessed 2022-5-11).

DataONE, P. P. i. S. R. W. G. (2013). Data management guide for public participation in scientific research. https://old.dataone.org/sites/all/documents/DataONE-PPSR-DataManagementGuide.pdf, (accessed 2022-5-11).

Edwards, P. N., M. S. Mayernik, A. L. Batcheller, G. C. Bowker, and C. L. Borgman (2011). Science friction: Data, metadata, and collaboration. *Social Studies of Science 41*(5), 667–690.

Faniel, I., E. Kansa, S. Whitcher Kansa, J. Barrera-Gomez, and E. Yakel (2013). The challenges of digging data: A study of context in archaeological data reuse. In *Proceedings of the 13th ACM/IEEE-CS Joint Conference on Digital Libraries*, JCDL '13, pp. 295–304. Association for Computing Machinery.

Faniel, I. M. and T. E. Jacobsen (2010). Reusing scientific data: How earthquake engineering researchers assess the reusability of colleagues' data. *Computer Supported Cooperative Work (CSCW) 19*(3–4), 355–375.

Faniel, I. M., A. Kriesberg, and E. Yakel (2016). Social scientists' satisfaction with data reuse. *Journal of the Association for Information Science and Technology 67*(6), 1404–1416.

Faundeen, J. L., T. E. Burley, J. A. Carlino, D. L. Govoni, H. S. Henkel, S. L. Holl, V. B. Hutchison, E. Martín, E. T. Montgomery, C. Ladino, S. Tessler, and L. S. Zolly (2014). The united states geological survey science data lifecycle model. https://doi.org/10.3133/ofr20131265, (accessed 2022-5-11).

Fecher, B., S. Friesike, and M. Hebing (2015). What drives academic data sharing? *PLOS ONE 10*(2), 1–25.

FORCE11 (2014). *Joint Declaration of Data Citation Principles.* https://www.force11.org/group/joint-declaration-data-citation-principles-final, (accessed 2022-5-11).

G8 (2013). G8 science ministers statement. https://www.gov.uk/government/news/g8-science-ministers-statement, (accessed 2022-5-11).

Grabus, S. and J. Greenberg (2019). The landscape of rights and licensing initiatives for data sharing. *Data Science Journal 18*(1), 29.

Gray, J., A. S. Szalay, A. R. Thakar, C. Stoughton, and J. vandenBerg (2002). Online scientific data curation, publication, and archiving. In A. S. Szalay (Ed.), *Virtual Observatories*, Volume 4846, pp. 103 – 107. International Society for Optics and Photonics: SPIE.

Griffin, P., J. Khadake, K. LeMay, S. Lewis, S. Orchard, A. Pask, B. Pope, U. Roessner, K. Russell, T. Seemann, A. Treloar, S. Tyagi, J. Christiansen, S. Dayalan, S. Gladman, S. Hangartner, H. Hayden, W. Ho, G. Keeble-Gagnere, P. Korhonen, P. Neish, P. Prestes, M. Richardson, N. Watson-Haigh, K. Wyres, N. Young, and M. Schneider (2018). Best practice data life cycle approaches for the life sciences [version 2; peer review: 2 approved]. *F1000Research 6*(1618).

Grootveld, M. and J. van Egmond (2012). Peer-reviewed open research data: Results of a pilot. *International Journal of Digital Curation 7*.

Hayashi, M., Y. Hayashi, M. Asaoka, M. Kawai, Y. Minamiyama, and K. Yamaji (2021). Development of national-level institutional repository cloud service for open science. https://doi.org/10.5281/zenodo.5442279, (accessed 2022-5-11).

Hemphill, L., A. Pienta, S. Lafia, D. Akmon, and D. A. Bleckley (2022). How do properties of data, their curation, and their funding relate to reuse? *Journal of the Association for Information Science and Technology*, 1–13.

Higgins, S. (2008). The dcc curation lifecycle model. *International Journal of Digital Curation 3*(1), 134–140.

Higgins, S. (2011). Digital curation: The emergence of a new discipline. *International Journal of Digital Curation 6*(2), 78–88.

Higgins, S. (2018). Digital curation: the development of a discipline within information science. *Journal of Documentation 74*(6), 1318–1338.

Howard, T., M. Darlington, A. Ball, S. Cully, and C. McMahon (2010). *Opportunities for and Barriers to Engineering Research Data Re-use.* Number ERIM Project Document erim3rep100805tjh10 in ERIM Project Document. UK United Kingdom: University of Bath.

Hrynaszkiewicz, I. and Y. Shintani (2014). Scientific data : An open access and open data publication to facilitate reproducible research. *Journal of Information Processing and Management 57*(9), 629–640.

Hudson-Vitale, C., I. Heidi, J. R. Lisa, C. Jake, K. Wendy, O. Robert, and S. Claire (2017). *Data Curation. SPEC Kit 354.* Association of Research Libraries.

Humphrey, C. (2006). *e-Science and the Life Cycle of Research.* https://era.library.ualberta.ca/items/3334684b-fa6a-4c9d-a74b-559fecd42f9f, (accessed 2022-5-11).

ICPSR (2021). Guide to social science data preparation and archiving : Best practice throughout the data life cycle 6th edition. https://www.icpsr.umich.edu/web/pages/deposit/guide/, (accessed 2022-5-11).

Ikeuchi, U. and Y. Minamiyama (2020). Data supporting "investigation and development of the workflow to clarify conditions of use for research data publishing in japan" version 1.0. https://doi.org/10.20736/00001468, (accessed 2022-5-11).

ISO (2012). Iso 14721:2012 space data and information transfer systems — open archival information system (oais) — reference model. https://www.iso.org/standard/57284.html, (accessed 2022-5-11).

Johnson, W. G. (2008). The icpsr and social science research. *Behavioral amp; Social Sciences Librarian 27*(3–4), 140–157.

Johnston, L. R. (Ed.) (2017). *Curating Research Data. Volume Two: A Handbook of Current Practice.* Association of College  Research Libraries.

Johnston, L. R., J. Carlson, C. Hudson-Vitale, H. Imker, W. Kozlowski, R. Olendorf, and C. Stewart (2016). Definitions of data curation activities used by the data curation network. https://hdl.handle.net/11299/188638, (accessed 2022-5-11).

Jones, M. (2009). The cedars project. *Library and Information Research 26*(84), 11–16.

Kindling, M., H. Pampel, S. van de Sandt, J. Rücknagel, P. Vierkant, G. Kloska, M. Witt, P. Schirmbacher, R. Bertelmann, and F. Scholze (2017). The landscape of research data repositories in 2015: A re3data analysis. *D-Lib Magazine 23*(3/4).

Klump, J., R. Bertelmann, J. Brase, M. Diepenbroek, H. Grobe, H. Höck, M. Lautenschlager, U. Schindler, I. Sens, and J. Wächter (2006). Data publication in the open access initiative. *Data Science Journal 5*, 79–83.

Koso, A. (2013). Data sharing in the life sciences: How far we have come and where we must go. *Journal of Information Processing and Management 56*(5), 294–301.

Kowalczyk, S. and K. Shankar (2011). Data sharing in the sciences. *Annual Review of Information Science and Technology 45*(1), 247–294.

Kowalczyk, S. T. (2018). Modelling the research data lifecycle. *International Journal of Digital Curation 12*(2), 331–361.

Kratz, J. and C. Strasser (2014). Data publication consensus and controversies [version 3; peer review: 3 approved]. *F1000Research 3*(94).

Lawrence, B., C. Jones, B. Matthews, S. Pepler, and S. Callaghan (2011). Citation and peer review of data: Moving towards formal data publication. *International Journal of Digital Curation 6*(2), 4–37.

Lebo, T., S. Sahoo, D. McGuinness, K. Belhajjame, J. Cheney, D. Corsar, D. Garijo, S. Soiland-Reyes, S. Zednik, and J. Zhao (2013). Prov-o: The prov ontology. https://www.w3.org/TR/prov-o/, (accessed 2022-5-11).

Lee, C. A. and H. R. Tibbo (2007). Digital curation and trusted repositories: Steps toward success. *Journal of Digital Information 8*(2).

Lord, P. and A. Macdonald (2003). e-science curation report: Data curation for e-science in the uk: an audit to establish requirements for future curation and provision. https://digitalpreservation.gov/news/2004/e-ScienceReportFinal.pdf, (accessed 2022-5-11).

Matsuzato, F. (1992). A brief history of the national institute of polar research library (in japanese). *2*(2).

Mayernik, M. S., K. Breseman, R. R. Downs, R. Duerr, A. Garretson, and C.-Y. S. Hou (2020). Risk assessment for scientific data. *Data Science Journal 19.*

Mayernik, M. S., S. Callaghan, R. Leigh, J. Tedds, and S. Worley (2015). Peer review of datasets: When, why, and how. *Bulletin of the American Meteorological Society 96*(2).

Mayernik, M. S., T. DiLauro, R. Duerr, E. Metsger, A. E. Thessen, and G. S. Choudhury (2013). Data conservancy provenance, context, and lineage services: Key components for data preservation and curation. *Data Science Journal 12.*

Minamiyama, Y., U. Ikeuchi, K. Ueshima, M. Suto, N. Okayama, I. Yamada, K. Ebisawa, H. Nakanishi, and Y. Kumazaki (2020). Guidelines for specifying conditions of use in research data publishing, ver. 1.0. https://doi.org/10.11502/rduf_license_guideline, (accessed 2022-5-11).

Murayama, Y. and K. Hayashi (2014). New trends in open science (part 1): International trends in science and technology and academic information sharing frameworks and open research data (in japanese). *Science  Technology Trends* (146).

NISTEP (2015). New trends in open science (part 3): Trends in research data publishing and the promotion for sharing basis data of research papers (in japanese). *Science Technology Trends* (148).

Noy, N. and D. McGuinness (2001). Ontology development 101: A guide to creating your first ontology. *Technical Report SMI-2001-0880.*

OECD (2015a). Enquiries into intellectual property's economic impact, chapter 7: Legal aspects of open access to publicly funded research in enquiries into intellectual property's economic impact. https://www.oecd.org/sti/ieconomy/KBC2-IP.Final.pdf, (accessed 2022-5-11).

OECD (2015b). Making open science a reality. *OECD Science, Technology and Industry Policy Papers*. http://dx.doi.org/10.1787/5jrs2f963zs1-en, (accessed 2022-5-11).

OECD (2017). Business models for sustainable research data repositories. http://dx.doi.org/10.1787/302b12bb-en, (accessed 2022-5-11).

of Illinois Urbana-Champaign, U. (2017). Foundations of data curation. https://ischool.illinois.edu/degrees-programs/courses/is547, (accessed 2022-6-30).

Oostdijk, N., H. Van den Heuvel, and M. Treurniet (2013). The clarin-nl data curation service: Bringing data to the foreground. *International Journal of Digital Curation 8*(2), 134–145.

Open Knowledge, F. (n.d.). Open definition. https://opendefinition.org/, (accessed 2022-7-8).

Pampel, H., H. Pfeiffenberger, A. Schäfer, E. Smit, S. Pröll, and C. Bruch (2012). Report on peer review of research data in scholarly communication (part a of d 33.1. *Trends in Biotechnology.*

Peer, L., A. Green, and E. Stephenson (2014). Committing to data quality review. *International Journal of Digital Curation 9*(1), 263–291.

Piwowar, H. A. (2011). Who shares? who doesn't? factors associated with openly archiving raw research data. *PLOS ONE 6*(7), 1–13.

RDA-CODATA, L. I. I. G. (2016). *Legal Interoperability Of Research Data: Principles And Implementation Guidelines.* https://zenodo.org/record/162241, (accessed 2022-5-11).

Robinson-García, N., E. Jiménez-Contreras, and D. Torres-Salinas (2015). Analyzing data citation practices using the data citation index. *Journal of the Association for Information Science and Technology 67*(12), 2964–2975.

Rolland, B. and C. P. Lee (2013). Beyond trust and reliability: Reusing data in collaborative cancer epidemiology research. In *Proceedings of the 2013 Conference on Computer Supported Cooperative Work*, CSCW '13, pp. 435–444. Association for Computing Machinery.

Ruggles, S. (2017). The importance of data curation. *The Palgrave Handbook of Survey Research*, 303–308.

Sands, A., C. L. Borgman, L. Wynholds, and S. Traweek (2012). Follow the data: How astronomers use and reuse data. *Proceedings of the American Society for Information Science and Technology 49*(1), 1–3.

Spier, R. (2002). The history of the peer-review process. *Trends in Biotechnology 20*.

Stuart, D., G. Baynes, I. Hrynaszkiewicz, K. Allin, D. Penny, M. Lucraft, and M. Astell (2018). Whitepaper: Practical challenges for researchers in data sharing. *figshare*. https://figshare.com/articles/Whitepaper_Practical_challenges_for_researchers_in _data_sharing/5975011/1, (accessed 2022-5-11).

Sun, G. and C. S. G. Khoo (2017). Social science research data curation: issues of reuse. *Libellarium: journal for the research of writing, books, and cultural heritage institutions 9*(2).

Tenopir, C., S. Allard, K. Douglass, A. U. Aydinoglu, L. Wu, E. Read, M. Manoff, and M. Frame (2011). Data sharing by scientists: Practices and perceptions. *PLOS ONE 6*(6), 1–21.

Tenopir, C., N. M. Rice, S. Allard, L. Baird, J. Borycz, L. Christian, B. Grant, R. Olendorf, and R. J. Sandusky (2020). Data sharing, management, use, and reuse: Practices and perceptions of scientists worldwide. *PLOS ONE 15*(3), 1–26.

Tibbo, H. R. (2012). Placing the horse before the cart: Conceptual and technical dimensions of digital curation. *Historical Social Research / Historische Sozialforschung 37*(3 (141)), 187–200.

Tsoukala, V., M. Angelaki, V. Kalaitzi, B. Wessels, L. Price, M. J. Taylor, R. Smallwood, P. Linde, J. Sondervan, S. Reilly, M. Noorman, S. Wyatt, L. Bigagli, R. Finn, T. Sveins-

dottir, and K. Wadhwa (2016). Recode: Policy recommendations for open access to research data. https://zenodo.org/record/50863, (accessed 2022-5-11).

UK Data Archive, U. (n.d.a). Curation process (webiste). https://www.data-archive.ac.uk/managing-data/digital-curation-and-data-publishing/curation-process/, (accessed 2022-5-11).

UK Data Archive, U. (n.d.b). Quality control (webiste). https://www.data-archive.ac.uk/managing-data/digital-curation-and-data-publishing/quality-control/, (accessed 2022-7-8).

UNESCO (2021). Unesco recommendation on open science. pp. 1–34. https://unesdoc.unesco.org/ark:/48223/pf0000379949.locale=en, (accessed 2022-5-11).

Uschold, M. and M. Gruninger (1996). Ontologies: principles, methods and applications. *The Knowledge Engineering Review 11*(2), 93–136.

Van de Sompel, H., S. Payette, J. Erickson, C. Lagoze, and S. Warner (2004). Rethinking scholarly communication. *D-Lib Magazine 10*(9).

Van Panhuis, W. G., P. Paul, C. Emerson, J. Grefenstette, R. Wilder, A. J. Herbst, D. Heymann, and D. S. Burke (2014). A systematic review of barriers to data sharing in public health. *BMC Public Health 14*(1).

Venkatesan, A., N. Karamanis, M. Ide-Smith, J. Hickford, and J. McEntyre (2019). Understanding life sciences data curation practices via user research [version 1; peer review: 1 approved, 1 approved with reservations]. *F1000Research 8*(1622).

Vines, T. H., A. Y. K. Albert, R. L. Andrew, F. Débarre, D. G. Bock, M. T. Franklin, K. J. Gilbert, J.-S. Moore, S. Renaut, and D. J. Rennison (2014). The availability of research data declines rapidly with article age. *Current Biology 24*(1), 94–97.

Wallis, J. C., C. L. Borgman, M. S. Mayernik, and A. Pepe (2008). Moving archival practices upstream: An exploration of the life cycle of ecological sensing data in collaborative field research. *International Journal of Digital Curation 3*(1), 114–126.

Wallis, J. C., E. Rolando, and C. L. Borgman (2013). If we share data, will anyone use them? data sharing and reuse in the long tail of science and technology. *PLOS ONE 8*(7), 1–17.

Wilkinson, M. D., M. Dumontier, I. J. Aalbersberg, G. Appleton, M. Axton, A. Baak, N. Blomberg, J.-W. Boiten, L. B. da Silva Santos, P. E. Bourne, J. Bouwman, A. J. Brookes, T. Clark, M. Crosas, I. Dillo, O. Dumon, S. Edmunds, C. T. Evelo, R. Finkers, A. Gonzalez-Beltran, A. J. Gray, P. Groth, C. Goble, J. S. Grethe, J. Heringa, P. A. 't Hoen, R. Hooft, T. Kuhn, R. Kok, J. Kok, S. J. Lusher, M. E. Martone, A. Mons, A. L. Packer, B. Persson, P. Rocca-Serra, M. Roos, R. van Schaik, S.-A. Sansone, E. Schultes, T. Sengstag, T. Slater, G. Strawn, M. A. Swertz, M. Thompson, J. van der Lei, E. van Mulligen, J. Velterop, A. Waagmeester, P. Wittenburg, K. Wolstencroft, J. Zhao, and B. Mons (2016). The fair guiding principles for scientific data management and stewardship. *Scientific Data 3*(1).

Wilkinson, M. D., M. Dumontier, S.-A. Sansone, L. O. Bonino da Silva Santos, M. Prieto, D. Batista, P. McQuilton, T. Kuhn, P. Rocca-Serra, M. Crosas, and E. Schultes (2019). Evaluating fair maturity through a scalable, automated, community-governed framework. *Scientific Data 6*(1).

Willis, C., J. Greenberg, and H. White (2012). Analysis and synthesis of metadata goals for scientific data. *Journal of the American Society for Information Science and Technology 63*(8), 1505–1520.

Yoon, A. (2014). "making a square fit into a circle" : Researchers' experiences reusing qualitative data. *Proceedings of the American Society for Information Science and Technology 51*(1), 1–4.

Yoon, A. (2016). Red flags in data: Learning from failed data reuse experiences. *Proceedings of the Association for Information Science and Technology 53*(1), 1–6.

Zimmerman, A. S. (2008). New knowledge from old data. *Science, Technology, amp; Human Values 33*(5), 631–652.

# Acknowledgments

I would like to express my sincere gratitude to my advisor Prof. Hideaki Takeda for his guidance. I got a lot of learning and awareness through daily discussions. I cannot count the times of progressing in my research with his rich knowledge and deep insights. This research would not have been possible without him.

I would like to show my deepest appreciation to Prof. Kazutsuna Yamaji for his dedicated support. He allowed me to balance research and work in my job, and I was able to lead a research life in a highly favorable environment. I also learned how to deal with my research as a researcher.

I would like to offer my special thanks to Prof. Ikki Ohmukai for his kindness. He paved the way for me to enter the SOKENDAI doctoral program, even though I have been conducting research in a different field from informatics. Furthermore, he provided valuable suggestions for defining the data curation process as the evaluation committee member.

I would like to express my gratitude to Prof. Asanobu Kitamoto, Prof. Hitoshi Okada, and Prof. Masao Takaku for serving on the evaluation committee for this thesis. Prof. Kitamoto especially commented on the structure and direction of this research. Prof. Okada especially commented on the consistency of this research with the multiple fields behind it. Prof. Takaku commented on the methodology of this research and the perspective of writing the thesis. All of their suggestions are insightful and enhance the completeness of this thesis.

I appreciate my two lab members, especially Prof. Yusuke Komiyama, Dr. Sungmin Joo, Dr. Phuc Nguyen, Dr. Masaharu Hayashi, Mr. Makoto Asaoka, and Mr. Koichi Ojiro, for their insightful discussions. Daily discussions with them helped me to deepen my research.

I appreciate Dr. Yukiko Sakai for her kindness and support. Not only did she give the opportunity to pursue my research in earnest, but she also encouraged me to enter the doctoral program.

I gratefully acknowledge my parents and family, who together decided to enter the doctoral program and have always been by my side to encourage me. Without their support, this challenge would not have been possible.

Last but not least, this research was made possible by collaborating with many people whose names could not be mentioned here. I would like to thank them all.

# List of Publications

## Journal Papers

(1) Yasuyuki Minamiyama, Hideaki Takeda, Masaharu Hayashi, Makoto Asaoka, and Kazutsuna Yamaji. A study on formalizing data curation activities across different fields. (submitted to Information Processing and Management).

(2) Yasuyuki Minamiyama, Ui Ikeuchi, Kunihiko Ueshima, Nobuya Okayama, and Hideaki Takeda (2020). Investigation and Development of the Workflow to Clarify Conditions of Use for Research Data Publishing in Japan. Data Science Journal, 19(1), p.53. http://doi.org/10.5334/dsj-2020-053

(3) Yasuyuki Minamiyama, Takeshi Terui, Yasuhiro Murayama, Hironori Yabuki, Kazutsuna Yamaji, and Masaki Kanao (2017). Launching a new data journal "Polar Data Journal": Toward a new data publishing framework for polar science. Journal of Information Processing and Management. 60(3) p.147-156. https://doi.org/10.1241/johokanri.60.147 (in Japanese)

## Technical and Poster Papers

(1) Yasuyuki Minamiyama, Hideaki Takeda, Masaharu Hayashi, Makoto Asaoka, Kazutsuna Yamaji (2022). Development and utilization of data curation process ontology in 17th International Digital Curation Conference, (Poster).

(2) Y. Minamiyama, H. Takeda, M. Hayashi, M. Asaoka, K. Ojiro and K. Yamaji (2021). Structuring data curation activities by using ontology-based data modeling in

16th International Digital Curation Conference, (Poster).

(3) Yasuyuki Minamiyama, Akira Kadokura, Masaki Kanao, Takeshi Terui, Hironori Yabuki, Kazutsuna Yamaji (2017). "Polar Data Journal" ; A new data publishing platform for polar science in International Workshop on Sharing, Citation and Publication of Scientific Data across Disciplines, Joint Support-Center for Data Science Research (DS), (Poster).

(4) Yasuyuki Minamiyama, Takeshi Terui, Akira Kadokura, Masaki Kanao, Hironori Yabuki, Kazutsuna Yamaji (2017). Polar Data Journal by National Institute of Polar Research in JpGU-AGU Joint Meeting 2017, (Poster).

(5) Yasuyuki Minamiyama (2015). Data journal : A new approach to research data management. Current-Awareness. 325, pp.19-22. https://doi.org/10.11501/9497651 (in Japanese), (technical report).

## Other works

(1) Yasuyuki Minamiyama (2022). Data Curation Process Ontology ver. 1.0. https://purl.archive.org/curation-ontology, (Dataset).

(2) Yasuyuki Minamiyama, Ui Ikeuchi, Kunihiko Ueshima, Misaki Suto, Nobuya Okayama, Issaku Yamada, Ken Ebisawa, Hodaka Nakanishi, and Yui Kumazaki (2020). Guideline for specifying conditions of use in research data publishing, ver. 1.0., https://doi.org/10.11502/rduf_license_guideline (Other).

# Appendix 1. Topic guide of "Questions related to data curation activities"

We used a topic guide in Section 3.4.2 to share the specific phase of data curation activities with the interviewee. In the topic guide, we set the following nine questions referred to the previous study categories (Johnston, 2017). Note that the topic guide sent to interviewees was prepared in Japanese.

**List of questions**

1) Ingest
In this section, we will ask you about your data ingestion process details. This includes the identification of datasets available for registration, the procedure for receiving a deposit agreement, available media, and procedures for obtaining metadata and documentation.

2) Appraisal and Selection
In this section, we will ask you about your data appraisal and selection process details. This includes the identification of legal risks, such as the handling of personal information, arising from the characteristics of the data to be accepted, and the operation of the collection policies established by each repository.

3) Data Processing
In this section, we will ask you about your data processing process details. This

includes the storage and work space environment for processing data files before publishing, the scope of work logging and who created it, matters related to the software used to handle the data, and other data processing policies.

4) Data Storing

In this section, we will ask you about the system for storing data that has been processed. This includes the workflow up to preservation and the history management in each phase.

5) Metadata generation

In this section, we will ask you about handling metadata details submitted by data providers. This includes the validation, modifications, and additions that the data curator will make to the submitted metadata, distribution of the metadata envisioned, and the handling of metadata tied to a specific study.

6) Access level

In this section, we will ask you about the status of your repository's data access levels. This includes the options of access restrictions (e.g., affiliation, IP range, specific circumstances) and the types of conditions of use that are granted with description examples.

7) Long-term preservation

In this section, we will ask you about efforts for long-term preservation. This includes statements on preferred file formats, migration, format standardization/restructuring, emulation support, etc.

8) Re-evaluation and disposal

In this section, we will ask you about the metrics used to evaluate published data and how they are operationalized. This includes metrics information used as evaluation criteria (e.g., number of accesses, downloads, citations, papers linked to published data, etc.) and how the collected information is used fall under this category.

9) Other

In this section, we will ask you about the operation for conducting data curation activities. This includes the actual number of people involved in data curation activities and the allocation of specialized personnel, in addition to the number of cases processed per year and the time each process takes.

# Appendix 2. List of data curation process description rationale

In Section 3.4.2, we read and referred each organization's data curation process manuals and related documents for the rationale for the activities to ensure consistency with the interview results. We mapped the specific description of the activities and the data curators' information onto the working framework for those activities for which we were able to identify a description of the rationale for the activities. The following tables show the description of the rationale for the activities. We classified each data curation process mapped to the working framework into three levels: 1. Implemented, 2. Partially implemented, and 3. Not implemented. The statements on which the interviewee relies are masked: recorded as "Survey participant". Note that two of the MDR's processes have undergone changes involving category revisions since the time of the interview. We changed the mappings to match the currently available sources of information.

Table 1: List of data curation process description rationale
(khirin).

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 1 | Authentication | The process of confirming the identity of a person, generally the depositor, who is contributing data to the data repository. (e.g., password authentication or authorization via digital signature). Used for tracking provenance of the data files. | 3 | Survey participant | – |
| 2 | Chain of custody | Intentional recording of provenance metadata of the files (e.g., metadata about who created the file, when it was last edited, etc.) in order to preserve file authenticity when data are transferred to third-parties. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 3 | Deposit agreement | The certification by the data author (or depositor) that the data conform to all policies and conditions (e.g., do not violate any legal restrictions placed on the data) and are fit for deposit into the repository. A deposit agreement may also include rights transfer to the repository for ongoing stewardship. | 3 | Survey participant | – |
| 4 | Document-ation | Information describing any necessary information to use and understand the data. Documentation may be structured (e.g., a code book) or unstructured (e.g., a plain text "Readme" file). | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 5 | File Validation | A computational process to ensure that the intended data transfer to a repository was perfect and complete using means such as generating and validating file checksums (e.g., test if a digital file has changed at the bit level) and format validation to ensure that file types match their extensions. | 1 | Survey participant | – |
| 6 | Metadata | Information about a data set that is structured (often in machine-readable format) for purposes of search and retrieval. Metadata elements may include basic information (e.g. title, author, date created, etc.) and/or specific elements inherent to datasets (e.g., spatial coverage, time periods). | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 7 | Rights management | The process of tracking and managing ownership and copyright inherent to a data set as well as monitoring conditions and policies for access and reuse (e.g., licenses and data use agreements). | 3 | Survey participant | – |
| 8 | Risk management | The process of reviewing data for known risks such as confidentiality issues inherent to human subjects data, sensitive information (e.g., sexual histories, credit card information) or data regulated by law (e.g. HIPAA, FERPA) and taking actions to reject or facilitate remediation (e.g., de-identification services) when necessary. | 3 | Survey participant | – |
| 9 | Selection | The result of a successful appraisal. The data are determined appropriate for acceptance and ingest into the repository according to local collection policy and practice. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 10 | Arrangement and description | The re-organization of files (e.g., new folder directory structure) in a dataset that may also involve the creation of new file names, file descriptions, and the recording of technical metadata inherent to the files (e.g., date last modified). | 3 | Survey participant | – |
| 11 | Code review | Run and validate computer code (e.g., look for missing files and/or errors) in order to find mistakes overlooked in the initial development phase, improving the overall quality of software. | 1 | Survey participant | – |
| 12 | Contextualize | Use metadata to link the data set to related publications, dissertations, and/or projects that provide added context to how the data were generated and why. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 13 | Conversion (Analog) | In effort to increase the usability of a data set, the information is transferred into digital file formats (e.g., analog data keyed into a database). Note: digital conversion is also used to convert "fixed" data (e.g., PDF formats) into machine-readable formats. | 3 | Survey participant | – |
| 14 | Curation log | A written record of any changes made to the data during the curation process and by whom. File is often preserved as part of the overall record. | 3 | Survey participant | – |
| 15 | Data cleaning | A process used to improve data quality by detecting and correcting (or removing) defects & errors in data. | 1 | Survey participant | – |
| 16 | Deidentification | Redacting or removing personally identifiable or protected information (e.g., sensitive geographic locations) from a dataset prior to sharing with third-parties. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 17 | File format transformations | Transform files into open, non-proprietary file formats that broaden the potential for long-term reuse and ensure that additional preservation actions might be taken in the future. Note: Retention of the original file formats may be necessary if data transfer is not perfect. | 1 | Survey participant | – |
| 18 | Transcoding | With audio and video files, detect technical metadata (min resolution, audio/video codec) and encode files in ways that optimize reuse and long-term preservation actions. (E.g, Convert QuickTime files to MPEG4). | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 19 | File inventory or manifest | The data files are inspected periodically and the number, file types (extensions), and file sizes of the data are understood and documented. Any missing, duplicate, or corrupt (e.g., unable to open) files are discovered. | 1 | Survey participant | – |
| 20 | File renaming | To rename files in a dataset, often to standardize and/or reflect important metadata. | 1 | Survey participant | – |
| 21 | Indexing | Verify all metadata provided by the author and crosswalk to descriptive and administrative metadata compliant with a standard format for repository interoperability. | 3 | Survey participant | – |
| 22 | Interoperability | Formatting the data using a disciplinary standard for better integration with other datasets and/or systems. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 23 | Peer-review | The review of a data set by an expert with similar credentials and subject knowledge as the data creator for the purposes of validating the soundness and trustworthiness of the file contents. | 1 | Survey participant | – |
| 24 | Persistent Identifier | A URL (or Uniform Resource Locator) that is monitored by an authority to ensure a stable web location for consistent citation and long-term discoverability. Provides redirection when necessary. E.g., a Digital Object Identifier or DOI. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 25 | Quality assurance | Ensure that all documentation and metadata are comprehensive and complete. Example actions might include: open and run the data files; inspect the contents in order to validate, clean, and/or enhance data for future use; look for missing documentation about codes used, the significance of "null" and "blank" values, or unclear acronyms. | 1 | Survey participant | – |
| 26 | Restructure | Organize and/or reformate poorly structured data files to clarify their meaning and importance. | 3 | Survey participant | – |
| 27 | Software registry | Maintain copies of modern and obsolete versions of software (and any relevant code libraries) so that data may be opened/used overtime. | 1 | Survey participant | – |

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 28 | Contact information | Keep up-to-date contact information for the data authors and/or the contact persons in order to facilitate connection with third-party users. Often involves managing ephemeral information that will change over time. | 3 | Survey participant | – |
| 29 | Data citation | Display of a recommended bibliographic citation for a dataset to enable appropriate attribution by third-party users in order to formally incorporate data reuse as part of the scholarly ecosystem. | 1 | Survey participant | – |
| 30 | Data visualization | The presentation of pictorial and/or graphical representations of a data set used to identify patterns, detect errors, and/or demonstrate the extent of a data set to third party users. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 31 | Discovery Services | Services that incorporate machine-based search and retrieval functionality that help users identify what data exist, where the data are located, and how can they be accessed (e.g., full-text indexing or web optimization). | 3 | Survey participant | – |
| 32 | File download | Allow access to the data materials by authorized third parties. | 3 | Survey participant | – |
| 33 | Full-text indexing | Enhance the data for discovery purposes by generating search-engine-optimized formats of the text inherent to the data. | 3 | Survey participant | – |
| 34 | Metadata brokerage | Active dissemination of a data set's metadata to search and discovery services (e.g., article databases, catalogs, web-based indexes) for federated search and discovery. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 35 | Restricted access | In order to maintain the privacy of research subjects without losing integral components of the data, some data access will be protected and/or mediated to individuals that meet predefined criteria. | 3 | Survey participant | – |
| 36 | Embargo | To restrict or mediate access to a data set, usually for a set period of time. In some cases an embargo may be used to protect not only access, but any knowledge that the data exist. | 1 | Survey participant | – |
| 37 | Terms of use | Information provided to end users of a data set that outline the requirements or conditions for use (e.g., a Creative Commons License). | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 38 | Use analytics | Monitor and record how often data are viewed, requested, and/or downloaded. Track and report reuse metrics, such as data citations and impact measures for the data over time. | 2 | Survey participant | – |
| 39 | Cease data curation | Plan for any contingencies that will ultimately terminate access to the data. For example, providing tombstones or metadata records for data that have been deselected and removed from stewardship. | 1 | Survey participant | – |
| 40 | Migration | Monitor and anticipate file format obsolescence and, as needed, transform obsolete file formats to new formats as standards and use dictate. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 41 | Emulation | Provide legacy system configurations in modern equipment in order to ensure long-term usability of data. (E.g., arcade games emulated on modern web-browsers) | 1 | Survey participant | – |
| 42 | Secure storage | Data files are properly stored in a well-configured (in terms of hardware and software) storage environment that is routinely backed-up and physically protected. Perform routine fixity checks (to detect degradation or loss) and provide recovery services as needed. | 3 | Survey participant | – |
| 43 | File audit | Periodic review of the digital integrity of the data files and taking action when needed to protect data from digital erosion (e.g., bitrot) and/or hardware failure. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 44 | Repository certification | The technical and administrative capacities of the repository undergo review through a transparent and well-documented process by a trusted third-party accreditation body (e.g., TRAC, or Data Seal of Approval). | 1 | Survey participant | – |
| 45 | Succession planning | Planning for contingency, and/or escrow arrangements, in the case that the repository (or other entity responsible) ceases to operate or the institution substantially changes its scope. | 3 | Survey participant | – |
| 46 | Technology monitoring and Refresh | Formal, periodic review and assessment to ensure responsiveness to technological developments and evolving requirements of the digital infrastructure and hardware storing the data. | 2 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 47 | Versioning | Provide mechanisms to ingest new versions of the data overtime that includes metadata describing the version history and any changes made for each version. | 1 | Survey participant | – |

Table 2: List of data curation process description rationale (DiPLAS).

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 1 | Authentication | The process of confirming the identity of a person, generally the depositor, who is contributing data to the data repository. (e.g., password authentication or authorization via digital signature). Used for tracking provenance of the data files. | 3 | Survey participant | Internal manual |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 2 | Chain of custody | Intentional recording of provenance metadata of the files (e.g., metadata about who created the file, when it was last edited, etc.) in order to preserve file authenticity when data are transferred to third-parties. | 3 | Survey participant | Internal manual |
| 3 | Deposit agreement | The certification by the data author (or depositor) that the data conform to all policies and conditions (e.g., do not violate any legal restrictions placed on the data) and are fit for deposit into the repository. A deposit agreement may also include rights transfer to the repository for ongoing stewardship. | 3 | The National Museum of Ethnology will sign a memorandum of license agreement for each photographer, copyright holder, and owner of photographic materials. | http://diplas. jp/pdf/ requirements. html |

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 4 | Document-ation | Information describing any necessary information to use and understand the data. Documentation may be structured (e.g., a code book) or unstructured (e.g., a plain text "Readme" file). | 3 | The members of the publicly solicited project will use this database to successively add to the available information and carry out the research plan for the Grant-in-Aid for Scientific Research. The database will contain the following items: (1) ID, (2) photographic image, (3) photographer, copyright holder, (4) time of photographing, (5) region of photographing (country and local name at the time of photographing), (6) ethnic name (if identifiable), (7) date of photographing, (8) image content tag, (9) related information (references, etc.), and (10) free text (Japanese and English in principle). | http://diplas. jp/pdf/ requirements. html |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 5 | File Validation | A computational process to ensure that the intended data transfer to a repository was perfect and complete using means such as generating and validating file checksums (e.g., test if a digital file has changed at the bit level) and format validation to ensure that file types match their extensions. | 2 | Survey participant | – |
| 6 | Metadata | Information about a data set that is structured (often in machine-readable format) for purposes of search and retrieval. Metadata elements may include basic information (e.g. title, author, date created, etc.) and/or specific elements inherent to datasets (e.g., spatial coverage, time periods). | 3 | We will also convert the basic information attached to the photographic materials into data. For data that has already been digitized, such as photographs taken with a digital camera, we will assign an ID (if necessary, we will rename the file and retrieve the basic information). | http://diplas. jp/pdf/ requirements. html |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 7 | Rights management | The process of tracking and managing ownership and copyright inherent to a data set as well as monitoring conditions and policies for access and reuse (e.g., licenses and data use agreements). | 3 | 2. Copyright of the academic materials themselves. ; 3. copyright of the image data (still and moving images) taken of the academic materials. | http://diplas.jp/pdf/outline_Guideline.pdf |
| 8 | Risk management | The process of reviewing data for known risks such as confidentiality issues inherent to human subjects data, sensitive information (e.g., sexual histories, credit card information) or data regulated by law (e.g. HIPAA, FERPA) and taking actions to reject or facilitate remediation (e.g., de-identification services) when necessary. | 3 | 4. Rights relating to persons and acts recorded in academic information. | http://diplas.jp/pdf/outline_Guideline.pdf |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 9 | Selection | The result of a successful appraisal. The data are determined appropriate for acceptance and ingest into the repository according to local collection policy and practice. | 3 | "8. Acceptance or rejection: (1) The review will be conducted by the Open Call Project Review Committee, which will be placed under the Platform Committee that oversees the entire ""Digital Picture Library for Area Studies"", and will decide which proposals to adopt." | http://diplas. jp/pdf/ requirements. html |
| 10 | Arrangement and description | The re-organization of files (e.g., new folder directory structure) in a dataset that may also involve the creation of new file names, file descriptions, and the recording of technical metadata inherent to the files (e.g., date last modified). | 3 | Survey participant | Internal manual |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 11 | Code review | Run and validate computer code (e.g., look for missing files and/or errors) in order to find mistakes overlooked in the initial development phase, improving the overall quality of software. | 1 | Survey participant | – |
| 12 | Contextual-ize | Use metadata to link the data set to related publications, dissertations, and/or projects that provide added context to how the data were generated and why. | 1 | Survey participant | – |
| 13 | Conversion (Analog) | In effort to increase the usability of a data set, the information is transferred into digital file formats (e.g., analog data keyed into a database). Note: digital conversion is also used to convert "fixed" data (e.g., PDF formats) into machine-readable formats. | 3 | 3. Digitization of images and registration of basic information: Serial numbers (IDs) are assigned to materials, and images are digitized and registered in the database. | http://diplas.jp/outline.html |

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 14 | Curation log | A written record of any changes made to the data during the curation process and by whom. File is often preserved as part of the overall record. | 3 | Survey participant | Internal manual |
| 15 | Data cleaning | A process used to improve data quality by detecting and correcting (or removing) defects & errors in data. | 1 | Survey participant | – |
| 16 | Deidentification | Redacting or removing personally identifiable or protected information (e.g., sensitive geographic locations) from a dataset prior to sharing with third-parties. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 17 | File format transformations | Transform files into open, non-proprietary file formats that broaden the potential for long-term reuse and ensure that additional preservation actions might be taken in the future. Note: Retention of the original file formats may be necessary if data transfer is not perfect. | 3 | Survey participant | Internal manual |
| 18 | Transcoding | With audio and video files, detect technical metadata (min resolution, audio/video codec) and encode files in ways that optimize reuse and long-term preservation actions. (E.g, Convert QuickTime files to MPEG4). | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 19 | File inventory or manifest | The data files are inspected periodically and the number, file types (extensions), and file sizes of the data are understood and documented. Any missing, duplicate, or corrupt (e.g., unable to open) files are discovered. | 3 | Survey participant | Internal manual |
| 20 | File renaming | To rename files in a dataset, often to standardize and/or reflect important metadata. | 3 | Survey participant | Internal manual |
| 21 | Indexing | Verify all metadata provided by the author and crosswalk to descriptive and administrative metadata compliant with a standard format for repository interoperability. | 3 | Survey participant | Internal manual |
| 22 | Interoperability | Formatting the data using a disciplinary standard for better integration with other datasets and/or systems. | 3 | Survey participant | Internal manual |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 23 | Peer-review | The review of a data set by an expert with similar credentials and subject knowledge as the data creator for the purposes of validating the soundness and trustworthiness of the file contents. | 1 | Survey participant | – |
| 24 | Persistent Identifier | A URL (or Uniform Resource Locator) that is monitored by an authority to ensure a stable web location for consistent citation and long-term discoverability. Provides redirection when necessary. E.g., a Digital Object Identifier or DOI. | 1 | Survey participant | – |

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 25 | Quality assurance | Ensure that all documentation and metadata are comprehensive and complete. Example actions might include: open and run the data files; inspect the contents in order to validate, clean, and/or enhance data for future use; look for missing documentation about codes used, the significance of "null" and "blank" values, or unclear acronyms. | 3 | Survey participant | Internal manual |
| 26 | Restructure | Organize and/or reformate poorly structured data files to clarify their meaning and importance. | 1 | Survey participant | – |
| 27 | Software registry | Maintain copies of modern and obsolete versions of software (and any relevant code libraries) so that data may be opened/used overtime. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 28 | Contact information | Keep up-to-date contact information for the data authors and/or the contact persons in order to facilitate connection with third-party users. Often involves managing ephemeral information that will change over time. | 3 | Survey participant | – |
| 29 | Data citation | Display of a recommended bibliographic citation for a dataset to enable appropriate attribution by third-party users in order to formally incorporate data reuse as part of the scholarly ecosystem. | 1 | Survey participant | – |
| 30 | Data visualization | The presentation of pictorial and/or graphical representations of a data set used to identify patterns, detect errors, and/or demonstrate the extent of a data set to third party users. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 31 | Discovery Services | Services that incorporate machine-based search and retrieval functionality that help users identify what data exist, where the data are located, and how can they be accessed (e.g., full-text indexing or web optimization). | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 32 | File download | Allow access to the data materials by authorized third parties. | 3 | "The content entered into the database will be made available to the public if possible, after confirming the scope of use and rights in accordance with relevant laws and regulations, the National Museum of Ethnology's ""Guidelines for the Disclosure of Academic Information on the Internet"" and its ""Guidelines,"" based on the copyright holder's decision to make each image public or private, and the possibility that a photograph of a person may lead to a violation of rights or damage to dignity." | http://diplas.jp/outline.html |
| 33 | Full-text indexing | Enhance the data for discovery purposes by generating search-engine-optimized formats of the text inherent to the data. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 34 | Metadata brokerage | Active dissemination of a data set's metadata to search and discovery services (e.g., article databases, catalogs, web-based indexes) for federated search and discovery. | 3 | Survey participant | Internal manual |
| 35 | Restricted access | In order to maintain the privacy of research subjects without losing integral components of the data, some data access will be protected and/or mediated to individuals that meet predefined criteria. | 3 | "The content entered into the database will be made available to the public if possible, after confirming the scope of use and rights in accordance with relevant laws and regulations, the National Museum of Ethnology's ""Guidelines for the Disclosure of Academic Information on the Internet"" and its ""Guidelines,"" based on the copyright holder's decision to make each image public or private, and the possibility that a photograph of a person may lead to a violation of rights or damage to dignity." | http: //diplas.jp/ outline.html |

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 36 | Embargo | To restrict or mediate access to a data set, usually for a set period of time. In some cases an embargo may be used to protect not only access, but any knowledge that the data exist. | 1 | Survey participant | – |
| 37 | Terms of use | Information provided to end users of a data set that outline the requirements or conditions for use (e.g., a Creative Commons License). | 3 | Survey participant | – |
| 38 | Use analytics | Monitor and record how often data are viewed, requested, and/or downloaded. Track and report reuse metrics, such as data citations and impact measures for the data over time. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 39 | Cease data curation | Plan for any contingencies that will ultimately terminate access to the data. For example, providing tombstones or metadata records for data that have been de-selected and removed from stewardship. | 1 | Survey participant | – |
| 40 | Migration | Monitor and anticipate file format obsolescence and, as needed, transform obsolete file formats to new formats as standards and use dictate. | 3 | Survey participant | – |
| 41 | Emulation | Provide legacy system configurations in modern equipment in order to ensure long-term usability of data. (E.g., arcade games emulated on modern web-browsers) | 1 | Survey participant | – |

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 42 | Secure storage | Data files are properly stored in a well-configured (in terms of hardware and software) storage environment that is routinely backed-up and physically protected. Perform routine fixity checks (to detect degradation or loss) and provide recovery services as needed. | 1 | Survey participant | – |
| 43 | File audit | Periodic review of the digital integrity of the data files and taking action when needed to protect data from digital erosion (e.g., bitrot) and/or hardware failure. | 2 | Survey participant | – |
| 44 | Repository certification | The technical and administrative capacities of the repository undergo review through a transparent and well-documented process by a trusted third-party accreditation body (e.g., TRAC, or Data Seal of Approval). | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 45 | Succession planning | Planning for contingency, and/or escrow arrangements, in the case that the repository (or other entity responsible) ceases to operate or the institution substantially changes its scope. | 1 | Survey participant | – |
| 46 | Technology monitoring and Refresh | Formal, periodic review and assessment to ensure responsiveness to technological developments and evolving requirements of the digital infrastructure and hardware storing the data. | 2 | Survey participant | – |
| 47 | Versioning | Provide mechanisms to ingest new versions of the data overtime that includes metadata describing the version history and any changes made for each version. | 1 | Survey participant | – |

*End of table*

Table 3: List of data curation process description rationale (Materials Data Repository).

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 1 | Authentication | The process of confirming the identity of a person, generally the depositor, who is contributing data to the data repository. (e.g., password authentication or authorization via digital signature). Used for tracking provenance of the data files. | 3 | Users who already have a NIMS account (the account that you use to log in to desknet's, Denbun, and other NIMS systems) can use MDR using the same account. | https://dice.nims.go.jp/services/MDR/manual/html/login.html |
| 2 | Chain of custody | Intentional recording of provenance metadata of the files (e.g., metadata about who created the file, when it was last edited, etc.) in order to preserve file authenticity when data are transferred to third-parties. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 3 | Deposit agreement | The certification by the data author (or depositor) that the data conform to all policies and conditions (e.g., do not violate any legal restrictions placed on the data) and are fit for deposit into the repository. A deposit agreement may also include rights transfer to the repository for ongoing stewardship. | 3 | Click the checkbox to indicate that you have understood and agree to the Deposit Agreement. If the box is checked, one of the items in the "Requirements" at the top right turns from a red "!" to a green check, indicating that you satisfied one of the requirements(Fig. 13). | https://dice.nims.go.jp/services/MDR/manual/html/deposit.html |
| 4 | Document-ation | Information describing any necessary information to use and understand the data. Documentation may be structured (e.g., a code book) or unstructured (e.g., a plain text "Readme" file). | 3 | Enter your metadata (title, creators, etc.) in the Descriptions tab. The following fields must be filled out if they are applicable. | https://dice.nims.go.jp/services/MDR/manual/html/metadata-data.html |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 5 | File Validation | A computational process to ensure that the intended data transfer to a repository was perfect and complete using means such as generating and validating file checksums (e.g., test if a digital file has changed at the bit level) and format validation to ensure that file types match their extensions. | 2 | Survey participant | – |
| 6 | Metadata | Information about a data set that is structured (often in machine-readable format) for purposes of search and retrieval. Metadata elements may include basic information (e.g. title, author, date created, etc.) and/or specific elements inherent to datasets (e.g., spatial coverage, time periods). | 3 | Enter your metadata (title, creators, etc.) in the Descriptions tab. The following fields must be filled out if they are applicable. | https: //dice.nims.go. jp/services/ MDR/manual/ html/deposit. html |

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 7 | Rights management | The process of tracking and managing ownership and copyright inherent to a data set as well as monitoring conditions and policies for access and reuse (e.g., licenses and data use agreements). | 3 | Copyright license or a statement of the condition that re-users of this work must adhere to (how others may reuse your work). For further details see the entry for license. If the work has already been released under a license elsewhere (e.g. from the publisher), select the one that applies. | https: //dice.nims.go. jp/services/ MDR/manual/ html/deposit. html |
| 8 | Risk management | The process of reviewing data for known risks such as confidentiality issues inherent to human subjects data, sensitive information (e.g., sexual histories, credit card information) or data regulated by law (e.g. HIPAA, FERPA) and taking actions to reject or facilitate remediation (e.g., de-identification services) when necessary. | 3 | Enter the URL where we can confirm that your work is already public. (E.g. conference website, external database, etc.) | https: //dice.nims.go. jp/services/ MDR/manual/ html/deposit. html |

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 9 | Selection | The result of a successful appraisal. The data are determined appropriate for acceptance and ingest into the repository according to local collection policy and practice. | 3 | After you click on SUBMIT, MDR staff will review your deposit. If it has no issues regarding the terms of the deposit agreement and copyright policies, your work will be made public. | https://dice.nims.go.jp/services/MDR/manual/html/deposit.html |
| 10 | Arrangement and description | The re-organization of files (e.g., new folder directory structure) in a dataset that may also involve the creation of new file names, file descriptions, and the recording of technical metadata inherent to the files (e.g., date last modified). | 3 | Survey participant | – |
| 11 | Code review | Run and validate computer code (e.g., look for missing files and/or errors) in order to find mistakes overlooked in the initial development phase, improving the overall quality of software. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 12 | Contextualize | Use metadata to link the data set to related publications, dissertations, and/or projects that provide added context to how the data were generated and why. | 3 | Your work only contains the article: Publication ; Your primary intention is to make your paper available, with data accompanying the paper: Publication ; Your primary intention is to make your data/software available, with your paper documenting your data/software: Dataset ; You wish to describe your work using detailed metadata: Dataset | https://dice.nims.go.jp/services/MDR/manual/html/faq.html |
| 13 | Conversion (Analog) | In effort to increase the usability of a data set, the information is transferred into digital file formats (e.g., analog data keyed into a database). Note: digital conversion is also used to convert "fixed" data (e.g., PDF formats) into machine-readable formats. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 14 | Curation log | A written record of any changes made to the data during the curation process and by whom. File is often preserved as part of the overall record. | 3 | Survey participant | – |
| 15 | Data cleaning | A process used to improve data quality by detecting and correcting (or removing) defects & errors in data. | 1 | Survey participant | – |
| 16 | Deidentification | Redacting or removing personally identifiable or protected information (e.g., sensitive geographic locations) from a dataset prior to sharing with third-parties. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 17 | File format transformations | Transform files into open, non-proprietary file formats that broaden the potential for long-term reuse and ensure that additional preservation actions might be taken in the future. Note: Retention of the original file formats may be necessary if data transfer is not perfect. | 1 | Survey participant | – |
| 18 | Transcoding | With audio and video files, detect technical metadata (min resolution, audio/video codec) and encode files in ways that optimize reuse and long-term preservation actions. (E.g, Convert QuickTime files to MPEG4). | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 19 | File inventory or manifest | The data files are inspected periodically and the number, file types (extensions), and file sizes of the data are understood and documented. Any missing, duplicate, or corrupt (e.g., unable to open) files are discovered. | 1 | Survey participant | – |
| 20 | File renaming | To rename files in a dataset, often to standardize and/or reflect important metadata. | 1 | Survey participant | – |
| 21 | Indexing | Verify all metadata provided by the author and crosswalk to descriptive and administrative metadata compliant with a standard format for repository interoperability. | 3 | Survey participant | – |
| 22 | Interoperability | Formatting the data using a disciplinary standard for better integration with other datasets and/or systems. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 23 | Peer-review | The review of a data set by an expert with similar credentials and subject knowledge as the data creator for the purposes of validating the soundness and trustworthiness of the file contents. | 1 | Survey participant | – |
| 24 | Persistent Identifier | A URL (or Uniform Resource Locator) that is monitored by an authority to ensure a stable web location for consistent citation and long-term discoverability. Provides redirection when necessary. E.g., a Digital Object Identifier or DOI. | 3 | Your work will be assigned a new DOI ,and will be made public. | https://dice.nims.go.jp/services/MDR/manual/html/deposit.html |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 25 | Quality assurance | Ensure that all documentation and metadata are comprehensive and complete. Example actions might include: open and run the data files; inspect the contents in order to validate, clean, and/or enhance data for future use; look for missing documentation about codes used, the significance of "null" and "blank" values, or unclear acronyms. | 3 | After you click on SUBMIT, MDR staff will review your deposit. If it has no issues regarding the terms of the deposit agreement and copyright policies, your work will be made public. In some circumstances, an admin may request changes before accepting your work. | https://dice.nims.go.jp/services/MDR/manual/html/deposit.html |
| 26 | Restructure | Organize and/or reformate poorly structured data files to clarify their meaning and importance. | 1 | Survey participant | – |
| 27 | Software registry | Maintain copies of modern and obsolete versions of software (and any relevant code libraries) so that data may be opened/used overtime. | 1 | Instruments used for the work. | https://dice.nims.go.jp/services/MDR/manual/html/metadata-data.html#instruments |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|---------------------|------|
| 28 | Contact information | Keep up-to-date contact information for the data authors and/or the contact persons in order to facilitate connection with third-party users. Often involves managing ephemeral information that will change over time. | 3 | Survey participant | – |
| 29 | Data citation | Display of a recommended bibliographic citation for a dataset to enable appropriate attribution by third-party users in order to formally incorporate data reuse as part of the scholarly ecosystem. | 3 | Survey participant | – |
| 30 | Data visualization | The presentation of pictorial and/or graphical representations of a data set used to identify patterns, detect errors, and/or demonstrate the extent of a data set to third party users. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 31 | Discovery Services | Services that incorporate machine-based search and retrieval functionality that help users identify what data exist, where the data are located, and how can they be accessed (e.g., full-text indexing or web optimization). | 1 | Survey participant | – |
| 32 | File download | Allow access to the data materials by authorized third parties. | 3 | Open the Files tab from the top of the form and upload your files. You can use any of the following methods: | https://dice.nims.go.jp/services/MDR/manual/html/deposit.html |
| 33 | Full-text indexing | Enhance the data for discovery purposes by generating search-engine-optimized formats of the text inherent to the data. | 3 | Users can discover publications and datasets using metadata tailored for materials or by a full-text search, and can view and download them. | https://mdr.nims.go.jp/about?locale=en |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 34 | Metadata brokerage | Active dissemination of a data set's metadata to search and discovery services (e.g., article databases, catalogs, web-based indexes) for federated search and discovery. | 1→3 | Datasets and software resource types on MDR are indexed by Clarivate Data Citation Index. | https://dice.nims.go.jp/services/MDR/manual/html/functions.html |
| 35 | Restricted access | In order to maintain the privacy of research subjects without losing integral components of the data, some data access will be protected and/or mediated to individuals that meet predefined criteria. | 1 | You do not need to change this from MDR Open. | https://dice.nims.go.jp/services/MDR/manual/html/deposit.html |
| 36 | Embargo | To restrict or mediate access to a data set, usually for a set period of time. In some cases an embargo may be used to protect not only access, but any knowledge that the data exist. | 1 | Embargoes are not available in MDR. Please deposit your work only after it is ready to be viewed without restrictions. | https://dice.nims.go.jp/services/MDR/manual/html/faq.html |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 37 | Terms of use | Information provided to end users of a data set that outline the requirements or conditions for use (e.g., a Creative Commons License). | 3 | Copyright license or a statement of the condition that re-users of this work must adhere to (how others may reuse your work). For further details see the entry for license. If the work has already been released under a license elsewhere (e.g. from the publisher), select the one that applies. | https://dice.nims.go.jp/services/MDR/manual/html/deposit.html |
| 38 | Use analytics | Monitor and record how often data are viewed, requested, and/or downloaded. Track and report reuse metrics, such as data citations and impact measures for the data over time. | 2 | Survey participant | – |
| 39 | Cease data curation | Plan for any contingencies that will ultimately terminate access to the data. For example, providing tombstones or metadata records for data that have been deselected and removed from stewardship. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 40 | Migration | Monitor and anticipate file format obsolescence and, as needed, transform obsolete file formats to new formats as standards and use dictate. | 1 | Survey participant | – |
| 41 | Emulation | Provide legacy system configurations in modern equipment in order to ensure long-term usability of data. (E.g., arcade games emulated on modern web-browsers) | 1 | Survey participant | – |
| 42 | Secure storage | Data files are properly stored in a well-configured (in terms of hardware and software) storage environment that is routinely backed-up and physically protected. Perform routine fixity checks (to detect degradation or loss) and provide recovery services as needed. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 43 | File audit | Periodic review of the digital integrity of the data files and taking action when needed to protect data from digital erosion (e.g., bitrot) and/or hardware failure. | 1 | Survey participant | – |
| 44 | Repository certification | The technical and administrative capacities of the repository undergo review through a transparent and well-documented process by a trusted third-party accreditation body (e.g., TRAC, or Data Seal of Approval). | 1 | Survey participant | – |
| 45 | Succession planning | Planning for contingency, and/or escrow arrangements, in the case that the repository (or other entity responsible) ceases to operate or the institution substantially changes its scope. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 46 | Technology monitoring and Refresh | Formal, periodic review and assessment to ensure responsiveness to technological developments and evolving requirements of the digital infrastructure and hardware storing the data. | 3 | Survey participant | – |
| 47 | Versioning | Provide mechanisms to ingest new versions of the data overtime that includes metadata describing the version history and any changes made for each version. | 1→3 | "We recommend uploading your new version as a new work and refer to your older version in ""Related item"" metadata. Select ""is new version of"" for the relationship and enter the old version's URL and link title." | https://dice.nims.go.jp/services/MDR/manual/html/faq.html |

*End of table*

Table 4: List of data curation process description rationale (DARWIN).

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|----------------------|------|
| 1 | Authentication | The process of confirming the identity of a person, generally the depositor, who is contributing data to the data repository. (e.g., password authentication or authorization via digital signature). Used for tracking provenance of the data files. | 3 | Survey participant | – |
| 2 | Chain of custody | Intentional recording of provenance metadata of the files (e.g., metadata about who created the file, when it was last edited, etc.) in order to preserve file authenticity when data are transferred to third-parties. | 3 | Survey participant | – |

*Continued on next page*

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 3 | Deposit agreement | The certification by the data author (or depositor) that the data conform to all policies and conditions (e.g., do not violate any legal restrictions placed on the data) and are fit for deposit into the repository. A deposit agreement may also include rights transfer to the repository for ongoing stewardship. | 3 | 4.1. Before the voyage "Pre-voyage preparations for the handling data and samples are listed in Table 4.1. The forms for the pledge/agreement, metadata sheets and reporting documents will be sent by the management department prior to the voyage." | https://www.jamstec.go.jp/ceist/e/datasample/JAM_DS_Handbook.pdf |
| 4 | Document-ation | Information describing any necessary information to use and understand the data. Documentation may be structured (e.g., a code book) or unstructured (e.g., a plain text "Readme" file). | 3 | 3.2. Summary of Submissions "The types and summary of reporting documents, metadata, data, and samples to be submitted are listed in Table 3.2." | https://www.jamstec.go.jp/ceist/e/datasample/JAM_DS_Handbook.pdf |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 5 | File Validation | A computational process to ensure that the intended data transfer to a repository was perfect and complete using means such as generating and validating file checksums (e.g., test if a digital file has changed at the bit level) and format validation to ensure that file types match their extensions. | 2 | 3.3 Exceptions to the Data or Sample to be Submitted "In the event of any of the reasons listed in Table 3.3, the principal investigator may decide that submission is not required." | https://www.jamstec.go.jp/ceist/e/datasample/JAM_DS_Handbook.pdf |
| 6 | Metadata | Information about a data set that is structured (often in machine-readable format) for purposes of search and retrieval. Metadata elements may include basic information (e.g. title, author, date created, etc.) and/or specific elements inherent to datasets (e.g., spatial coverage, time periods). | 3 | 3.2. Summary of Submissions "The types and summary of reporting documents, metadata, data, and samples to be submitted are listed in Table 3.2." | https://www.jamstec.go.jp/ceist/e/datasample/JAM_DS_Handbook.pdf |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 7 | Rights management | The process of tracking and managing ownership and copyright inherent to a data set as well as monitoring conditions and policies for access and reuse (e.g., licenses and data use agreements). | 3 | Survey participant | – |
| 8 | Risk management | The process of reviewing data for known risks such as confidentiality issues inherent to human subjects data, sensitive information (e.g., sexual histories, credit card information) or data regulated by law (e.g. HIPAA, FERPA) and taking actions to reject or facilitate remediation (e.g., de-identification services) when necessary. | 3 | Survey participant | – |
| 9 | Selection | The result of a successful appraisal. The data are determined appropriate for acceptance and ingest into the repository according to local collection policy and practice. | 3 | 2.2 Confirmation of data and sample handling regulations | https://www.jamstec.go.jp/ceist/e/datasample/JAM_DS_Handbook.pdf |

*Continued on next page*

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|---------------------|------|
| 10 | Arrangement and description | The re-organization of files (e.g., new folder directory structure) in a dataset that may also involve the creation of new file names, file descriptions, and the recording of technical metadata inherent to the files (e.g., date last modified). | 3 | Survey participant | – |
| 11 | Code review | Run and validate computer code (e.g., look for missing files and/or errors) in order to find mistakes overlooked in the initial development phase, improving the overall quality of software. | 3 | Survey participant | – |
| 12 | Contextualize | Use metadata to link the data set to related publications, dissertations, and/or projects that provide added context to how the data were generated and why. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|----------------------|------|
| 13 | Conversion (Analog) | In effort to increase the usability of a data set, the information is transferred into digital file formats (e.g., analog data keyed into a database). Note: digital conversion is also used to convert "fixed" data (e.g., PDF formats) into machine-readable formats. | 1 | Survey participant | – |
| 14 | Curation log | A written record of any changes made to the data during the curation process and by whom. File is often preserved as part of the overall record. | 3 | 1. Introduction "This system has functions for inputting and outputting information via a Web browser, and enables management of metadata data, assistance with data publication tasks, and searching for past research voyages." | https://doi.org/10.5918/jamstecr.18.53 |
| 15 | Data cleaning | A process used to improve data quality by detecting and correcting (or removing) defects & errors in data. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 16 | Deidentific- ation | Redacting or removing per- sonally identifiable or pro- tected information (e.g., sensitive geographic lo- cations) from a dataset prior to sharing with third- parties. | 3 | Survey participant | – |
| 17 | File format transforma- tions | Transform files into open, non-proprietary file for- mats that broaden the po- tential for long-term reuse and ensure that additional preservation actions might be taken in the future. Note: Retention of the orig- inal file formats may be necessary if data transfer is not perfect. | 3 | Survey participant | – |
| 18 | Transcoding | With audio and video files, detect technical meta- data (min resolution, au- dio/video codec) and en- code files in ways that op- timize reuse and long-term preservation actions. (E.g, Convert QuickTime files to MPEG4). | 3 | Survey participant | – |

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 19 | File inventory or manifest | The data files are inspected periodically and the number, file types (extensions), and file sizes of the data are understood and documented. Any missing, duplicate, or corrupt (e.g., unable to open) files are discovered. | 3 | Survey participant | – |
| 20 | File renaming | To rename files in a dataset, often to standardize and/or reflect important metadata. | 3 | Survey participant | – |
| 21 | Indexing | Verify all metadata provided by the author and crosswalk to descriptive and administrative metadata compliant with a standard format for repository interoperability. | 3 | Survey participant | – |
| 22 | Interoperability | Formatting the data using a disciplinary standard for better integration with other datasets and/or systems. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 23 | Peer-review | The review of a data set by an expert with similar credentials and subject knowledge as the data creator for the purposes of validating the soundness and trustworthiness of the file contents. | 1 | Survey participant | – |
| 24 | Persistent Identifier | A URL (or Uniform Resource Locator) that is monitored by an authority to ensure a stable web location for consistent citation and long-term discoverability. Provides redirection when necessary. E.g., a Digital Object Identifier or DOI. | 3 | Minting DOIs for Information of Research Cruises Disseminated through DARWIN | http://www.jamstec.go.jp/e/database/darwin_doi.html |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|----------------------|------|
| 25 | Quality assurance | Ensure that all documentation and metadata are comprehensive and complete. Example actions might include: open and run the data files; inspect the contents in order to validate, clean, and/or enhance data for future use; look for missing documentation about codes used, the significance of "null" and "blank" values, or unclear acronyms. | 3 | Survey participant | – |
| 26 | Restructure | Organize and/or reformate poorly structured data files to clarify their meaning and importance. | 3 | Survey participant | – |
| 27 | Software registry | Maintain copies of modern and obsolete versions of software (and any relevant code libraries) so that data may be opened/used overtime. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 28 | Contact information | Keep up-to-date contact information for the data authors and/or the contact persons in order to facilitate connection with third-party users. Often involves managing ephemeral information that will change over time. | 3 | 6.1 Information for Proposal Applicants "The researcher who adopted the proposal will contact the department in charge whenever there are matters listed in Table 6.1.; (Metadata sheet content changes)" | https://www.jamstec.go.jp/ceist/e/datasample/JAM_DS_Handbook.pdf |
| 29 | Data citation | Display of a recommended bibliographic citation for a dataset to enable appropriate attribution by third-party users in order to formally incorporate data reuse as part of the scholarly ecosystem. | 3 | Survey participant | – |
| 30 | Data visualization | The presentation of pictorial and/or graphical representations of a data set used to identify patterns, detect errors, and/or demonstrate the extent of a data set to third party users. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 31 | Discovery Services | Services that incorporate machine-based search and retrieval functionality that help users identify what data exist, where the data are located, and how can they be accessed (e.g., full-text indexing or web optimization). | 3 | 5.3. Routine data reporting "We report the acquired data and observation information based on laws and regulations and agency collaboration, as listed in Table 5.3." | https://www.jamstec.go.jp/ceist/e/datasample/JAM_DS_Handbook.pdf |
| 32 | File download | Allow access to the data materials by authorized third parties. | 3 | 3.1 Submission period and embargo period "The Principal Investigator or Project Leader will compile the data acquired during the voyage and the samples for submission and submit them to the responsible department." | https://www.jamstec.go.jp/ceist/e/datasample/JAM_DS_Handbook.pdf |
| 33 | Full-text indexing | Enhance the data for discovery purposes by generating search-engine-optimized formats of the text inherent to the data. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 34 | Metadata brokerage | Active dissemination of a data set's metadata to search and discovery services (e.g., article databases, catalogs, web-based indexes) for federated search and discovery. | 3 | Survey participant | – |
| 35 | Restricted access | In order to maintain the privacy of research subjects without losing integral components of the data, some data access will be protected and/or mediated to individuals that meet predefined criteria. | 3 | Survey participant | – |
| 36 | Embargo | To restrict or mediate access to a data set, usually for a set period of time. In some cases an embargo may be used to protect not only access, but any knowledge that the data exist. | 3 | 3.1 Submission period and embargo period "The Principal Investigator or Project Leader will compile the data acquired during the voyage and the samples for submission and submit them to the responsible department." | https://www.jamstec.go.jp/ceist/e/datasample/JAM_DS_Handbook.pdf |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 37 | Terms of use | Information provided to end users of a data set that outline the requirements or conditions for use (e.g., a Creative Commons License). | 3 | Method of Provision of Data and Samples "JAMSTEC will provide the information on the Data and Samples in an easily accessible manner, and will try to satisfy the needs and requirements of users." ; Pricing Policies of Data and Samples "The Data and Samples belonging to JAMSTEC will be available free of charge for scientific and educational uses in principle except for the actual providing costs. Industrial uses of Data and Samples will be subject to be charged appropriately in principle and they will be depending on the nature of the use." | http://www.jamstec.go.jp/e/database/data_policy.html |

*Continued on next page*

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|---------------------|------|
| 38 | Use analytics | Monitor and record how often data are viewed, requested, and/or downloaded. Track and report reuse metrics, such as data citations and impact measures for the data over time. | 2 | Survey participant | – |
| 39 | Cease data curation | Plan for any contingencies that will ultimately terminate access to the data. For example, providing tombstones or metadata records for data that have been deselected and removed from stewardship. | 1 | Survey participant | – |
| 40 | Migration | Monitor and anticipate file format obsolescence and, as needed, transform obsolete file formats to new formats as standards and use dictate. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 41 | Emulation | Provide legacy system configurations in modern equipment in order to ensure long-term usability of data. (E.g., arcade games emulated on modern web-browsers) | 1 | Survey participant | – |
| 42 | Secure storage | Data files are properly stored in a well-configured (in terms of hardware and software) storage environment that is routinely backed-up and physically protected. Perform routine fixity checks (to detect degradation or loss) and provide recovery services as needed. | 3 | Survey participant | – |
| 43 | File audit | Periodic review of the digital integrity of the data files and taking action when needed to protect data from digital erosion (e.g., bitrot) and/or hardware failure. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|---------------------|------|
| 44 | Repository certification | The technical and administrative capacities of the repository undergo review through a transparent and well-documented process by a trusted third-party accreditation body (e.g., TRAC, or Data Seal of Approval). | 1 | Survey participant | – |
| 45 | Succession planning | Planning for contingency, and/or escrow arrangements, in the case that the repository (or other entity responsible) ceases to operate or the institution substantially changes its scope. | 3 | Survey participant | – |
| 46 | Technology monitoring and Refresh | Formal, periodic review and assessment to ensure responsiveness to technological developments and evolving requirements of the digital infrastructure and hardware storing the data. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 47 | Versioning | Provide mechanisms to ingest new versions of the data overtime that includes metadata describing the version history and any changes made for each version. | 3 | Survey participant | – |

*End of table*

Table 5: List of data curation process description rationale (Global Environmental Database).

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 1 | Authentication | The process of confirming the identity of a person, generally the depositor, who is contributing data to the data repository. (e.g., password authentication or authorization via digital signature). Used for tracking provenance of the data files. | 3 | Survey participant | – |

*Continued on next page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 2 | Chain of custody | Intentional recording of provenance metadata of the files (e.g., metadata about who created the file, when it was last edited, etc.) in order to preserve file authenticity when data are transferred to third-parties. | 3 | Survey participant | – |
| 3 | Deposit agreement | The certification by the data author (or depositor) that the data conform to all policies and conditions (e.g., do not violate any legal restrictions placed on the data) and are fit for deposit into the repository. A deposit agreement may also include rights transfer to the repository for ongoing stewardship. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 4 | Document-ation | Information describing any necessary information to use and understand the data. Documentation may be structured (e.g., a code book) or unstructured (e.g., a plain text "Readme" file). | 3 | Survey participant | – |
| 5 | File Valida-tion | A computational process to ensure that the intended data transfer to a repository was perfect and complete using means such as generating and validating file checksums (e.g., test if a digital file has changed at the bit level) and format validation to ensure that file types match their extensions. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 6 | Metadata | Information about a data set that is structured (often in machine-readable format) for purposes of search and retrieval. Metadata elements may include basic information (e.g. title, author, date created, etc.) and/or specific elements inherent to datasets (e.g., spatial coverage, time periods). | 3 | Survey participant | – |
| 7 | Rights management | The process of tracking and managing ownership and copyright inherent to a data set as well as monitoring conditions and policies for access and reuse (e.g., licenses and data use agreements). | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 8 | Risk management | The process of reviewing data for known risks such as confidentiality issues inherent to human subjects data, sensitive information (e.g., sexual histories, credit card information) or data regulated by law (e.g. HIPAA, FERPA) and taking actions to reject or facilitate remediation (e.g., de-identification services) when necessary. | 3 | Survey participant | – |
| 9 | Selection | The result of a successful appraisal. The data are determined appropriate for acceptance and ingest into the repository according to local collection policy and practice. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 10 | Arrangement and description | The re-organization of files (e.g., new folder directory structure) in a dataset that may also involve the creation of new file names, file descriptions, and the recording of technical metadata inherent to the files (e.g., date last modified). | 3 | Survey participant | – |
| 11 | Code review | Run and validate computer code (e.g., look for missing files and/or errors) in order to find mistakes overlooked in the initial development phase, improving the overall quality of software. | 1 | Survey participant | – |
| 12 | Contextualize | Use metadata to link the data set to related publications, dissertations, and/or projects that provide added context to how the data were generated and why. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 13 | Conversion (Analog) | In effort to increase the usability of a data set, the information is transferred into digital file formats (e.g., analog data keyed into a database). Note: digital conversion is also used to convert "fixed" data (e.g., PDF formats) into machine-readable formats. | 1 | Survey participant | – |
| 14 | Curation log | A written record of any changes made to the data during the curation process and by whom. File is often preserved as part of the overall record. | 3 | Survey participant | – |
| 15 | Data cleaning | A process used to improve data quality by detecting and correcting (or removing) defects & errors in data. | 3 | Survey participant | – |
| 16 | Deidentification | Redacting or removing personally identifiable or protected information (e.g., sensitive geographic locations) from a dataset prior to sharing with third-parties. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 17 | File format transformations | Transform files into open, non-proprietary file formats that broaden the potential for long-term reuse and ensure that additional preservation actions might be taken in the future. Note: Retention of the original file formats may be necessary if data transfer is not perfect. | 3 | Survey participant | – |
| 18 | Transcoding | With audio and video files, detect technical metadata (min resolution, audio/video codec) and encode files in ways that optimize reuse and long-term preservation actions. (E.g, Convert QuickTime files to MPEG4). | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|----------------------|------|
| 19 | File inventory or manifest | The data files are inspected periodically and the number, file types (extensions), and file sizes of the data are understood and documented. Any missing, duplicate, or corrupt (e.g., unable to open) files are discovered. | 3 | Survey participant | – |
| 20 | File renaming | To rename files in a dataset, often to standardize and/or reflect important metadata. | 3 | Survey participant | – |
| 21 | Indexing | Verify all metadata provided by the author and crosswalk to descriptive and administrative metadata compliant with a standard format for repository interoperability. | 3 | Survey participant | – |
| 22 | Interoperability | Formatting the data using a disciplinary standard for better integration with other datasets and/or systems. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 23 | Peer-review | The review of a data set by an expert with similar credentials and subject knowledge as the data creator for the purposes of validating the soundness and trustworthiness of the file contents. | 1 | Survey participant | – |
| 24 | Persistent Identifier | A URL (or Uniform Resource Locator) that is monitored by an authority to ensure a stable web location for consistent citation and long-term discoverability. Provides redirection when necessary. E.g., a Digital Object Identifier or DOI. | 3 | Survey participant | – |

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 25 | Quality assurance | Ensure that all documentation and metadata are comprehensive and complete. Example actions might include: open and run the data files; inspect the contents in order to validate, clean, and/or enhance data for future use; look for missing documentation about codes used, the significance of "null" and "blank" values, or unclear acronyms. | 3 | Survey participant | – |
| 26 | Restructure | Organize and/or reformate poorly structured data files to clarify their meaning and importance. | 3 | Survey participant | – |
| 27 | Software registry | Maintain copies of modern and obsolete versions of software (and any relevant code libraries) so that data may be opened/used overtime. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 28 | Contact information | Keep up-to-date contact information for the data authors and/or the contact persons in order to facilitate connection with third-party users. Often involves managing ephemeral information that will change over time. | 3 | Survey participant | – |
| 29 | Data citation | Display of a recommended bibliographic citation for a dataset to enable appropriate attribution by third-party users in order to formally incorporate data reuse as part of the scholarly ecosystem. | 3 | Survey participant | – |
| 30 | Data visualization | The presentation of pictorial and/or graphical representations of a data set used to identify patterns, detect errors, and/or demonstrate the extent of a data set to third party users. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 31 | Discovery Services | Services that incorporate machine-based search and retrieval functionality that help users identify what data exist, where the data are located, and how can they be accessed (e.g., full-text indexing or web optimization). | 3 | Survey participant | – |
| 32 | File download | Allow access to the data materials by authorized third parties. | 3 | Survey participant | – |
| 33 | Full-text indexing | Enhance the data for discovery purposes by generating search-engine-optimized formats of the text inherent to the data. | 1 | Survey participant | – |
| 34 | Metadata brokerage | Active dissemination of a data set's metadata to search and discovery services (e.g., article databases, catalogs, web-based indexes) for federated search and discovery. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 35 | Restricted access | In order to maintain the privacy of research subjects without losing integral components of the data, some data access will be protected and/or mediated to individuals that meet predefined criteria. | 3 | Survey participant | – |
| 36 | Embargo | To restrict or mediate access to a data set, usually for a set period of time. In some cases an embargo may be used to protect not only access, but any knowledge that the data exist. | 1 | Survey participant | – |
| 37 | Terms of use | Information provided to end users of a data set that outline the requirements or conditions for use (e.g., a Creative Commons License). | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 38 | Use analytics | Monitor and record how often data are viewed, requested, and/or downloaded. Track and report reuse metrics, such as data citations and impact measures for the data over time. | 2 | Survey participant | – |
| 39 | Cease data curation | Plan for any contingencies that will ultimately terminate access to the data. For example, providing tombstones or metadata records for data that have been deselected and removed from stewardship. | 1 | Survey participant | – |
| 40 | Migration | Monitor and anticipate file format obsolescence and, as needed, transform obsolete file formats to new formats as standards and use dictate. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 41 | Emulation | Provide legacy system configurations in modern equipment in order to ensure long-term usability of data. (E.g., arcade games emulated on modern web-browsers) | 1 | Survey participant | – |
| 42 | Secure storage | Data files are properly stored in a well-configured (in terms of hardware and software) storage environment that is routinely backed-up and physically protected. Perform routine fixity checks (to detect degradation or loss) and provide recovery services as needed. | 1 | Survey participant | – |
| 43 | File audit | Periodic review of the digital integrity of the data files and taking action when needed to protect data from digital erosion (e.g., bitrot) and/or hardware failure. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|------------|---------------------|------|
| 44 | Repository certification | The technical and administrative capacities of the repository undergo review through a transparent and well-documented process by a trusted third-party accreditation body (e.g., TRAC, or Data Seal of Approval). | 1 | Survey participant | – |
| 45 | Succession planning | Planning for contingency, and/or escrow arrangements, in the case that the repository (or other entity responsible) ceases to operate or the institution substantially changes its scope. | 1 | Survey participant | – |
| 46 | Technology monitoring and Refresh | Formal, periodic review and assessment to ensure responsiveness to technological developments and evolving requirements of the digital infrastructure and hardware storing the data. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 47 | Versioning | Provide mechanisms to ingest new versions of the data overtime that includes metadata describing the version history and any changes made for each version. | 3 | Survey participant | – |

*End of table*

Table 6: List of data curation process description rationale (Rikkyo University Data Archive).

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 1 | Authentication | The process of confirming the identity of a person, generally the depositor, who is contributing data to the data repository. (e.g., password authentication or authorization via digital signature). Used for tracking provenance of the data files. | 3 | 1. Contact about depositing the data "In the first instance, please contact RUDA if you are considering depositing data. Details of the deposition process will be explained by the relevant staff (In addition, the staff may ask you some questions regarding the data)." | https://spirit. rikkyo.ac.jp/ csi/RUDA/ SitePages/ index.aspx |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 2 | Chain of custody | Intentional recording of provenance metadata of the files (e.g., metadata about who created the file, when it was last edited, etc.) in order to preserve file authenticity when data are transferred to third-parties. | 3 | Record of the work | Internal manual |
| 3 | Deposit agreement | The certification by the data author (or depositor) that the data conform to all policies and conditions (e.g., do not violate any legal restrictions placed on the data) and are fit for deposit into the repository. A deposit agreement may also include rights transfer to the repository for ongoing stewardship. | 3 | 4. Concluding the Deposit Agreement "Once RUDA receives the materials to be deposited, it will forward to you two copies of the "Deposit Agreement." After you carefully read the form, please sign and stamp (if possible) both forms and return one of them to RUDA via post (Please keep the other as a receipt)." | https://spirit. rikkyo.ac.jp/ csi/RUDA/ SitePages/ index.aspx |

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 4 | Document-ation | Information describing any necessary information to use and understand the data. Documentation may be structured (e.g., a code book) or unstructured (e.g., a plain text "Readme" file). | 3 | 3. Transfer of deposition items "On the basis of the content written in the "Deposition Check-list," RUDA will determine whether or not to accept the data. If it is deemed acceptable, RUDA will ask the depositor for basic information about the survey by forwarding a "Metadata Sheet" The information provided in the "Metadata Sheet" will be publicized on our website after the data has been deposited. After filling the "Metadata Sheet," please send it to RUDA via either post or email. After receiving the "Metadata Sheet," our staff will provide the details for transferring the data and related material." | https://spirit. rikkyo.ac.jp/ csi/RUDA/ SitePages/ index.aspx |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 5 | File Validation | A computational process to ensure that the intended data transfer to a repository was perfect and complete using means such as generating and validating file checksums (e.g., test if a digital file has changed at the bit level) and format validation to ensure that file types match their extensions. | 3 | How to transfer with Pros-elf | Internal manual |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 6 | Metadata | Information about a data set that is structured (often in machine-readable format) for purposes of search and retrieval. Metadata elements may include basic information (e.g. title, author, date created, etc.) and/or specific elements inherent to datasets (e.g., spatial coverage, time periods). | 3 | 3. Transfer of deposition items "On the basis of the content written in the "Deposition Checklist," RUDA will determine whether or not to accept the data. If it is deemed acceptable, RUDA will ask the depositor for basic information about the survey by forwarding a "Metadata Sheet" The information provided in the "Metadata Sheet" will be publicized on our website after the data has been deposited. After filling the "Metadata Sheet," please send it to RUDA via either post or email. After receiving the "Metadata Sheet," our staff will provide the details for transferring the data and related material." | https://spirit. rikkyo.ac.jp/ csi/RUDA/ SitePages/ index.aspx |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 7 | Rights management | The process of tracking and managing ownership and copyright inherent to a data set as well as monitoring conditions and policies for access and reuse (e.g., licenses and data use agreements). | 3 | Prior to depositing your data. "Before proceeding to deposit the data, please confirm the following requirements: The depositor holds the copyright, ownership rights, and the license of the data. ; No other party's right are violated by depositing your data. ; Unless these requirements are satisfied, RUDA may refuse any offers to deposit the data." | https://spirit.rikkyo.ac.jp/csi/RUDA/SitePages/index.aspx |
| 8 | Risk management | The process of reviewing data for known risks such as confidentiality issues inherent to human subjects data, sensitive information (e.g., sexual histories, credit card information) or data regulated by law (e.g. HIPAA, FERPA) and taking actions to reject or facilitate remediation (e.g., de-identification services) when necessary. | 3 | Prior to depositing your data. "Before proceeding to deposit the data, please confirm the following requirements: The depositor holds the copyright, ownership rights, and the license of the data. ; No other party's right are violated by depositing your data. ; Unless these requirements are satisfied, RUDA may refuse any offers to deposit the data." | https://spirit.rikkyo.ac.jp/csi/RUDA/SitePages/index.aspx |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 9 | Selection | The result of a successful appraisal. The data are determined appropriate for acceptance and ingest into the repository according to local collection policy and practice. | 3 | 3. Transfer of deposition items "On the basis of the content written in the "Deposition Checklist," RUDA will determine whether or not to accept the data. If it is deemed acceptable, RUDA will ask the depositor for basic information about the survey by forwarding a "Metadata Sheet" The information provided in the "Metadata Sheet" will be publicized on our website after the data has been deposited. After filling the "Metadata Sheet," please send it to RUDA via either post or email. After receiving the "Metadata Sheet," our staff will provide the details for transferring the data and related material." | https://spirit. rikkyo.ac.jp/ csi/RUDA/ SitePages/ index.aspx |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 10 | Arrangement and description | The re-organization of files (e.g., new folder directory structure) in a dataset that may also involve the creation of new file names, file descriptions, and the recording of technical metadata inherent to the files (e.g., date last modified). | 3 | Survey participant | Internal manual |
| 11 | Code review | Run and validate computer code (e.g., look for missing files and/or errors) in order to find mistakes overlooked in the initial development phase, improving the overall quality of software. | 3 | 3 Check with published materials | Internal manual |
| 12 | Contextualize | Use metadata to link the data set to related publications, dissertations, and/or projects that provide added context to how the data were generated and why. | 3 | Survey participant | – |

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 13 | Conversion (Analog) | In effort to increase the usability of a data set, the information is transferred into digital file formats (e.g., analog data keyed into a database). Note: digital conversion is also used to convert "fixed" data (e.g., PDF formats) into machine-readable formats. | 3 | Cases of refusal (did) / consider (did) | Internal manual |
| 14 | Curation log | A written record of any changes made to the data during the curation process and by whom. File is often preserved as part of the overall record. | 3 | Record of the work | Internal manual |
| 15 | Data cleaning | A process used to improve data quality by detecting and correcting (or removing) defects & errors in data. | 3 | Cleaning manual | Internal manual |
| 16 | Deidentification | Redacting or removing personally identifiable or protected information (e.g., sensitive geographic locations) from a dataset prior to sharing with third-parties. | 3 | 0.8 Delete unnecessary variables | Internal manual |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 17 | File format transformations | Transform files into open, non-proprietary file formats that broaden the potential for long-term reuse and ensure that additional preservation actions might be taken in the future. Note: Retention of the original file formats may be necessary if data transfer is not perfect. | 3 | 2. After data cleaning | Internal manual |
| 18 | Transcoding | With audio and video files, detect technical metadata (min resolution, audio/video codec) and encode files in ways that optimize reuse and long-term preservation actions. (E.g, Convert QuickTime files to MPEG4). | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 19 | File inventory or manifest | The data files are inspected periodically and the number, file types (extensions), and file sizes of the data are understood and documented. Any missing, duplicate, or corrupt (e.g., unable to open) files are discovered. | 3 | Double check | Internal manual |
| 20 | File renaming | To rename files in a dataset, often to standardize and/or reflect important metadata. | 3 | Survey participant | – |
| 21 | Indexing | Verify all metadata provided by the author and crosswalk to descriptive and administrative metadata compliant with a standard format for repository interoperability. | 3 | Survey participant | – |
| 22 | Interoperability | Formatting the data using a disciplinary standard for better integration with other datasets and/or systems. | 3 | 6.3 Creating and Writing a Modified Syntax | Internal manual |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 23 | Peer-review | The review of a data set by an expert with similar credentials and subject knowledge as the data creator for the purposes of validating the soundness and trustworthiness of the file contents. | 3 | 7 Logical check | Internal manual |
| 24 | Persistent Identifier | A URL (or Uniform Resource Locator) that is monitored by an authority to ensure a stable web location for consistent citation and long-term discoverability. Provides redirection when necessary. E.g., a Digital Object Identifier or DOI. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|---------------------|------|
| 25 | Quality assurance | Ensure that all documentation and metadata are comprehensive and complete. Example actions might include: open and run the data files; inspect the contents in order to validate, clean, and/or enhance data for future use; look for missing documentation about codes used, the significance of "null" and "blank" values, or unclear acronyms. | 3 | 3 Check with published materials | Internal manual |
| 26 | Restructure | Organize and/or reformate poorly structured data files to clarify their meaning and importance. | 3 | 1.5 Reordering Variables | Internal manual |
| 27 | Software registry | Maintain copies of modern and obsolete versions of software (and any relevant code libraries) so that data may be opened/used overtime. | 1 | Survey participant | – |

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 28 | Contact information | Keep up-to-date contact information for the data authors and/or the contact persons in order to facilitate connection with third-party users. Often involves managing ephemeral information that will change over time. | 3 | Survey participant | – |
| 29 | Data citation | Display of a recommended bibliographic citation for a dataset to enable appropriate attribution by third-party users in order to formally incorporate data reuse as part of the scholarly ecosystem. | 3 | Survey participant | – |
| 30 | Data visualization | The presentation of pictorial and/or graphical representations of a data set used to identify patterns, detect errors, and/or demonstrate the extent of a data set to third party users. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 31 | Discovery Services | Services that incorporate machine-based search and retrieval functionality that help users identify what data exist, where the data are located, and how can they be accessed (e.g., full-text indexing or web optimization). | 3 | Survey participant | – |
| 32 | File download | Allow access to the data materials by authorized third parties. | 3 | Guide to Downloading Data | https://spirit. rikkyo.ac.jp/ csi/RUDA/ userguidance/ english/ datadownload/ Home.aspx |
| 33 | Full-text indexing | Enhance the data for discovery purposes by generating search-engine-optimized formats of the text inherent to the data. | 1 | Survey participant | – |

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 34 | Metadata brokerage | Active dissemination of a data set's metadata to search and discovery services (e.g., article databases, catalogs, web-based indexes) for federated search and discovery. | 3 | Survey participant | – |
| 35 | Restricted access | In order to maintain the privacy of research subjects without losing integral components of the data, some data access will be protected and/or mediated to individuals that meet predefined criteria. | 3 | Article 7. Prohibition of Use of User Account by Third Party "The User shall not cause or allow a third party to use the User account issued in response to his/her User registration. Further, the User shall make every effort to take necessary measures (such as password management) for preventing such use by a third party. " | https://spirit.rikkyo.ac.jp/csi/RUDA/userguidance/english/download/doc/agreement_e.pdf |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 36 | Embargo | To restrict or mediate access to a data set, usually for a set period of time. In some cases an embargo may be used to protect not only access, but any knowledge that the data exist. | 1 | Article 11. Data Access Permit Period "1 The User shall promptly erase any Data for which the data access permit period has elapsed, and thereafter shall not access such Data." "2 If the User desires to continue accessing the Data, the User shall submit a new application for such access. " | https://spirit. rikkyo.ac.jp/ csi/RUDA/ userguidance/ english/ download/ doc/ agreement_e. pdf |
| 37 | Terms of use | Information provided to end users of a data set that outline the requirements or conditions for use (e.g., a Creative Commons License). | 3 | Article 3. Agreement to Terms of Use "The User shall be deemed to apply for access to the services provided by the Center and access the same after giving consent to the terms and conditions of these Terms of Use. " | https://spirit. rikkyo.ac.jp/ csi/RUDA/ userguidance/ english/ download/ doc/ agreement_e. pdf |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 38 | Use analytics | Monitor and record how often data are viewed, requested, and/or downloaded. Track and report reuse metrics, such as data citations and impact measures for the data over time. | 2 | Usage report | https://spirit.rikkyo.ac.jp/csi/RUDA/userguidance/english/report/Home.aspx |
| 39 | Cease data curation | Plan for any contingencies that will ultimately terminate access to the data. For example, providing tombstones or metadata records for data that have been deselected and removed from stewardship. | 1 | Survey participant | – |
| 40 | Migration | Monitor and anticipate file format obsolescence and, as needed, transform obsolete file formats to new formats as standards and use dictate. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 41 | Emulation | Provide legacy system configurations in modern equipment in order to ensure long-term usability of data. (E.g., arcade games emulated on modern web-browsers) | 1 | Survey participant | – |
| 42 | Secure storage | Data files are properly stored in a well-configured (in terms of hardware and software) storage environment that is routinely backed-up and physically protected. Perform routine fixity checks (to detect degradation or loss) and provide recovery services as needed. | 1 | Survey participant | – |
| 43 | File audit | Periodic review of the digital integrity of the data files and taking action when needed to protect data from digital erosion (e.g., bitrot) and/or hardware failure. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 44 | Repository certification | The technical and administrative capacities of the repository undergo review through a transparent and well-documented process by a trusted third-party accreditation body (e.g., TRAC, or Data Seal of Approval). | 1 | Survey participant | – |
| 45 | Succession planning | Planning for contingency, and/or escrow arrangements, in the case that the repository (or other entity responsible) ceases to operate or the institution substantially changes its scope. | 1 | Survey participant | – |
| 46 | Technology monitoring and Refresh | Formal, periodic review and assessment to ensure responsiveness to technological developments and evolving requirements of the digital infrastructure and hardware storing the data. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 47 | Versioning | Provide mechanisms to ingest new versions of the data overtime that includes metadata describing the version history and any changes made for each version. | 3 | Survey participant | – |

*End of table*

Table 7: List of data curation process description rationale (IUGONET).

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 1 | Authentication | The process of confirming the identity of a person, generally the depositor, who is contributing data to the data repository. (e.g., password authentication or authorization via digital signature). Used for tracking provenance of the data files. | 3 | If you are interested, please contact with the IUGONET members. | http://www.iugonet.org/product/metadata.jsp?lang=en |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 2 | Chain of custody | Intentional recording of provenance metadata of the files (e.g., metadata about who created the file, when it was last edited, etc.) in order to preserve file authenticity when data are transferred to third-parties. | 3 | Survey participant | – |
| 3 | Deposit agreement | The certification by the data author (or depositor) that the data conform to all policies and conditions (e.g., do not violate any legal restrictions placed on the data) and are fit for deposit into the repository. A deposit agreement may also include rights transfer to the repository for ongoing stewardship. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 4 | Document-ation | Information describing any necessary information to use and understand the data. Documentation may be structured (e.g., a code book) or unstructured (e.g., a plain text "Readme" file). | 3 | "When you first contact us, we will ask you the following questions regarding the data you plan to register: The name of the data set you are registering, A brief description of the data set, Name and affiliation of the Principal Investigator, Name and affiliation of the person responsible for creating the metadata" | http://www.iugonet.org/data/manual/IUGONET_metadata_manual_v2_20170418.pdf |
| 5 | File Valida-tion | A computational process to ensure that the intended data transfer to a repository was perfect and complete using means such as generating and validating file checksums (e.g., test if a digital file has changed at the bit level) and format validation to ensure that file types match their extensions. | 2 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 6 | Metadata | Information about a data set that is structured (often in machine-readable format) for purposes of search and retrieval. Metadata elements may include basic information (e.g. title, author, date created, etc.) and/or specific elements inherent to datasets (e.g., spatial coverage, time periods). | 3 | The IUGONET Common Metadata Format (http://www.iugonet.org/mdformat.html) allows for the creation of metadata describing not only a single dataset, but also instruments, observation sites, human resources, and databases of real data. In addition to metadata describing a single data set, metadata describing instruments, observation sites, human resources, and databases of actual data are created independently and linked to each other. | http://www.iugonet.org/data/manual/IUGONET_metadata_manual_v2_20170418.pdf |
| 7 | Rights management | The process of tracking and managing ownership and copyright inherent to a data set as well as monitoring conditions and policies for access and reuse (e.g., licenses and data use agreements). | 3 | Acknowledgment statement required when using the data set. | http://www.iugonet.org/data/manual/IUGONET_metadata_manual_v2_20170418.pdf |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 8 | Risk management | The process of reviewing data for known risks such as confidentiality issues inherent to human subjects data, sensitive information (e.g., sexual histories, credit card information) or data regulated by law (e.g. HIPAA, FERPA) and taking actions to reject or facilitate remediation (e.g., de-identification services) when necessary. | 3 | Access policy for the actual dataset (open to the public, restricted access, closed, etc.). | http://www.iugonet.org/data/manual/IUGONET_metadata_manual_v2_20170418.pdf |
| 9 | Selection | The result of a successful appraisal. The data are determined appropriate for acceptance and ingest into the repository according to local collection policy and practice. | 3 | IUGONET welcomes your registration of the metadata of your observation data to our metadata database. | http://www.iugonet.org/product/metadata.jsp?lang=en |

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 10 | Arrangement and description | The re-organization of files (e.g., new folder directory structure) in a dataset that may also involve the creation of new file names, file descriptions, and the recording of technical metadata inherent to the files (e.g., date last modified). | 3 | Survey participant | – |
| 11 | Code review | Run and validate computer code (e.g., look for missing files and/or errors) in order to find mistakes overlooked in the initial development phase, improving the overall quality of software. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 12 | Contextual-ize | Use metadata to link the data set to related publications, dissertations, and/or projects that provide added context to how the data were generated and why. | 3 | The IUGONET Common Metadata Format (http://www.iugonet.org/ mdformat.html) allows for the creation of metadata describing not only a single dataset, but also instruments, observation sites, human resources, and databases of real data. In addition to metadata describing a single data set, metadata describing instruments, observation sites, human resources, and databases of actual data are created independently and linked to each other. | http://www. iugonet.org/ data/manual/ IUGONET_ metadata_ manual_v2_ 20170418.pdf |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 13 | Conversion (Analog) | In effort to increase the usability of a data set, the information is transferred into digital file formats (e.g., analog data keyed into a database). Note: digital conversion is also used to convert "fixed" data (e.g., PDF formats) into machine-readable formats. | 1 | Survey participant | – |
| 14 | Curation log | A written record of any changes made to the data during the curation process and by whom. File is often preserved as part of the overall record. | 1 | Survey participant | – |
| 15 | Data cleaning | A process used to improve data quality by detecting and correcting (or removing) defects & errors in data. | 3 | Survey participant | – |
| 16 | Deidentification | Redacting or removing personally identifiable or protected information (e.g., sensitive geographic locations) from a dataset prior to sharing with third-parties. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 17 | File format transformations | Transform files into open, non-proprietary file formats that broaden the potential for long-term reuse and ensure that additional preservation actions might be taken in the future. Note: Retention of the original file formats may be necessary if data transfer is not perfect. | 3 | Survey participant | – |
| 18 | Transcoding | With audio and video files, detect technical metadata (min resolution, audio/video codec) and encode files in ways that optimize reuse and long-term preservation actions. (E.g, Convert QuickTime files to MPEG4). | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 19 | File inventory or manifest | The data files are inspected periodically and the number, file types (extensions), and file sizes of the data are understood and documented. Any missing, duplicate, or corrupt (e.g., unable to open) files are discovered. | 3 | Survey participant | – |
| 20 | File renaming | To rename files in a dataset, often to standardize and/or reflect important metadata. | 3 | Survey participant | – |
| 21 | Indexing | Verify all metadata provided by the author and crosswalk to descriptive and administrative metadata compliant with a standard format for repository interoperability. | 3 | 3 How to submit your created metadata and confirm registration | http://www.iugonet.org/data/manual/IUGONET_metadata_manual_v2_20170418.pdf |
| 22 | Interoperability | Formatting the data using a disciplinary standard for better integration with other datasets and/or systems. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 23 | Peer-review | The review of a data set by an expert with similar credentials and subject knowledge as the data creator for the purposes of validating the soundness and trustworthiness of the file contents. | 1 | Survey participant | – |
| 24 | Persistent Identifier | A URL (or Uniform Resource Locator) that is monitored by an authority to ensure a stable web location for consistent citation and long-term discoverability. Provides redirection when necessary. E.g., a Digital Object Identifier or DOI. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 25 | Quality assurance | Ensure that all documentation and metadata are comprehensive and complete. Example actions might include: open and run the data files; inspect the contents in order to validate, clean, and/or enhance data for future use; look for missing documentation about codes used, the significance of "null" and "blank" values, or unclear acronyms. | 3 | The sent metadata will be run through a check script on the IUGONET server. | http://www.iugonet.org/data/manual/IUGONET_metadata_manual_v2_20170418.pdf |
| 26 | Restructure | Organize and/or reformate poorly structured data files to clarify their meaning and importance. | 1 | Survey participant | – |
| 27 | Software registry | Maintain copies of modern and obsolete versions of software (and any relevant code libraries) so that data may be opened/used overtime. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 28 | Contact information | Keep up-to-date contact information for the data authors and/or the contact persons in order to facilitate connection with third-party users. Often involves managing ephemeral information that will change over time. | 3 | Survey participant | – |
| 29 | Data citation | Display of a recommended bibliographic citation for a dataset to enable appropriate attribution by third-party users in order to formally incorporate data reuse as part of the scholarly ecosystem. | 1 | Survey participant | – |
| 30 | Data visualization | The presentation of pictorial and/or graphical representations of a data set used to identify patterns, detect errors, and/or demonstrate the extent of a data set to third party users. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 31 | Discovery Services | Services that incorporate machine-based search and retrieval functionality that help users identify what data exist, where the data are located, and how can they be accessed (e.g., full-text indexing or web optimization). | 1 | Survey participant | – |
| 32 | File download | Allow access to the data materials by authorized third parties. | 3 | 3 How to submit your created metadata and confirm registration | http://www.iugonet.org/data/manual/IUGONET_metadata_manual_v2_20170418.pdf |
| 33 | Full-text indexing | Enhance the data for discovery purposes by generating search-engine-optimized formats of the text inherent to the data. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 34 | Metadata brokerage | Active dissemination of a data set's metadata to search and discovery services (e.g., article databases, catalogs, web-based indexes) for federated search and discovery. | 1 | Survey participant | – |
| 35 | Restricted access | In order to maintain the privacy of research subjects without losing integral components of the data, some data access will be protected and/or mediated to individuals that meet predefined criteria. | 1 | Survey participant | – |
| 36 | Embargo | To restrict or mediate access to a data set, usually for a set period of time. In some cases an embargo may be used to protect not only access, but any knowledge that the data exist. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 37 | Terms of use | Information provided to end users of a data set that outline the requirements or conditions for use (e.g., a Creative Commons License). | 3 | Acknowledgment statement required when using the data set. | http://www.iugonet.org/data/manual/IUGONET_metadata_manual_v2_20170418.pdf |
| 38 | Use analytics | Monitor and record how often data are viewed, requested, and/or downloaded. Track and report reuse metrics, such as data citations and impact measures for the data over time. | 2 | Survey participant | – |
| 39 | Cease data curation | Plan for any contingencies that will ultimately terminate access to the data. For example, providing tombstones or metadata records for data that have been deselected and removed from stewardship. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 40 | Migration | Monitor and anticipate file format obsolescence and, as needed, transform obsolete file formats to new formats as standards and use dictate. | 1 | Survey participant | – |
| 41 | Emulation | Provide legacy system configurations in modern equipment in order to ensure long-term usability of data. (E.g., arcade games emulated on modern web-browsers) | 1 | Survey participant | – |
| 42 | Secure storage | Data files are properly stored in a well-configured (in terms of hardware and software) storage environment that is routinely backed-up and physically protected. Perform routine fixity checks (to detect degradation or loss) and provide recovery services as needed. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 43 | File audit | Periodic review of the digital integrity of the data files and taking action when needed to protect data from digital erosion (e.g., bitrot) and/or hardware failure. | 1 | Survey participant | – |
| 44 | Repository certification | The technical and administrative capacities of the repository undergo review through a transparent and well-documented process by a trusted third-party accreditation body (e.g., TRAC, or Data Seal of Approval). | 1 | Survey participant | – |
| 45 | Succession planning | Planning for contingency, and/or escrow arrangements, in the case that the repository (or other entity responsible) ceases to operate or the institution substantially changes its scope. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 46 | Technology monitoring and Refresh | Formal, periodic review and assessment to ensure responsiveness to technological developments and evolving requirements of the digital infrastructure and hardware storing the data. | 2 | Survey participant | – |
| 47 | Versioning | Provide mechanisms to ingest new versions of the data overtime that includes metadata describing the version history and any changes made for each version. | 3 | Survey participant | – |

*End of table*

Table 8: List of data curation process description rationale (NBDC Archive).

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 1 | Authentication | The process of confirming the identity of a person, generally the depositor, who is contributing data to the data repository. (e.g., password authentication or authorization via digital signature). Used for tracking provenance of the data files. | 3 | If you have any question or request about database archive and database deposit, please contact the follwing: | https://dbarchive.biosciencedbc.jp/contents-en/contact/contact.html |
| 2 | Chain of custody | Intentional recording of provenance metadata of the files (e.g., metadata about who created the file, when it was last edited, etc.) in order to preserve file authenticity when data are transferred to third-parties. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 3 | Deposit agreement | The certification by the data author (or depositor) that the data conform to all policies and conditions (e.g., do not violate any legal restrictions placed on the data) and are fit for deposit into the repository. A deposit agreement may also include rights transfer to the repository for ongoing stewardship. | 3 | The NBDC will check the archived data, metadata, and the wording of the license agreement, and approve the release of the archive. | https: //dbarchive. biosciencedbc. jp/files/nbdc_ dbarchive_ guidelines.pdf |
| 4 | Document-ation | Information describing any necessary information to use and understand the data. Documentation may be structured (e.g., a code book) or unstructured (e.g., a plain text "Readme" file). | 3 | 7 Metadata creation | https: //dbarchive. biosciencedbc. jp/files/nbdc_ dbarchive_ guidelines.pdf |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 5 | File Validation | A computational process to ensure that the intended data transfer to a repository was perfect and complete using means such as generating and validating file checksums (e.g., test if a digital file has changed at the bit level) and format validation to ensure that file types match their extensions. | 1 | Survey participant | – |
| 6 | Metadata | Information about a data set that is structured (often in machine-readable format) for purposes of search and retrieval. Metadata elements may include basic information (e.g. title, author, date created, etc.) and/or specific elements inherent to datasets (e.g., spatial coverage, time periods). | 3 | 7 Metadata creation | https: //dbarchive. biosciencedbc. jp/files/nbdc_ dbarchive_ guidelines.pdf |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 7 | Rights management | The process of tracking and managing ownership and copyright inherent to a data set as well as monitoring conditions and policies for access and reuse (e.g., licenses and data use agreements). | 3 | 8 License decision | https://dbarchive.biosciencedbc.jp/files/nbdc_dbarchive_guidelines.pdf |
| 8 | Risk management | The process of reviewing data for known risks such as confidentiality issues inherent to human subjects data, sensitive information (e.g., sexual histories, credit card information) or data regulated by law (e.g. HIPAA, FERPA) and taking actions to reject or facilitate remediation (e.g., de-identification services) when necessary. | 3 | 3.1.2 Databases that are not accepted for archiving | https://dbarchive.biosciencedbc.jp/files/nbdc_dbarchive_guidelines.pdf |
| 9 | Selection | The result of a successful appraisal. The data are determined appropriate for acceptance and ingest into the repository according to local collection policy and practice. | 3 | 3.1 Archive Acceptance Policy | https://dbarchive.biosciencedbc.jp/files/nbdc_dbarchive_guidelines.pdf |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 10 | Arrangem-ent and description | The re-organization of files (e.g., new folder direc-tory structure) in a dataset that may also involve the creation of new file names, file descriptions, and the recording of tech-nical metadata inherent to the files (e.g., date last mod-ified). | 3 | 5.1.2 Designing the File Structure | https://dbarchive.biosciencedbc.jp/files/nbdc_dbarchive_guidelines.pdf |
| 11 | Code review | Run and validate computer code (e.g., look for missing files and/or errors) in or-der to find mistakes over-looked in the initial devel-opment phase, improving the overall quality of soft-ware. | 3 | Survey participant | – |
| 12 | Contextual-ize | Use metadata to link the data set to related publications, dissertations, and/or projects that provide added context to how the data were generated and why. | 3 | Create detailed informa-tion (metadata) about the contents of the database and each item in the archive file. | https://dbarchive.biosciencedbc.jp/files/nbdc_dbarchive_guidelines.pdf |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 13 | Conversion (Analog) | In effort to increase the usability of a data set, the information is transferred into digital file formats (e.g., analog data keyed into a database). Note: digital conversion is also used to convert "fixed" data (e.g., PDF formats) into machine-readable formats. | 1 | Survey participant | – |
| 14 | Curation log | A written record of any changes made to the data during the curation process and by whom. File is often preserved as part of the overall record. | 3 | Survey participant | – |
| 15 | Data cleaning | A process used to improve data quality by detecting and correcting (or removing) defects & errors in data. | 3 | Survey participant | – |
| 16 | Deidentification | Redacting or removing personally identifiable or protected information (e.g., sensitive geographic locations) from a dataset prior to sharing with third-parties. | 3 | 3.1 Archive Acceptance Policy | https://dbarchive.biosciencedbc.jp/files/nbdc_dbarchive_guidelines.pdf |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|---------------------|------|
| 17 | File format transformations | Transform files into open, non-proprietary file formats that broaden the potential for long-term reuse and ensure that additional preservation actions might be taken in the future. Note: Retention of the original file formats may be necessary if data transfer is not perfect. | 3 | Convert the original data into the archive file designed in (1). | https://dbarchive.biosciencedbc.jp/files/nbdc_dbarchive_guidelines.pdf |
| 18 | Transcoding | With audio and video files, detect technical metadata (min resolution, audio/video codec) and encode files in ways that optimize reuse and long-term preservation actions. (E.g, Convert QuickTime files to MPEG4). | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 19 | File inventory or manifest | The data files are inspected periodically and the number, file types (extensions), and file sizes of the data are understood and documented. Any missing, duplicate, or corrupt (e.g., unable to open) files are discovered. | 3 | Survey participant | – |
| 20 | File renaming | To rename files in a dataset, often to standardize and/or reflect important metadata. | 3 | 5.2.2 Naming conventions for archive files | https://dbarchive.biosciencedbc.jp/files/nbdc_dbarchive_guidelines.pdf |
| 21 | Indexing | Verify all metadata provided by the author and crosswalk to descriptive and administrative metadata compliant with a standard format for repository interoperability. | 3 | Survey participant | – |
| 22 | Interoperability | Formatting the data using a disciplinary standard for better integration with other datasets and/or systems. | 3 | 5.1.3 Determination of data items | https://dbarchive.biosciencedbc.jp/files/nbdc_dbarchive_guidelines.pdf |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|---------------------|------|
| 23 | Peer-review | The review of a data set by an expert with similar credentials and subject knowledge as the data creator for the purposes of validating the soundness and trustworthiness of the file contents. | 1 | Survey participant | – |
| 24 | Persistent Identifier | A URL (or Uniform Resource Locator) that is monitored by an authority to ensure a stable web location for consistent citation and long-term discoverability. Provides redirection when necessary. E.g., a Digital Object Identifier or DOI. | 3 | DOI in Life Science Database Archive | https://dbarchive.biosciencedbc.jp/contents-en/doi/list.html |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 25 | Quality assurance | Ensure that all documentation and metadata are comprehensive and complete. Example actions might include: open and run the data files; inspect the contents in order to validate, clean, and/or enhance data for future use; look for missing documentation about codes used, the significance of "null" and "blank" values, or unclear acronyms. | 3 | The NBDC will check the archived data, metadata, and the wording of the license agreement, and approve the release of the archive. | https://dbarchive.biosciencedbc.jp/files/nbdc_dbarchive_guidelines.pdf |
| 26 | Restructure | Organize and/or reformate poorly structured data files to clarify their meaning and importance. | 3 | "5.1.3 Determination of data items ; 5.1.3.4 Integration and division of items" | https://dbarchive.biosciencedbc.jp/files/nbdc_dbarchive_guidelines.pdf |
| 27 | Software registry | Maintain copies of modern and obsolete versions of software (and any relevant code libraries) so that data may be opened/used overtime. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 28 | Contact information | Keep up-to-date contact information for the data authors and/or the contact persons in order to facilitate connection with third-party users. Often involves managing ephemeral information that will change over time. | 3 | 7.1 Database Metadata | https://dbarchive.biosciencedbc.jp/files/nbdc_dbarchive_guidelines.pdf |
| 29 | Data citation | Display of a recommended bibliographic citation for a dataset to enable appropriate attribution by third-party users in order to formally incorporate data reuse as part of the scholarly ecosystem. | 1 | Survey participant | – |
| 30 | Data visualization | The presentation of pictorial and/or graphical representations of a data set used to identify patterns, detect errors, and/or demonstrate the extent of a data set to third party users. | 1 | Survey participant | – |

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 31 | Discovery Services | Services that incorporate machine-based search and retrieval functionality that help users identify what data exist, where the data are located, and how can they be accessed (e.g., full-text indexing or web optimization). | 3 | 6 Create a simple search site (optional) | https://dbarchive.biosciencedbc.jp/files/nbdc_dbarchive_guidelines.pdf |
| 32 | File download | Allow access to the data materials by authorized third parties. | 3 | About the Life Science Database Archive | https://dbarchive.biosciencedbc.jp/contents-en/about/about.html |
| 33 | Full-text indexing | Enhance the data for discovery purposes by generating search-engine-optimized formats of the text inherent to the data. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|-----------|-------------|---------------------|------|
| 34 | Metadata brokerage | Active dissemination of a data set's metadata to search and discovery services (e.g., article databases, catalogs, web-based indexes) for federated search and discovery. | 3 | 6 Create a simple search site (optional) | https://dbarchive.biosciencedbc.jp/files/nbdc_dbarchive_guidelines.pdf |
| 35 | Restricted access | In order to maintain the privacy of research subjects without losing integral components of the data, some data access will be protected and/or mediated to individuals that meet predefined criteria. | 3 | 3.1.2 Databases that are not accepted for archiving | https://dbarchive.biosciencedbc.jp/files/nbdc_dbarchive_guidelines.pdf |
| 36 | Embargo | To restrict or mediate access to a data set, usually for a set period of time. In some cases an embargo may be used to protect not only access, but any knowledge that the data exist. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 37 | Terms of use | Information provided to end users of a data set that outline the requirements or conditions for use (e.g., a Creative Commons License). | 3 | 8 License decision | https://dbarchive.biosciencedbc.jp/files/nbdc_dbarchive_guidelines.pdf |
| 38 | Use analytics | Monitor and record how often data are viewed, requested, and/or downloaded. Track and report reuse metrics, such as data citations and impact measures for the data over time. | 3 | Survey participant | – |
| 39 | Cease data curation | Plan for any contingencies that will ultimately terminate access to the data. For example, providing tombstones or metadata records for data that have been deselected and removed from stewardship. | 1 | Survey participant | – |
| 40 | Migration | Monitor and anticipate file format obsolescence and, as needed, transform obsolete file formats to new formats as standards and use dictate. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 41 | Emulation | Provide legacy system configurations in modern equipment in order to ensure long-term usability of data. (E.g., arcade games emulated on modern web-browsers) | 1 | Survey participant | – |
| 42 | Secure storage | Data files are properly stored in a well-configured (in terms of hardware and software) storage environment that is routinely backed-up and physically protected. Perform routine fixity checks (to detect degradation or loss) and provide recovery services as needed. | 3 | Survey participant | – |
| 43 | File audit | Periodic review of the digital integrity of the data files and taking action when needed to protect data from digital erosion (e.g., bitrot) and/or hardware failure. | 1 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|----|----------|------------|-------------|----------------------|------|
| 44 | Repository certification | The technical and administrative capacities of the repository undergo review through a transparent and well-documented process by a trusted third-party accreditation body (e.g., TRAC, or Data Seal of Approval). | 1 | Survey participant | – |
| 45 | Succession planning | Planning for contingency, and/or escrow arrangements, in the case that the repository (or other entity responsible) ceases to operate or the institution substantially changes its scope. | 3 | Survey participant | – |
| 46 | Technology monitoring and Refresh | Formal, periodic review and assessment to ensure responsiveness to technological developments and evolving requirements of the digital infrastructure and hardware storing the data. | 3 | Survey participant | – |

*Continued from previous page*

| No | Activity | Definition | Imp. levels | Rationale description | Link |
|---|---|---|---|---|---|
| 47 | Versioning | Provide mechanisms to ingest new versions of the data overtime that includes metadata describing the version history and any changes made for each version. | 3 | Survey participant | – |

*End of table*

# Appendix 3. Data Curation Process Ontology

In Section 3.5, we developed the Data Curation Process Ontology. The ontology aims to represent commonalities and differences in data curation activity's structure. The domain to be covered by this ontology is that of data curation. Providing the structured data curation activity in a machine-readable format can support the knowledge-sharing process between humans and information systems in a scalable manner. It would be desirable to maintain the ontology through the collaboration of data curators in each field and ontologists who deal with knowledge sharing in information systems. The description of this ontology is based on the vocabulary of the PROV ontology (https://www.w3.org/TR/prov-o/) endorsed by W3C. The described activities were defined by the Data Curation Network, with some additions. The ontology is available at the following URL (https://purl.archive.org/curation-ontology). To open the ontology the program Protégé (https://protege.stanford.edu/) is recommended.

lodac / **curation-ontology** Public

<> Code  ⊙ Issues  ⅋ Pull requests  ▷ Actions  ⊞ Projects  ▭ Wiki  ⊘ Security  ∿ Insights  ⚙ Settings

ⵥ main ▾  curation-ontology / src / ontology / **data-curation-process-ontology-1.0.0.owl**  Go to file  ···

y-minamiyama Add files via upload  Latest commit 187930c on 20 Apr  ⏱ History

🗅 1 contributor

2560 lines (2556 sloc)  105 KB  Raw  Blame  ✎ ▾  🗍 🗑

```
 1   <?xml version="1.0"?>
 2   <Ontology xmlns="http://www.w3.org/2002/07/owl#"
 3       xml:base="http://www.semanticweb.org/2022/04/data-curation-process"
 4       xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
 5       xmlns:xml="http://www.w3.org/XML/1998/namespace"
 6       xmlns:xsd="http://www.w3.org/2001/XMLSchema#"
 7       xmlns:rdfs="http://www.w3.org/2000/01/rdf-schema#"
 8       ontologyIRI="http://www.semanticweb.org/2022/04/data-curation-process"
 9       versionIRI="http://www.semanticweb.org/2022/04/data-curation-process/1.0.0">
10   <Prefix name="" IRI="http://www.semanticweb.org/2022/04/data-curation-process"/>
11   <Prefix name="owl" IRI="http://www.w3.org/2002/07/owl#"/>
12   <Prefix name="rdf" IRI="http://www.w3.org/1999/02/22-rdf-syntax-ns#"/>
13   <Prefix name="xml" IRI="http://www.w3.org/XML/1998/namespace"/>
14   <Prefix name="xsd" IRI="http://www.w3.org/2001/XMLSchema#"/>
15   <Prefix name="rdfs" IRI="http://www.w3.org/2000/01/rdf-schema#"/>
16   <Prefix name="prov-o-20130430" IRI="http://www.w3.org/ns/prov-o-20130430#"/>
17   <Import>http://www.w3.org/ns/prov-o-20130430</Import>
18   <Import>http://xmlns.com/foaf/0.1/</Import>
19   <Declaration>
20   <Class IRI="#ActivatingMetadataBrokerage"/>
21   </Declaration>
22   <Declaration>
23   <Class IRI="#ActualDataProcessing"/>
24   </Declaration>
25   <Declaration>
26   <Class IRI="#AllowingFileDownload"/>
27   </Declaration>
28   <Declaration>
29   <Class IRI="#Appraisal"/>
30   </Declaration>
31   <Declaration>
32   <Class IRI="#ArrangementAndDescription"/>
33   </Declaration>
34   <Declaration>
35   <Class IRI="#Authentication"/>
36   </Declaration>
37   <Declaration>
38   <Class IRI="#CeaseDataCuration"/>
39   </Declaration>
40   <Declaration>
41   <Class IRI="#ChainOfCustody"/>
42   </Declaration>
43   <Declaration>
44   <Class IRI="#CodeReview"/>
45   </Declaration>
46   <Declaration>
47   <Class IRI="#ConnectingDiscoveryServices"/>
48   </Declaration>
49   <Declaration>
50   <Class IRI="#Contextualization"/>
51   </Declaration>
52   <Declaration>
53   <Class IRI="#Contracts"/>
54   </Declaration>
55   <Declaration>
56   <Class IRI="#Conversion"/>
57   </Declaration>
```

```
58   <Declaration>
59   <Class IRI="#Copyright"/>
60   </Declaration>
61   <Declaration>
62   <Class IRI="#CreatingLandingPage"/>
63   </Declaration>
64   <Declaration>
65   <Class IRI="#DataCleaning"/>
66   </Declaration>
67   <Declaration>
68   <Class IRI="#DataCurationActivities"/>
69   </Declaration>
70   <Declaration>
71   <Class IRI="#DataEvaluation"/>
72   </Declaration>
73   <Declaration>
74   <Class IRI="#DataLifecycleActivities"/>
75   </Declaration>
76   <Declaration>
77   <Class IRI="#DataPolicy"/>
78   </Declaration>
79   <Declaration>
80   <Class IRI="#DataPreservationActivities"/>
81   </Declaration>
82   <Declaration>
83   <Class IRI="#DataProcessing"/>
84   </Declaration>
85   <Declaration>
86   <Class IRI="#DataPublishing"/>
87   </Declaration>
88   <Declaration>
89   <Class IRI="#DataVisualization"/>
90   </Declaration>
91   <Declaration>
92   <Class IRI="#Deidentification"/>
93   </Declaration>
94   <Declaration>
95   <Class IRI="#DepositAgreement"/>
96   </Declaration>
97   <Declaration>
98   <Class IRI="#Digitization"/>
99   </Declaration>
100  <Declaration>
101  <Class IRI="#DiplomaticOrNationalSecurity"/>
102  </Declaration>
103  <Declaration>
104  <Class IRI="#DisciplinaryCustoms"/>
105  </Declaration>
106  <Declaration>
107  <Class IRI="#DisplayingDataCitation"/>
108  </Declaration>
109  <Declaration>
110  <Class IRI="#Documentation"/>
111  </Declaration>
112  <Declaration>
113  <Class IRI="#Embargo"/>
114  </Declaration>
115  <Declaration>
116  <Class IRI="#Emulation"/>
117  </Declaration>
118  <Declaration>
119  <Class IRI="#FileAuditing"/>
120  </Declaration>
121  <Declaration>
122  <Class IRI="#FileFormatTransformation"/>
123  </Declaration>
124  <Declaration>
125  <Class IRI="#FileInventoryOrManifest"/>
126  </Declaration>
127  <Declaration>
128  <Class IRI="#FileRenaming"/>
129  </Declaration>
130  <Declaration>
131  <Class IRI="#FileValidation"/>
```

```
132   </Declaration>
133   <Declaration>
134   <Class IRI="#FormulateSuccessionPlanning"/>
135   </Declaration>
136   <Declaration>
137   <Class IRI="#GeneratingFulltextIndexing"/>
138   </Declaration>
139   <Declaration>
140   <Class IRI="#IPR"/>
141   </Declaration>
142   <Declaration>
143   <Class IRI="#Indexing"/>
144   </Declaration>
145   <Declaration>
146   <Class IRI="#InformedConsent"/>
147   </Declaration>
148   <Declaration>
149   <Class IRI="#Ingest"/>
150   </Declaration>
151   <Declaration>
152   <Class IRI="#Interoperability"/>
153   </Declaration>
154   <Declaration>
155   <Class IRI="#MaintainingContactInformation"/>
156   </Declaration>
157   <Declaration>
158   <Class IRI="#MetadataGeneration"/>
159   </Declaration>
160   <Declaration>
161   <Class IRI="#MetadataProcessing"/>
162   </Declaration>
163   <Declaration>
164   <Class IRI="#Migration"/>
165   </Declaration>
166   <Declaration>
167   <Class IRI="#MintingPersistentIdentifier"/>
168   </Declaration>
169   <Declaration>
170   <Class IRI="#PeerReview"/>
171   </Declaration>
172   <Declaration>
173   <Class IRI="#PersonalInformation"/>
174   </Declaration>
175   <Declaration>
176   <Class IRI="#ProvidingRestrictedAccess"/>
177   </Declaration>
178   <Declaration>
179   <Class IRI="#QualityAssurance"/>
180   </Declaration>
181   <Declaration>
182   <Class IRI="#RegisteringSoftware"/>
183   </Declaration>
184   <Declaration>
185   <Class IRI="#Restructure"/>
186   </Declaration>
187   <Declaration>
188   <Class IRI="#RightsManagement"/>
189   </Declaration>
190   <Declaration>
191   <Class IRI="#RiskManagement"/>
192   </Declaration>
193   <Declaration>
194   <Class IRI="#SecuringStorage"/>
195   </Declaration>
196   <Declaration>
197   <Class IRI="#Selection"/>
198   </Declaration>
199   <Declaration>
200   <Class IRI="#SensitiveData"/>
201   </Declaration>
202   <Declaration>
203   <Class IRI="#SettingTermsOfUse"/>
204   </Declaration>
205   <Declaration>
```

```
206   <Class IRI="#SubmitData"/>
207   </Declaration>
208   <Declaration>
209   <Class IRI="#TechnologyMonitoringAndRefreshing"/>
210   </Declaration>
211   <Declaration>
212   <Class IRI="#TrackingUseAnalytics"/>
213   </Declaration>
214   <Declaration>
215   <Class IRI="#Transcoding"/>
216   </Declaration>
217   <Declaration>
218   <Class IRI="#TranscribingAndTranslatingData"/>
219   </Declaration>
220   <Declaration>
221   <Class IRI="#Versioning"/>
222   </Declaration>
223   <Declaration>
224   <Class IRI="#accessRestriction"/>
225   </Declaration>
226   <Declaration>
227   <Class IRI="#administrator"/>
228   </Declaration>
229   <Declaration>
230   <Class IRI="#agreement"/>
231   </Declaration>
232   <Declaration>
233   <Class IRI="#appropriateRepository"/>
234   </Declaration>
235   <Declaration>
236   <Class IRI="#authenticationResults"/>
237   </Declaration>
238   <Declaration>
239   <Class IRI="#backupData"/>
240   </Declaration>
241   <Declaration>
242   <Class IRI="#contactInformation"/>
243   </Declaration>
244   <Declaration>
245   <Class IRI="#costAndLeadTime"/>
246   </Declaration>
247   <Declaration>
248   <Class IRI="#creditOnTheResults"/>
249   </Declaration>
250   <Declaration>
251   <Class IRI="#curatedData"/>
252   </Declaration>
253   <Declaration>
254   <Class IRI="#curationLog"/>
255   </Declaration>
256   <Declaration>
257   <Class IRI="#curationRecord"/>
258   </Declaration>
259   <Declaration>
260   <Class IRI="#dataCurator"/>
261   </Declaration>
262   <Declaration>
263   <Class IRI="#dataDepositer"/>
264   </Declaration>
265   <Declaration>
266   <Class IRI="#dataDocument"/>
267   </Declaration>
268   <Declaration>
269   <Class IRI="#dataProcessingPolicy"/>
270   </Declaration>
271   <Declaration>
272   <Class IRI="#dataUser"/>
273   </Declaration>
274   <Declaration>
275   <Class IRI="#documentationPolicy"/>
276   </Declaration>
277   <Declaration>
278   <Class IRI="#evaluationResults"/>
279   </Declaration>
```

```
280    <Declaration>
281    <Class IRI="#expertiseNecessity"/>
282    </Declaration>
283    <Declaration>
284    <Class IRI="#externalConstraints"/>
285    </Declaration>
286    <Declaration>
287    <Class IRI="#externalServiceProvider"/>
288    </Declaration>
289    <Declaration>
290    <Class IRI="#feasibility"/>
291    </Declaration>
292    <Declaration>
293    <Class IRI="#feeForUse"/>
294    </Declaration>
295    <Declaration>
296    <Class IRI="#fileLocation"/>
297    </Declaration>
298    <Declaration>
299    <Class IRI="#fileValidationResults"/>
300    </Declaration>
301    <Declaration>
302    <Class IRI="#fullTextInformation"/>
303    </Declaration>
304    <Declaration>
305    <Class IRI="#imposeTheSameConditions"/>
306    </Declaration>
307    <Declaration>
308    <Class IRI="#improperUse"/>
309    </Declaration>
310    <Declaration>
311    <Class IRI="#indexingInformation"/>
312    </Declaration>
313    <Declaration>
314    <Class IRI="#landingPage"/>
315    </Declaration>
316    <Declaration>
317    <Class IRI="#linkingInformation"/>
318    </Declaration>
319    <Declaration>
320    <Class IRI="#metadata"/>
321    </Declaration>
322    <Declaration>
323    <Class IRI="#metadataCurator"/>
324    </Declaration>
325    <Declaration>
326    <Class IRI="#metadataSchema"/>
327    </Declaration>
328    <Declaration>
329    <Class IRI="#noDelivs"/>
330    </Declaration>
331    <Declaration>
332    <Class IRI="#nonCommercial"/>
333    </Declaration>
334    <Declaration>
335    <Class IRI="#peerReviewer"/>
336    </Declaration>
337    <Declaration>
338    <Class IRI="#persistentIdentifier"/>
339    </Declaration>
340    <Declaration>
341    <Class IRI="#policy"/>
342    </Declaration>
343    <Declaration>
344    <Class IRI="#preservationPolicy"/>
345    </Declaration>
346    <Declaration>
347    <Class IRI="#quality"/>
348    </Declaration>
349    <Declaration>
350    <Class IRI="#reporting"/>
351    </Declaration>
352    <Declaration>
353    <Class IRI="#repositorySystem"/>
```

```
354    </Declaration>
355    <Declaration>
356    <Class IRI="#researchData"/>
357    </Declaration>
358    <Declaration>
359    <Class IRI="#retrievalMetadata"/>
360    </Declaration>
361    <Declaration>
362    <Class IRI="#reusability"/>
363    </Declaration>
364    <Declaration>
365    <Class IRI="#secondaryUseProhibited"/>
366    </Declaration>
367    <Declaration>
368    <Class IRI="#selectionPolicy"/>
369    </Declaration>
370    <Declaration>
371    <Class IRI="#selectionResults"/>
372    </Declaration>
373    <Declaration>
374    <Class IRI="#size"/>
375    </Declaration>
376    <Declaration>
377    <Class IRI="#softwareRegistry"/>
378    </Declaration>
379    <Declaration>
380    <Class IRI="#sourceCode"/>
381    </Declaration>
382    <Declaration>
383    <Class IRI="#submittedData"/>
384    </Declaration>
385    <Declaration>
386    <Class IRI="#successionPlan"/>
387    </Declaration>
388    <Declaration>
389    <Class IRI="#systemAdministrator"/>
390    </Declaration>
391    <Declaration>
392    <Class IRI="#technicalInformation"/>
393    </Declaration>
394    <Declaration>
395    <Class IRI="#termsOfUse"/>
396    </Declaration>
397    <Declaration>
398    <Class IRI="#typeOfResource"/>
399    </Declaration>
400    <Declaration>
401    <Class IRI="#usageResults"/>
402    </Declaration>
403    <Declaration>
404    <Class IRI="#useLatestVersion"/>
405    </Declaration>
406    <Declaration>
407    <Class IRI="#versionInformation"/>
408    </Declaration>
409    <Declaration>
410    <Class IRI="#visualizedData"/>
411    </Declaration>
412    <Declaration>
413    <Class IRI="#waiver"/>
414    </Declaration>
415    <Declaration>
416    <ObjectProperty IRI="http://www.w3.org/ns/prov#Revision"/>
417    </Declaration>
418    <Declaration>
419    <AnnotationProperty IRI="http://www.w3.org/2004/02/skos/core#altLabel"/>
420    </Declaration>
421    <SubClassOf>
422    <Class IRI="#ActivatingMetadataBrokerage"/>
423    <Class IRI="#DataPublishing"/>
424    </SubClassOf>
425    <SubClassOf>
426    <Class IRI="#ActivatingMetadataBrokerage"/>
427    <ObjectSomeValuesFrom>
```

```
428  <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
429  <Class IRI="#indexingInformation"/>
430  </ObjectSomeValuesFrom>
431  </SubClassOf>
432  <SubClassOf>
433  <Class IRI="#ActivatingMetadataBrokerage"/>
434  <ObjectSomeValuesFrom>
435  <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
436  <Class IRI="#landingPage"/>
437  </ObjectSomeValuesFrom>
438  </SubClassOf>
439  <SubClassOf>
440  <Class IRI="#ActivatingMetadataBrokerage"/>
441  <ObjectSomeValuesFrom>
442  <ObjectProperty IRI="http://www.w3.org/ns/prov#wasAssociatedWith"/>
443  <Class IRI="#externalServiceProvider"/>
444  </ObjectSomeValuesFrom>
445  </SubClassOf>
446  <SubClassOf>
447  <Class IRI="#ActualDataProcessing"/>
448  <Class IRI="#DataProcessing"/>
449  </SubClassOf>
450  <SubClassOf>
451  <Class IRI="#ActualDataProcessing"/>
452  <ObjectSomeValuesFrom>
453  <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
454  <Class IRI="#visualizedData"/>
455  </ObjectSomeValuesFrom>
456  </SubClassOf>
457  <SubClassOf>
458  <Class IRI="#ActualDataProcessing"/>
459  <ObjectSomeValuesFrom>
460  <ObjectProperty IRI="http://www.w3.org/ns/prov#influenced"/>
461  <Class IRI="#accessRestriction"/>
462  </ObjectSomeValuesFrom>
463  </SubClassOf>
464  <SubClassOf>
465  <Class IRI="#ActualDataProcessing"/>
466  <ObjectSomeValuesFrom>
467  <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
468  <Class IRI="#dataDocument"/>
469  </ObjectSomeValuesFrom>
470  </SubClassOf>
471  <SubClassOf>
472  <Class IRI="#ActualDataProcessing"/>
473  <ObjectSomeValuesFrom>
474  <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
475  <Class IRI="#dataProcessingPolicy"/>
476  </ObjectSomeValuesFrom>
477  </SubClassOf>
478  <SubClassOf>
479  <Class IRI="#ActualDataProcessing"/>
480  <ObjectMinCardinality cardinality="1">
481  <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
482  <Class IRI="#curatedData"/>
483  </ObjectMinCardinality>
484  </SubClassOf>
485  <SubClassOf>
486  <Class IRI="#ActualDataProcessing"/>
487  <ObjectMinCardinality cardinality="1">
488  <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
489  <Class IRI="#submittedData"/>
490  </ObjectMinCardinality>
491  </SubClassOf>
492  <SubClassOf>
493  <Class IRI="#ActualDataProcessing"/>
494  <ObjectMinCardinality cardinality="1">
495  <ObjectProperty IRI="http://www.w3.org/ns/prov#wasAssociatedWith"/>
496  <Class IRI="#dataCurator"/>
497  </ObjectMinCardinality>
498  </SubClassOf>
499  <SubClassOf>
500  <Class IRI="#AllowingFileDownload"/>
501  <Class IRI="#DataPublishing"/>
```

```
502   </SubClassOf>
503   <SubClassOf>
504   <Class IRI="#AllowingFileDownload"/>
505   <ObjectSomeValuesFrom>
506   <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
507   <Class IRI="#versionInformation"/>
508   </ObjectSomeValuesFrom>
509   </SubClassOf>
510   <SubClassOf>
511   <Class IRI="#AllowingFileDownload"/>
512   <ObjectSomeValuesFrom>
513   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
514   <Class IRI="#accessRestriction"/>
515   </ObjectSomeValuesFrom>
516   </SubClassOf>
517   <SubClassOf>
518   <Class IRI="#AllowingFileDownload"/>
519   <ObjectSomeValuesFrom>
520   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
521   <Class IRI="#dataDocument"/>
522   </ObjectSomeValuesFrom>
523   </SubClassOf>
524   <SubClassOf>
525   <Class IRI="#AllowingFileDownload"/>
526   <ObjectSomeValuesFrom>
527   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
528   <Class IRI="#researchData"/>
529   </ObjectSomeValuesFrom>
530   </SubClassOf>
531   <SubClassOf>
532   <Class IRI="#AllowingFileDownload"/>
533   <ObjectSomeValuesFrom>
534   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
535   <Class IRI="#sourceCode"/>
536   </ObjectSomeValuesFrom>
537   </SubClassOf>
538   <SubClassOf>
539   <Class IRI="#AllowingFileDownload"/>
540   <ObjectMinCardinality cardinality="1">
541   <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
542   <Class IRI="#fileLocation"/>
543   </ObjectMinCardinality>
544   </SubClassOf>
545   <SubClassOf>
546   <Class IRI="#Appraisal"/>
547   <Class IRI="#DataCurationActivities"/>
548   </SubClassOf>
549   <SubClassOf>
550   <Class IRI="#Appraisal"/>
551   <ObjectMinCardinality cardinality="1">
552   <ObjectProperty IRI="http://www.w3.org/ns/prov#wasAssociatedWith"/>
553   <ObjectUnionOf>
554   <Class IRI="#administrator"/>
555   <Class IRI="#dataCurator"/>
556   <Class IRI="#repositorySystem"/>
557   </ObjectUnionOf>
558   </ObjectMinCardinality>
559   </SubClassOf>
560   <SubClassOf>
561   <Class IRI="#ArrangementAndDescription"/>
562   <Class IRI="#ActualDataProcessing"/>
563   </SubClassOf>
564   <SubClassOf>
565   <Class IRI="#Authentication"/>
566   <Class IRI="#Appraisal"/>
567   </SubClassOf>
568   <SubClassOf>
569   <Class IRI="#Authentication"/>
570   <ObjectSomeValuesFrom>
571   <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
572   <Class IRI="#authenticationResults"/>
573   </ObjectSomeValuesFrom>
574   </SubClassOf>
575   <SubClassOf>
```

```
576   <Class IRI="#Authentication"/>
577   <ObjectSomeValuesFrom>
578   <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
579   <Class IRI="#contactInformation"/>
580   </ObjectSomeValuesFrom>
581   </SubClassOf>
582   <SubClassOf>
583   <Class IRI="#Authentication"/>
584   <ObjectMinCardinality cardinality="1">
585   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
586   <Class IRI="#dataDepositer"/>
587   </ObjectMinCardinality>
588   </SubClassOf>
589   <SubClassOf>
590   <Class IRI="#CeaseDataCuration"/>
591   <Class IRI="#DataPreservationActivities"/>
592   </SubClassOf>
593   <SubClassOf>
594   <Class IRI="#CeaseDataCuration"/>
595   <ObjectSomeValuesFrom>
596   <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
597   <Class IRI="#preservationPolicy"/>
598   </ObjectSomeValuesFrom>
599   </SubClassOf>
600   <SubClassOf>
601   <Class IRI="#CeaseDataCuration"/>
602   <ObjectSomeValuesFrom>
603   <ObjectProperty IRI="http://www.w3.org/ns/prov#wasAssociatedWith"/>
604   <Class IRI="#dataCurator"/>
605   </ObjectSomeValuesFrom>
606   </SubClassOf>
607   <SubClassOf>
608   <Class IRI="#ChainOfCustody"/>
609   <Class IRI="#DataProcessing"/>
610   </SubClassOf>
611   <SubClassOf>
612   <Class IRI="#ChainOfCustody"/>
613   <ObjectMinCardinality cardinality="1">
614   <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
615   <Class IRI="#curationLog"/>
616   </ObjectMinCardinality>
617   </SubClassOf>
618   <SubClassOf>
619   <Class IRI="#ChainOfCustody"/>
620   <ObjectMinCardinality cardinality="1">
621   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
622   <Class IRI="#submittedData"/>
623   </ObjectMinCardinality>
624   </SubClassOf>
625   <SubClassOf>
626   <Class IRI="#CodeReview"/>
627   <Class IRI="#DataEvaluation"/>
628   </SubClassOf>
629   <SubClassOf>
630   <Class IRI="#CodeReview"/>
631   <ObjectSomeValuesFrom>
632   <ObjectProperty IRI="http://www.w3.org/ns/prov#wasAssociatedWith"/>
633   <Class IRI="#dataDepositer"/>
634   </ObjectSomeValuesFrom>
635   </SubClassOf>
636   <SubClassOf>
637   <Class IRI="#CodeReview"/>
638   <ObjectMinCardinality cardinality="1">
639   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
640   <Class IRI="#sourceCode"/>
641   </ObjectMinCardinality>
642   </SubClassOf>
643   <SubClassOf>
644   <Class IRI="#ConnectingDiscoveryServices"/>
645   <Class IRI="#DataPublishing"/>
646   </SubClassOf>
647   <SubClassOf>
648   <Class IRI="#ConnectingDiscoveryServices"/>
649   <ObjectSomeValuesFrom>
```

```
650   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
651   <Class IRI="#fullTextInformation"/>
652   </ObjectSomeValuesFrom>
653   </SubClassOf>
654   <SubClassOf>
655   <Class IRI="#ConnectingDiscoveryServices"/>
656   <ObjectSomeValuesFrom>
657   <ObjectProperty IRI="http://www.w3.org/ns/prov#wasAssociatedWith"/>
658   <Class IRI="#externalServiceProvider"/>
659   </ObjectSomeValuesFrom>
660   </SubClassOf>
661   <SubClassOf>
662   <Class IRI="#ConnectingDiscoveryServices"/>
663   <ObjectMinCardinality cardinality="1">
664   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
665   <Class IRI="#indexingInformation"/>
666   </ObjectMinCardinality>
667   </SubClassOf>
668   <SubClassOf>
669   <Class IRI="#Contextualization"/>
670   <Class IRI="#MetadataProcessing"/>
671   </SubClassOf>
672   <SubClassOf>
673   <Class IRI="#Contextualization"/>
674   <ObjectSomeValuesFrom>
675   <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
676   <Class IRI="#linkingInformation"/>
677   </ObjectSomeValuesFrom>
678   </SubClassOf>
679   <SubClassOf>
680   <Class IRI="#Contracts"/>
681   <Class IRI="#externalConstraints"/>
682   </SubClassOf>
683   <SubClassOf>
684   <Class IRI="#Conversion"/>
685   <Class IRI="#ActualDataProcessing"/>
686   </SubClassOf>
687   <SubClassOf>
688   <Class IRI="#Copyright"/>
689   <Class IRI="#Contracts"/>
690   </SubClassOf>
691   <SubClassOf>
692   <Class IRI="#Copyright"/>
693   <ObjectExactCardinality cardinality="1">
694   <ObjectProperty IRI="http://www.w3.org/ns/prov#qualifiedUsage"/>
695   <Class IRI="#termsOfUse"/>
696   </ObjectExactCardinality>
697   </SubClassOf>
698   <SubClassOf>
699   <Class IRI="#CreatingLandingPage"/>
700   <Class IRI="#DataPublishing"/>
701   </SubClassOf>
702   <SubClassOf>
703   <Class IRI="#CreatingLandingPage"/>
704   <ObjectSomeValuesFrom>
705   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
706   <Class IRI="#agreement"/>
707   </ObjectSomeValuesFrom>
708   </SubClassOf>
709   <SubClassOf>
710   <Class IRI="#CreatingLandingPage"/>
711   <ObjectSomeValuesFrom>
712   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
713   <Class IRI="#contactInformation"/>
714   </ObjectSomeValuesFrom>
715   </SubClassOf>
716   <SubClassOf>
717   <Class IRI="#CreatingLandingPage"/>
718   <ObjectSomeValuesFrom>
719   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
720   <Class IRI="#fileLocation"/>
721   </ObjectSomeValuesFrom>
722   </SubClassOf>
723   <SubClassOf>
```

```
724  <Class IRI="#CreatingLandingPage"/>
725  <ObjectSomeValuesFrom>
726  <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
727  <Class IRI="#indexingInformation"/>
728  </ObjectSomeValuesFrom>
729  </SubClassOf>
730  <SubClassOf>
731  <Class IRI="#CreatingLandingPage"/>
732  <ObjectSomeValuesFrom>
733  <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
734  <Class IRI="#persistentIdentifier"/>
735  </ObjectSomeValuesFrom>
736  </SubClassOf>
737  <SubClassOf>
738  <Class IRI="#CreatingLandingPage"/>
739  <ObjectSomeValuesFrom>
740  <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
741  <Class IRI="#selectionResults"/>
742  </ObjectSomeValuesFrom>
743  </SubClassOf>
744  <SubClassOf>
745  <Class IRI="#CreatingLandingPage"/>
746  <ObjectSomeValuesFrom>
747  <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
748  <Class IRI="#versionInformation"/>
749  </ObjectSomeValuesFrom>
750  </SubClassOf>
751  <SubClassOf>
752  <Class IRI="#CreatingLandingPage"/>
753  <ObjectSomeValuesFrom>
754  <ObjectProperty IRI="http://www.w3.org/ns/prov#wasAssociatedWith"/>
755  <Class IRI="#repositorySystem"/>
756  </ObjectSomeValuesFrom>
757  </SubClassOf>
758  <SubClassOf>
759  <Class IRI="#CreatingLandingPage"/>
760  <ObjectMinCardinality cardinality="1">
761  <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
762  <Class IRI="#landingPage"/>
763  </ObjectMinCardinality>
764  </SubClassOf>
765  <SubClassOf>
766  <Class IRI="#DataCleaning"/>
767  <Class IRI="#ActualDataProcessing"/>
768  </SubClassOf>
769  <SubClassOf>
770  <Class IRI="#DataCurationActivities"/>
771  <Class IRI="#DataLifecycleActivities"/>
772  </SubClassOf>
773  <SubClassOf>
774  <Class IRI="#DataCurationActivities"/>
775  <ObjectSomeValuesFrom>
776  <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
777  <Class IRI="#curationRecord"/>
778  </ObjectSomeValuesFrom>
779  </SubClassOf>
780  <SubClassOf>
781  <Class IRI="#DataEvaluation"/>
782  <Class IRI="#DataCurationActivities"/>
783  </SubClassOf>
784  <SubClassOf>
785  <Class IRI="#DataEvaluation"/>
786  <ObjectSomeValuesFrom>
787  <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
788  <Class IRI="#evaluationResults"/>
789  </ObjectSomeValuesFrom>
790  </SubClassOf>
791  <SubClassOf>
792  <Class IRI="#DataEvaluation"/>
793  <ObjectMinCardinality cardinality="1">
794  <ObjectProperty IRI="http://www.w3.org/ns/prov#wasAssociatedWith"/>
795  <ObjectUnionOf>
796  <Class IRI="#dataCurator"/>
797  <Class IRI="#peerReviewer"/>
```

```
798   </ObjectUnionOf>
799   </ObjectMinCardinality>
800   </SubClassOf>
801   <SubClassOf>
802   <Class IRI="#DataLifecycleActivities"/>
803   <Class IRI="http://www.w3.org/ns/prov#Activity"/>
804   </SubClassOf>
805   <SubClassOf>
806   <Class IRI="#DataPolicy"/>
807   <Class IRI="#externalConstraints"/>
808   </SubClassOf>
809   <SubClassOf>
810   <Class IRI="#DataPreservationActivities"/>
811   <Class IRI="#DataLifecycleActivities"/>
812   </SubClassOf>
813   <SubClassOf>
814   <Class IRI="#DataProcessing"/>
815   <Class IRI="#DataCurationActivities"/>
816   </SubClassOf>
817   <SubClassOf>
818   <Class IRI="#DataProcessing"/>
819   <ObjectMinCardinality cardinality="1">
820   <ObjectProperty IRI="http://www.w3.org/ns/prov#wasAssociatedWith"/>
821   <ObjectUnionOf>
822   <Class IRI="#dataCurator"/>
823   <Class IRI="#dataDepositer"/>
824   <Class IRI="#metadataCurator"/>
825   </ObjectUnionOf>
826   </ObjectMinCardinality>
827   </SubClassOf>
828   <SubClassOf>
829   <Class IRI="#DataPublishing"/>
830   <Class IRI="#DataCurationActivities"/>
831   </SubClassOf>
832   <SubClassOf>
833   <Class IRI="#DataPublishing"/>
834   <ObjectSomeValuesFrom>
835   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
836   <Class IRI="#evaluationResults"/>
837   </ObjectSomeValuesFrom>
838   </SubClassOf>
839   <SubClassOf>
840   <Class IRI="#DataPublishing"/>
841   <ObjectMinCardinality cardinality="1">
842   <ObjectProperty IRI="http://www.w3.org/ns/prov#wasAssociatedWith"/>
843   <ObjectUnionOf>
844   <Class IRI="#dataCurator"/>
845   <Class IRI="#repositorySystem"/>
846   <Class IRI="#systemAdministrator"/>
847   </ObjectUnionOf>
848   </ObjectMinCardinality>
849   </SubClassOf>
850   <SubClassOf>
851   <Class IRI="#DataVisualization"/>
852   <Class IRI="#ActualDataProcessing"/>
853   </SubClassOf>
854   <SubClassOf>
855   <Class IRI="#Deidentification"/>
856   <Class IRI="#ActualDataProcessing"/>
857   </SubClassOf>
858   <SubClassOf>
859   <Class IRI="#DepositAgreement"/>
860   <Class IRI="#Ingest"/>
861   </SubClassOf>
862   <SubClassOf>
863   <Class IRI="#DepositAgreement"/>
864   <ObjectSomeValuesFrom>
865   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
866   <Class IRI="#selectionPolicy"/>
867   </ObjectSomeValuesFrom>
868   </SubClassOf>
869   <SubClassOf>
870   <Class IRI="#DepositAgreement"/>
871   <ObjectExactCardinality cardinality="1">
```

```
872  <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
873  <Class IRI="#agreement"/>
874  </ObjectExactCardinality>
875  </SubClassOf>
876  <SubClassOf>
877  <Class IRI="#Digitization"/>
878  <Class IRI="#Conversion"/>
879  </SubClassOf>
880  <SubClassOf>
881  <Class IRI="#DiplomaticOrNationalSecurity"/>
882  <Class IRI="#externalConstraints"/>
883  </SubClassOf>
884  <SubClassOf>
885  <Class IRI="#DisciplinaryCustoms"/>
886  <Class IRI="#externalConstraints"/>
887  </SubClassOf>
888  <SubClassOf>
889  <Class IRI="#DisplayingDataCitation"/>
890  <Class IRI="#CreatingLandingPage"/>
891  </SubClassOf>
892  <SubClassOf>
893  <Class IRI="#Documentation"/>
894  <Class IRI="#Ingest"/>
895  </SubClassOf>
896  <SubClassOf>
897  <Class IRI="#Documentation"/>
898  <ObjectSomeValuesFrom>
899  <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
900  <Class IRI="#documentationPolicy"/>
901  </ObjectSomeValuesFrom>
902  </SubClassOf>
903  <SubClassOf>
904  <Class IRI="#Documentation"/>
905  <ObjectMinCardinality cardinality="1">
906  <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
907  <Class IRI="#dataDocument"/>
908  </ObjectMinCardinality>
909  </SubClassOf>
910  <SubClassOf>
911  <Class IRI="#Embargo"/>
912  <Class IRI="#AllowingFileDownload"/>
913  </SubClassOf>
914  <SubClassOf>
915  <Class IRI="#Emulation"/>
916  <Class IRI="#DataPreservationActivities"/>
917  </SubClassOf>
918  <SubClassOf>
919  <Class IRI="#Emulation"/>
920  <ObjectSomeValuesFrom>
921  <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
922  <Class IRI="#softwareRegistry"/>
923  </ObjectSomeValuesFrom>
924  </SubClassOf>
925  <SubClassOf>
926  <Class IRI="#Emulation"/>
927  <ObjectMinCardinality cardinality="1">
928  <ObjectProperty IRI="http://www.w3.org/ns/prov#wasAssociatedWith"/>
929  <Class IRI="#systemAdministrator"/>
930  </ObjectMinCardinality>
931  </SubClassOf>
932  <SubClassOf>
933  <Class IRI="#FileAuditing"/>
934  <Class IRI="#DataPreservationActivities"/>
935  </SubClassOf>
936  <SubClassOf>
937  <Class IRI="#FileAuditing"/>
938  <ObjectSomeValuesFrom>
939  <ObjectProperty IRI="http://www.w3.org/ns/prov#influenced"/>
940  <Class IRI="#fileValidationResults"/>
941  </ObjectSomeValuesFrom>
942  </SubClassOf>
943  <SubClassOf>
944  <Class IRI="#FileAuditing"/>
945  <ObjectMinCardinality cardinality="1">
```

```
946   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
947   <Class IRI="#curatedData"/>
948   </ObjectMinCardinality>
949   </SubClassOf>
950   <SubClassOf>
951   <Class IRI="#FileAuditing"/>
952   <ObjectMinCardinality cardinality="1">
953   <ObjectProperty IRI="http://www.w3.org/ns/prov#wasAssociatedWith"/>
954   <Class IRI="#systemAdministrator"/>
955   </ObjectMinCardinality>
956   </SubClassOf>
957   <SubClassOf>
958   <Class IRI="#FileAuditing"/>
959   <ObjectExactCardinality cardinality="1">
960   <ObjectProperty IRI="http://www.w3.org/ns/prov#wasInformedBy"/>
961   <Class IRI="#FormulateSuccessionPlanning"/>
962   </ObjectExactCardinality>
963   </SubClassOf>
964   <SubClassOf>
965   <Class IRI="#FileFormatTransformation"/>
966   <Class IRI="#ActualDataProcessing"/>
967   </SubClassOf>
968   <SubClassOf>
969   <Class IRI="#FileInventoryOrManifest"/>
970   <Class IRI="#DataProcessing"/>
971   </SubClassOf>
972   <SubClassOf>
973   <Class IRI="#FileInventoryOrManifest"/>
974   <ObjectSomeValuesFrom>
975   <ObjectProperty IRI="http://www.w3.org/ns/prov#influenced"/>
976   <Class IRI="#curationLog"/>
977   </ObjectSomeValuesFrom>
978   </SubClassOf>
979   <SubClassOf>
980   <Class IRI="#FileInventoryOrManifest"/>
981   <ObjectSomeValuesFrom>
982   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
983   <Class IRI="#researchData"/>
984   </ObjectSomeValuesFrom>
985   </SubClassOf>
986   <SubClassOf>
987   <Class IRI="#FileRenaming"/>
988   <Class IRI="#ActualDataProcessing"/>
989   </SubClassOf>
990   <SubClassOf>
991   <Class IRI="#FileValidation"/>
992   <Class IRI="#Appraisal"/>
993   </SubClassOf>
994   <SubClassOf>
995   <Class IRI="#FileValidation"/>
996   <ObjectSomeValuesFrom>
997   <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
998   <Class IRI="#fileValidationResults"/>
999   </ObjectSomeValuesFrom>
1000  </SubClassOf>
1001  <SubClassOf>
1002  <Class IRI="#FileValidation"/>
1003  <ObjectSomeValuesFrom>
1004  <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1005  <Class IRI="#sourceCode"/>
1006  </ObjectSomeValuesFrom>
1007  </SubClassOf>
1008  <SubClassOf>
1009  <Class IRI="#FileValidation"/>
1010  <ObjectMinCardinality cardinality="1">
1011  <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1012  <Class IRI="#submittedData"/>
1013  </ObjectMinCardinality>
1014  </SubClassOf>
1015  <SubClassOf>
1016  <Class IRI="#FormulateSuccessionPlanning"/>
1017  <Class IRI="#DataPreservationActivities"/>
1018  </SubClassOf>
1019  <SubClassOf>
```

```
1020   <Class IRI="#FormulateSuccessionPlanning"/>
1021   <ObjectSomeValuesFrom>
1022   <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
1023   <Class IRI="#successionPlan"/>
1024   </ObjectSomeValuesFrom>
1025   </SubClassOf>
1026   <SubClassOf>
1027   <Class IRI="#FormulateSuccessionPlanning"/>
1028   <ObjectSomeValuesFrom>
1029   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1030   <Class IRI="#preservationPolicy"/>
1031   </ObjectSomeValuesFrom>
1032   </SubClassOf>
1033   <SubClassOf>
1034   <Class IRI="#FormulateSuccessionPlanning"/>
1035   <ObjectMinCardinality cardinality="1">
1036   <ObjectProperty IRI="http://www.w3.org/ns/prov#wasAssociatedWith"/>
1037   <ObjectUnionOf>
1038   <Class IRI="#administrator"/>
1039   <Class IRI="#dataCurator"/>
1040   </ObjectUnionOf>
1041   </ObjectMinCardinality>
1042   </SubClassOf>
1043   <SubClassOf>
1044   <Class IRI="#GeneratingFulltextIndexing"/>
1045   <Class IRI="#DataPublishing"/>
1046   </SubClassOf>
1047   <SubClassOf>
1048   <Class IRI="#GeneratingFulltextIndexing"/>
1049   <ObjectSomeValuesFrom>
1050   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1051   <Class IRI="#curatedData"/>
1052   </ObjectSomeValuesFrom>
1053   </SubClassOf>
1054   <SubClassOf>
1055   <Class IRI="#GeneratingFulltextIndexing"/>
1056   <ObjectSomeValuesFrom>
1057   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1058   <Class IRI="#dataDocument"/>
1059   </ObjectSomeValuesFrom>
1060   </SubClassOf>
1061   <SubClassOf>
1062   <Class IRI="#GeneratingFulltextIndexing"/>
1063   <ObjectSomeValuesFrom>
1064   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1065   <Class IRI="#submittedData"/>
1066   </ObjectSomeValuesFrom>
1067   </SubClassOf>
1068   <SubClassOf>
1069   <Class IRI="#GeneratingFulltextIndexing"/>
1070   <ObjectMinCardinality cardinality="1">
1071   <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
1072   <Class IRI="#fullTextInformation"/>
1073   </ObjectMinCardinality>
1074   </SubClassOf>
1075   <SubClassOf>
1076   <Class IRI="#IPR"/>
1077   <Class IRI="#Contracts"/>
1078   </SubClassOf>
1079   <SubClassOf>
1080   <Class IRI="#Indexing"/>
1081   <Class IRI="#DataPublishing"/>
1082   </SubClassOf>
1083   <SubClassOf>
1084   <Class IRI="#Indexing"/>
1085   <ObjectSomeValuesFrom>
1086   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1087   <Class IRI="#dataDocument"/>
1088   </ObjectSomeValuesFrom>
1089   </SubClassOf>
1090   <SubClassOf>
1091   <Class IRI="#Indexing"/>
1092   <ObjectSomeValuesFrom>
1093   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
```

```
1094    <Class IRI="#linkingInformation"/>
1095    </ObjectSomeValuesFrom>
1096    </SubClassOf>
1097    <SubClassOf>
1098    <Class IRI="#Indexing"/>
1099    <ObjectMinCardinality cardinality="1">
1100    <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
1101    <Class IRI="#indexingInformation"/>
1102    </ObjectMinCardinality>
1103    </SubClassOf>
1104    <SubClassOf>
1105    <Class IRI="#Indexing"/>
1106    <ObjectMinCardinality cardinality="1">
1107    <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1108    <Class IRI="#retrievalMetadata"/>
1109    </ObjectMinCardinality>
1110    </SubClassOf>
1111    <SubClassOf>
1112    <Class IRI="#InformedConsent"/>
1113    <Class IRI="#Contracts"/>
1114    </SubClassOf>
1115    <SubClassOf>
1116    <Class IRI="#Ingest"/>
1117    <Class IRI="#DataCurationActivities"/>
1118    </SubClassOf>
1119    <SubClassOf>
1120    <Class IRI="#Ingest"/>
1121    <ObjectMinCardinality cardinality="1">
1122    <ObjectProperty IRI="http://www.w3.org/ns/prov#wasAssociatedWith"/>
1123    <Class IRI="#dataDepositer"/>
1124    </ObjectMinCardinality>
1125    </SubClassOf>
1126    <SubClassOf>
1127    <Class IRI="#Interoperability"/>
1128    <Class IRI="#ActualDataProcessing"/>
1129    </SubClassOf>
1130    <SubClassOf>
1131    <Class IRI="#MaintainingContactInformation"/>
1132    <Class IRI="#DataPreservationActivities"/>
1133    </SubClassOf>
1134    <SubClassOf>
1135    <Class IRI="#MaintainingContactInformation"/>
1136    <ObjectSomeValuesFrom>
1137    <ObjectProperty IRI="http://www.w3.org/ns/prov#influenced"/>
1138    <Class IRI="#contactInformation"/>
1139    </ObjectSomeValuesFrom>
1140    </SubClassOf>
1141    <SubClassOf>
1142    <Class IRI="#MaintainingContactInformation"/>
1143    <ObjectSomeValuesFrom>
1144    <ObjectProperty IRI="http://www.w3.org/ns/prov#wasAssociatedWith"/>
1145    <Class IRI="#dataCurator"/>
1146    </ObjectSomeValuesFrom>
1147    </SubClassOf>
1148    <SubClassOf>
1149    <Class IRI="#MaintainingContactInformation"/>
1150    <ObjectSomeValuesFrom>
1151    <ObjectProperty IRI="http://www.w3.org/ns/prov#wasAttributedTo"/>
1152    <Class IRI="#dataDepositer"/>
1153    </ObjectSomeValuesFrom>
1154    </SubClassOf>
1155    <SubClassOf>
1156    <Class IRI="#MaintainingContactInformation"/>
1157    <ObjectMinCardinality cardinality="1">
1158    <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1159    <Class IRI="#contactInformation"/>
1160    </ObjectMinCardinality>
1161    </SubClassOf>
1162    <SubClassOf>
1163    <Class IRI="#MetadataGeneration"/>
1164    <Class IRI="#MetadataProcessing"/>
1165    </SubClassOf>
1166    <SubClassOf>
1167    <Class IRI="#MetadataGeneration"/>
```

```
1168   <ObjectExactCardinality cardinality="1">
1169   <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
1170   <Class IRI="#retrievalMetadata"/>
1171   </ObjectExactCardinality>
1172   </SubClassOf>
1173   <SubClassOf>
1174   <Class IRI="#MetadataProcessing"/>
1175   <Class IRI="#DataProcessing"/>
1176   </SubClassOf>
1177   <SubClassOf>
1178   <Class IRI="#MetadataProcessing"/>
1179   <ObjectSomeValuesFrom>
1180   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1181   <Class IRI="#dataDocument"/>
1182   </ObjectSomeValuesFrom>
1183   </SubClassOf>
1184   <SubClassOf>
1185   <Class IRI="#MetadataProcessing"/>
1186   <ObjectSomeValuesFrom>
1187   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1188   <Class IRI="#metadataSchema"/>
1189   </ObjectSomeValuesFrom>
1190   </SubClassOf>
1191   <SubClassOf>
1192   <Class IRI="#MetadataProcessing"/>
1193   <ObjectSomeValuesFrom>
1194   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1195   <Class IRI="#researchData"/>
1196   </ObjectSomeValuesFrom>
1197   </SubClassOf>
1198   <SubClassOf>
1199   <Class IRI="#MetadataProcessing"/>
1200   <ObjectSomeValuesFrom>
1201   <ObjectProperty IRI="http://www.w3.org/ns/prov#wasAssociatedWith"/>
1202   <ObjectUnionOf>
1203   <Class IRI="#dataDepositer"/>
1204   <Class IRI="#metadataCurator"/>
1205   </ObjectUnionOf>
1206   </ObjectSomeValuesFrom>
1207   </SubClassOf>
1208   <SubClassOf>
1209   <Class IRI="#Migration"/>
1210   <Class IRI="#DataPreservationActivities"/>
1211   </SubClassOf>
1212   <SubClassOf>
1213   <Class IRI="#Migration"/>
1214   <ObjectSomeValuesFrom>
1215   <ObjectProperty IRI="http://www.w3.org/ns/prov#influenced"/>
1216   <Class IRI="#curatedData"/>
1217   </ObjectSomeValuesFrom>
1218   </SubClassOf>
1219   <SubClassOf>
1220   <Class IRI="#Migration"/>
1221   <ObjectMinCardinality cardinality="1">
1222   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1223   <Class IRI="#curatedData"/>
1224   </ObjectMinCardinality>
1225   </SubClassOf>
1226   <SubClassOf>
1227   <Class IRI="#Migration"/>
1228   <ObjectMinCardinality cardinality="1">
1229   <ObjectProperty IRI="http://www.w3.org/ns/prov#wasAssociatedWith"/>
1230   <ObjectUnionOf>
1231   <Class IRI="#dataCurator"/>
1232   <Class IRI="#repositorySystem"/>
1233   </ObjectUnionOf>
1234   </ObjectMinCardinality>
1235   </SubClassOf>
1236   <SubClassOf>
1237   <Class IRI="#MintingPersistentIdentifier"/>
1238   <Class IRI="#DataPublishing"/>
1239   </SubClassOf>
1240   <SubClassOf>
1241   <Class IRI="#MintingPersistentIdentifier"/>
```

```
1242    <ObjectSomeValuesFrom>
1243    <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
1244    <Class IRI="#persistentIdentifier"/>
1245    </ObjectSomeValuesFrom>
1246    </SubClassOf>
1247    <SubClassOf>
1248    <Class IRI="#MintingPersistentIdentifier"/>
1249    <ObjectSomeValuesFrom>
1250    <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1251    <Class IRI="#indexingInformation"/>
1252    </ObjectSomeValuesFrom>
1253    </SubClassOf>
1254    <SubClassOf>
1255    <Class IRI="#MintingPersistentIdentifier"/>
1256    <ObjectSomeValuesFrom>
1257    <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1258    <Class IRI="#researchData"/>
1259    </ObjectSomeValuesFrom>
1260    </SubClassOf>
1261    <SubClassOf>
1262    <Class IRI="#PeerReview"/>
1263    <Class IRI="#DataEvaluation"/>
1264    </SubClassOf>
1265    <SubClassOf>
1266    <Class IRI="#PeerReview"/>
1267    <ObjectSomeValuesFrom>
1268    <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1269    <Class IRI="#curationLog"/>
1270    </ObjectSomeValuesFrom>
1271    </SubClassOf>
1272    <SubClassOf>
1273    <Class IRI="#PeerReview"/>
1274    <ObjectSomeValuesFrom>
1275    <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1276    <Class IRI="#metadata"/>
1277    </ObjectSomeValuesFrom>
1278    </SubClassOf>
1279    <SubClassOf>
1280    <Class IRI="#PeerReview"/>
1281    <ObjectSomeValuesFrom>
1282    <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1283    <Class IRI="#sourceCode"/>
1284    </ObjectSomeValuesFrom>
1285    </SubClassOf>
1286    <SubClassOf>
1287    <Class IRI="#PeerReview"/>
1288    <ObjectMinCardinality cardinality="1">
1289    <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1290    <Class IRI="#curatedData"/>
1291    </ObjectMinCardinality>
1292    </SubClassOf>
1293    <SubClassOf>
1294    <Class IRI="#PeerReview"/>
1295    <ObjectMinCardinality cardinality="1">
1296    <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1297    <Class IRI="#dataDocument"/>
1298    </ObjectMinCardinality>
1299    </SubClassOf>
1300    <SubClassOf>
1301    <Class IRI="#PeerReview"/>
1302    <ObjectMinCardinality cardinality="1">
1303    <ObjectProperty IRI="http://www.w3.org/ns/prov#wasAssociatedWith"/>
1304    <Class IRI="#peerReviewer"/>
1305    </ObjectMinCardinality>
1306    </SubClassOf>
1307    <SubClassOf>
1308    <Class IRI="#PersonalInformation"/>
1309    <Class IRI="#externalConstraints"/>
1310    </SubClassOf>
1311    <SubClassOf>
1312    <Class IRI="#ProvidingRestrictedAccess"/>
1313    <Class IRI="#AllowingFileDownload"/>
1314    </SubClassOf>
1315    <SubClassOf>
```

```
1316   <Class IRI="#QualityAssurance"/>
1317   <Class IRI="#DataEvaluation"/>
1318   </SubClassOf>
1319   <SubClassOf>
1320   <Class IRI="#QualityAssurance"/>
1321   <ObjectSomeValuesFrom>
1322   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1323   <Class IRI="#curatedData"/>
1324   </ObjectSomeValuesFrom>
1325   </SubClassOf>
1326   <SubClassOf>
1327   <Class IRI="#QualityAssurance"/>
1328   <ObjectSomeValuesFrom>
1329   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1330   <Class IRI="#dataDocument"/>
1331   </ObjectSomeValuesFrom>
1332   </SubClassOf>
1333   <SubClassOf>
1334   <Class IRI="#QualityAssurance"/>
1335   <ObjectSomeValuesFrom>
1336   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1337   <Class IRI="#metadata"/>
1338   </ObjectSomeValuesFrom>
1339   </SubClassOf>
1340   <SubClassOf>
1341   <Class IRI="#QualityAssurance"/>
1342   <ObjectSomeValuesFrom>
1343   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1344   <Class IRI="#sourceCode"/>
1345   </ObjectSomeValuesFrom>
1346   </SubClassOf>
1347   <SubClassOf>
1348   <Class IRI="#QualityAssurance"/>
1349   <ObjectAllValuesFrom>
1350   <ObjectProperty IRI="http://www.w3.org/ns/prov#wasAssociatedWith"/>
1351   <Class IRI="#dataCurator"/>
1352   </ObjectAllValuesFrom>
1353   </SubClassOf>
1354   <SubClassOf>
1355   <Class IRI="#RegisteringSoftware"/>
1356   <Class IRI="#DataPreservationActivities"/>
1357   </SubClassOf>
1358   <SubClassOf>
1359   <Class IRI="#RegisteringSoftware"/>
1360   <ObjectSomeValuesFrom>
1361   <ObjectProperty IRI="http://www.w3.org/ns/prov#influenced"/>
1362   <Class IRI="#softwareRegistry"/>
1363   </ObjectSomeValuesFrom>
1364   </SubClassOf>
1365   <SubClassOf>
1366   <Class IRI="#RegisteringSoftware"/>
1367   <ObjectSomeValuesFrom>
1368   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1369   <Class IRI="#dataDocument"/>
1370   </ObjectSomeValuesFrom>
1371   </SubClassOf>
1372   <SubClassOf>
1373   <Class IRI="#RegisteringSoftware"/>
1374   <ObjectSomeValuesFrom>
1375   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1376   <Class IRI="#retrievalMetadata"/>
1377   </ObjectSomeValuesFrom>
1378   </SubClassOf>
1379   <SubClassOf>
1380   <Class IRI="#RegisteringSoftware"/>
1381   <ObjectSomeValuesFrom>
1382   <ObjectProperty IRI="http://www.w3.org/ns/prov#wasAssociatedWith"/>
1383   <Class IRI="#dataCurator"/>
1384   </ObjectSomeValuesFrom>
1385   </SubClassOf>
1386   <SubClassOf>
1387   <Class IRI="#Restructure"/>
1388   <Class IRI="#ActualDataProcessing"/>
1389   </SubClassOf>
```

```
1390   <SubClassOf>
1391   <Class IRI="#RightsManagement"/>
1392   <Class IRI="#Appraisal"/>
1393   </SubClassOf>
1394   <SubClassOf>
1395   <Class IRI="#RightsManagement"/>
1396   <ObjectSomeValuesFrom>
1397   <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
1398   <Class IRI="#termsOfUse"/>
1399   </ObjectSomeValuesFrom>
1400   </SubClassOf>
1401   <SubClassOf>
1402   <Class IRI="#RightsManagement"/>
1403   <ObjectSomeValuesFrom>
1404   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1405   <Class IRI="#dataDocument"/>
1406   </ObjectSomeValuesFrom>
1407   </SubClassOf>
1408   <SubClassOf>
1409   <Class IRI="#RiskManagement"/>
1410   <Class IRI="#Appraisal"/>
1411   </SubClassOf>
1412   <SubClassOf>
1413   <Class IRI="#RiskManagement"/>
1414   <ObjectSomeValuesFrom>
1415   <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
1416   <Class IRI="#accessRestriction"/>
1417   </ObjectSomeValuesFrom>
1418   </SubClassOf>
1419   <SubClassOf>
1420   <Class IRI="#RiskManagement"/>
1421   <ObjectSomeValuesFrom>
1422   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1423   <Class IRI="#externalConstraints"/>
1424   </ObjectSomeValuesFrom>
1425   </SubClassOf>
1426   <SubClassOf>
1427   <Class IRI="#RiskManagement"/>
1428   <ObjectMinCardinality cardinality="1">
1429   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1430   <ObjectUnionOf>
1431   <Class IRI="#dataDocument"/>
1432   <Class IRI="#submittedData"/>
1433   </ObjectUnionOf>
1434   </ObjectMinCardinality>
1435   </SubClassOf>
1436   <SubClassOf>
1437   <Class IRI="#SecuringStorage"/>
1438   <Class IRI="#DataPreservationActivities"/>
1439   </SubClassOf>
1440   <SubClassOf>
1441   <Class IRI="#SecuringStorage"/>
1442   <ObjectSomeValuesFrom>
1443   <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
1444   <Class IRI="#backupData"/>
1445   </ObjectSomeValuesFrom>
1446   </SubClassOf>
1447   <SubClassOf>
1448   <Class IRI="#SecuringStorage"/>
1449   <ObjectSomeValuesFrom>
1450   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1451   <Class IRI="#metadata"/>
1452   </ObjectSomeValuesFrom>
1453   </SubClassOf>
1454   <SubClassOf>
1455   <Class IRI="#SecuringStorage"/>
1456   <ObjectSomeValuesFrom>
1457   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1458   <Class IRI="#researchData"/>
1459   </ObjectSomeValuesFrom>
1460   </SubClassOf>
1461   <SubClassOf>
1462   <Class IRI="#SecuringStorage"/>
1463   <ObjectMinCardinality cardinality="1">
```

```
1464    <ObjectProperty IRI="http://www.w3.org/ns/prov#wasAssociatedWith"/>
1465    <Class IRI="#systemAdministrator"/>
1466    </ObjectMinCardinality>
1467    </SubClassOf>
1468    <SubClassOf>
1469    <Class IRI="#Selection"/>
1470    <Class IRI="#Appraisal"/>
1471    </SubClassOf>
1472    <SubClassOf>
1473    <Class IRI="#Selection"/>
1474    <ObjectSomeValuesFrom>
1475    <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1476    <Class IRI="#accessRestriction"/>
1477    </ObjectSomeValuesFrom>
1478    </SubClassOf>
1479    <SubClassOf>
1480    <Class IRI="#Selection"/>
1481    <ObjectSomeValuesFrom>
1482    <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1483    <Class IRI="#authenticationResults"/>
1484    </ObjectSomeValuesFrom>
1485    </SubClassOf>
1486    <SubClassOf>
1487    <Class IRI="#Selection"/>
1488    <ObjectSomeValuesFrom>
1489    <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1490    <Class IRI="#dataDocument"/>
1491    </ObjectSomeValuesFrom>
1492    </SubClassOf>
1493    <SubClassOf>
1494    <Class IRI="#Selection"/>
1495    <ObjectSomeValuesFrom>
1496    <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1497    <Class IRI="#fileValidationResults"/>
1498    </ObjectSomeValuesFrom>
1499    </SubClassOf>
1500    <SubClassOf>
1501    <Class IRI="#Selection"/>
1502    <ObjectSomeValuesFrom>
1503    <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1504    <Class IRI="#selectionPolicy"/>
1505    </ObjectSomeValuesFrom>
1506    </SubClassOf>
1507    <SubClassOf>
1508    <Class IRI="#Selection"/>
1509    <ObjectSomeValuesFrom>
1510    <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1511    <Class IRI="#termsOfUse"/>
1512    </ObjectSomeValuesFrom>
1513    </SubClassOf>
1514    <SubClassOf>
1515    <Class IRI="#Selection"/>
1516    <ObjectAllValuesFrom>
1517    <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
1518    <Class IRI="#selectionResults"/>
1519    </ObjectAllValuesFrom>
1520    </SubClassOf>
1521    <SubClassOf>
1522    <Class IRI="#SensitiveData"/>
1523    <Class IRI="#PersonalInformation"/>
1524    </SubClassOf>
1525    <SubClassOf>
1526    <Class IRI="#SettingTermsOfUse"/>
1527    <Class IRI="#CreatingLandingPage"/>
1528    </SubClassOf>
1529    <SubClassOf>
1530    <Class IRI="#SettingTermsOfUse"/>
1531    <ObjectSomeValuesFrom>
1532    <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1533    <Class IRI="#termsOfUse"/>
1534    </ObjectSomeValuesFrom>
1535    </SubClassOf>
1536    <SubClassOf>
1537    <Class IRI="#SubmitData"/>
```

```
1538   <Class IRI="#Ingest"/>
1539   </SubClassOf>
1540   <SubClassOf>
1541   <Class IRI="#SubmitData"/>
1542   <ObjectSomeValuesFrom>
1543   <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
1544   <Class IRI="#sourceCode"/>
1545   </ObjectSomeValuesFrom>
1546   </SubClassOf>
1547   <SubClassOf>
1548   <Class IRI="#SubmitData"/>
1549   <ObjectMinCardinality cardinality="1">
1550   <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
1551   <Class IRI="#submittedData"/>
1552   </ObjectMinCardinality>
1553   </SubClassOf>
1554   <SubClassOf>
1555   <Class IRI="#TechnologyMonitoringAndRefreshing"/>
1556   <Class IRI="#DataPreservationActivities"/>
1557   </SubClassOf>
1558   <SubClassOf>
1559   <Class IRI="#TechnologyMonitoringAndRefreshing"/>
1560   <ObjectSomeValuesFrom>
1561   <ObjectProperty IRI="http://www.w3.org/ns/prov#influenced"/>
1562   <Class IRI="#technicalInformation"/>
1563   </ObjectSomeValuesFrom>
1564   </SubClassOf>
1565   <SubClassOf>
1566   <Class IRI="#TechnologyMonitoringAndRefreshing"/>
1567   <ObjectSomeValuesFrom>
1568   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1569   <Class IRI="#technicalInformation"/>
1570   </ObjectSomeValuesFrom>
1571   </SubClassOf>
1572   <SubClassOf>
1573   <Class IRI="#TechnologyMonitoringAndRefreshing"/>
1574   <ObjectMinCardinality cardinality="1">
1575   <ObjectProperty IRI="http://www.w3.org/ns/prov#wasAssociatedWith"/>
1576   <Class IRI="#systemAdministrator"/>
1577   </ObjectMinCardinality>
1578   </SubClassOf>
1579   <SubClassOf>
1580   <Class IRI="#TrackingUseAnalytics"/>
1581   <Class IRI="#DataPreservationActivities"/>
1582   </SubClassOf>
1583   <SubClassOf>
1584   <Class IRI="#TrackingUseAnalytics"/>
1585   <ObjectSomeValuesFrom>
1586   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1587   <Class IRI="#dataDocument"/>
1588   </ObjectSomeValuesFrom>
1589   </SubClassOf>
1590   <SubClassOf>
1591   <Class IRI="#TrackingUseAnalytics"/>
1592   <ObjectSomeValuesFrom>
1593   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1594   <Class IRI="#landingPage"/>
1595   </ObjectSomeValuesFrom>
1596   </SubClassOf>
1597   <SubClassOf>
1598   <Class IRI="#TrackingUseAnalytics"/>
1599   <ObjectSomeValuesFrom>
1600   <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1601   <Class IRI="#researchData"/>
1602   </ObjectSomeValuesFrom>
1603   </SubClassOf>
1604   <SubClassOf>
1605   <Class IRI="#TrackingUseAnalytics"/>
1606   <ObjectMinCardinality cardinality="1">
1607   <ObjectProperty IRI="http://www.w3.org/ns/prov#generated"/>
1608   <Class IRI="#usageResults"/>
1609   </ObjectMinCardinality>
1610   </SubClassOf>
1611   <SubClassOf>
```

```
1612  <Class IRI="#Transcoding"/>
1613  <Class IRI="#FileFormatTransformation"/>
1614  </SubClassOf>
1615  <SubClassOf>
1616  <Class IRI="#TranscribingAndTranslatingData"/>
1617  <Class IRI="#Conversion"/>
1618  </SubClassOf>
1619  <SubClassOf>
1620  <Class IRI="#Versioning"/>
1621  <Class IRI="#DataPreservationActivities"/>
1622  </SubClassOf>
1623  <SubClassOf>
1624  <Class IRI="#Versioning"/>
1625  <ObjectSomeValuesFrom>
1626  <ObjectProperty IRI="http://www.w3.org/ns/prov#influenced"/>
1627  <Class IRI="#versionInformation"/>
1628  </ObjectSomeValuesFrom>
1629  </SubClassOf>
1630  <SubClassOf>
1631  <Class IRI="#Versioning"/>
1632  <ObjectSomeValuesFrom>
1633  <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1634  <Class IRI="#curatedData"/>
1635  </ObjectSomeValuesFrom>
1636  </SubClassOf>
1637  <SubClassOf>
1638  <Class IRI="#Versioning"/>
1639  <ObjectSomeValuesFrom>
1640  <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1641  <Class IRI="#submittedData"/>
1642  </ObjectSomeValuesFrom>
1643  </SubClassOf>
1644  <SubClassOf>
1645  <Class IRI="#Versioning"/>
1646  <ObjectExactCardinality cardinality="1">
1647  <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1648  <Class IRI="#versionInformation"/>
1649  </ObjectExactCardinality>
1650  </SubClassOf>
1651  <SubClassOf>
1652  <Class IRI="#accessRestriction"/>
1653  <Class IRI="#fileLocation"/>
1654  </SubClassOf>
1655  <SubClassOf>
1656  <Class IRI="#administrator"/>
1657  <Class IRI="http://www.w3.org/ns/prov#Person"/>
1658  </SubClassOf>
1659  <SubClassOf>
1660  <Class IRI="#agreement"/>
1661  <Class IRI="#curationRecord"/>
1662  </SubClassOf>
1663  <SubClassOf>
1664  <Class IRI="#appropriateRepository"/>
1665  <Class IRI="#selectionPolicy"/>
1666  </SubClassOf>
1667  <SubClassOf>
1668  <Class IRI="#authenticationResults"/>
1669  <Class IRI="#curationRecord"/>
1670  </SubClassOf>
1671  <SubClassOf>
1672  <Class IRI="#backupData"/>
1673  <Class IRI="http://www.w3.org/ns/prov#Entity"/>
1674  </SubClassOf>
1675  <SubClassOf>
1676  <Class IRI="#contactInformation"/>
1677  <Class IRI="#metadata"/>
1678  </SubClassOf>
1679  <SubClassOf>
1680  <Class IRI="#costAndLeadTime"/>
1681  <Class IRI="#selectionPolicy"/>
1682  </SubClassOf>
1683  <SubClassOf>
1684  <Class IRI="#creditOnTheResults"/>
1685  <Class IRI="#termsOfUse"/>
```

```
1686    </SubClassOf>
1687    <SubClassOf>
1688    <Class IRI="#curatedData"/>
1689    <Class IRI="#researchData"/>
1690    </SubClassOf>
1691    <SubClassOf>
1692    <Class IRI="#curatedData"/>
1693    <ObjectMinCardinality cardinality="1">
1694    <ObjectProperty IRI="http://www.w3.org/ns/prov#wasDerivedFrom"/>
1695    <Class IRI="#submittedData"/>
1696    </ObjectMinCardinality>
1697    </SubClassOf>
1698    <SubClassOf>
1699    <Class IRI="#curationLog"/>
1700    <Class IRI="#curationRecord"/>
1701    </SubClassOf>
1702    <SubClassOf>
1703    <Class IRI="#curationRecord"/>
1704    <Class IRI="http://www.w3.org/ns/prov#Entity"/>
1705    </SubClassOf>
1706    <SubClassOf>
1707    <Class IRI="#dataCurator"/>
1708    <Class IRI="http://www.w3.org/ns/prov#Person"/>
1709    </SubClassOf>
1710    <SubClassOf>
1711    <Class IRI="#dataDepositer"/>
1712    <Class IRI="http://www.w3.org/ns/prov#Person"/>
1713    </SubClassOf>
1714    <SubClassOf>
1715    <Class IRI="#dataDocument"/>
1716    <Class IRI="http://www.w3.org/ns/prov#Entity"/>
1717    </SubClassOf>
1718    <SubClassOf>
1719    <Class IRI="#dataProcessingPolicy"/>
1720    <Class IRI="#policy"/>
1721    </SubClassOf>
1722    <SubClassOf>
1723    <Class IRI="#dataUser"/>
1724    <Class IRI="http://www.w3.org/ns/prov#Person"/>
1725    </SubClassOf>
1726    <SubClassOf>
1727    <Class IRI="#documentationPolicy"/>
1728    <Class IRI="#policy"/>
1729    </SubClassOf>
1730    <SubClassOf>
1731    <Class IRI="#evaluationResults"/>
1732    <Class IRI="#curationRecord"/>
1733    </SubClassOf>
1734    <SubClassOf>
1735    <Class IRI="#expertiseNecessity"/>
1736    <Class IRI="#selectionPolicy"/>
1737    </SubClassOf>
1738    <SubClassOf>
1739    <Class IRI="#externalConstraints"/>
1740    <Class IRI="http://www.w3.org/ns/prov#Entity"/>
1741    </SubClassOf>
1742    <SubClassOf>
1743    <Class IRI="#externalServiceProvider"/>
1744    <Class IRI="http://www.w3.org/ns/prov#SoftwareAgent"/>
1745    </SubClassOf>
1746    <SubClassOf>
1747    <Class IRI="#feasibility"/>
1748    <Class IRI="#selectionPolicy"/>
1749    </SubClassOf>
1750    <SubClassOf>
1751    <Class IRI="#feeForUse"/>
1752    <Class IRI="#termsOfUse"/>
1753    </SubClassOf>
1754    <SubClassOf>
1755    <Class IRI="#fileLocation"/>
1756    <Class IRI="#metadata"/>
1757    </SubClassOf>
1758    <SubClassOf>
1759    <Class IRI="#fileValidationResults"/>
```

```
1760    <Class IRI="#curationRecord"/>
1761    </SubClassOf>
1762    <SubClassOf>
1763    <Class IRI="#fullTextInformation"/>
1764    <Class IRI="#metadata"/>
1765    </SubClassOf>
1766    <SubClassOf>
1767    <Class IRI="#imposeTheSameConditions"/>
1768    <Class IRI="#termsOfUse"/>
1769    </SubClassOf>
1770    <SubClassOf>
1771    <Class IRI="#improperUse"/>
1772    <Class IRI="#termsOfUse"/>
1773    </SubClassOf>
1774    <SubClassOf>
1775    <Class IRI="#indexingInformation"/>
1776    <Class IRI="#metadata"/>
1777    </SubClassOf>
1778    <SubClassOf>
1779    <Class IRI="#landingPage"/>
1780    <Class IRI="http://www.w3.org/ns/prov#Entity"/>
1781    </SubClassOf>
1782    <SubClassOf>
1783    <Class IRI="#landingPage"/>
1784    <ObjectSomeValuesFrom>
1785    <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1786    <Class IRI="#visualizedData"/>
1787    </ObjectSomeValuesFrom>
1788    </SubClassOf>
1789    <SubClassOf>
1790    <Class IRI="#landingPage"/>
1791    <ObjectAllValuesFrom>
1792    <ObjectProperty IRI="http://www.w3.org/ns/prov#wasAttributedTo"/>
1793    <Class IRI="#repositorySystem"/>
1794    </ObjectAllValuesFrom>
1795    </SubClassOf>
1796    <SubClassOf>
1797    <Class IRI="#landingPage"/>
1798    <ObjectAllValuesFrom>
1799    <ObjectProperty IRI="http://www.w3.org/ns/prov#wasStartedBy"/>
1800    <Class IRI="#DataPublishing"/>
1801    </ObjectAllValuesFrom>
1802    </SubClassOf>
1803    <SubClassOf>
1804    <Class IRI="#landingPage"/>
1805    <ObjectMinCardinality cardinality="1">
1806    <ObjectProperty IRI="http://www.w3.org/ns/prov#used"/>
1807    <Class IRI="#retrievalMetadata"/>
1808    </ObjectMinCardinality>
1809    </SubClassOf>
1810    <SubClassOf>
1811    <Class IRI="#landingPage"/>
1812    <ObjectExactCardinality cardinality="1">
1813    <ObjectProperty IRI="http://www.w3.org/ns/prov#wasEndedBy"/>
1814    <Class IRI="#CeaseDataCuration"/>
1815    </ObjectExactCardinality>
1816    </SubClassOf>
1817    <SubClassOf>
1818    <Class IRI="#linkingInformation"/>
1819    <Class IRI="#metadata"/>
1820    </SubClassOf>
1821    <SubClassOf>
1822    <Class IRI="#metadata"/>
1823    <Class IRI="http://www.w3.org/ns/prov#Entity"/>
1824    </SubClassOf>
1825    <SubClassOf>
1826    <Class IRI="#metadata"/>
1827    <ObjectSomeValuesFrom>
1828    <ObjectProperty IRI="http://xmlns.com/foaf/0.1/primaryTopic"/>
1829    <Class IRI="#researchData"/>
1830    </ObjectSomeValuesFrom>
1831    </SubClassOf>
1832    <SubClassOf>
1833    <Class IRI="#metadataCurator"/>
```

```
1834    <Class IRI="http://www.w3.org/ns/prov#Person"/>
1835    </SubClassOf>
1836    <SubClassOf>
1837    <Class IRI="#metadataSchema"/>
1838    <Class IRI="http://www.w3.org/ns/prov#Entity"/>
1839    </SubClassOf>
1840    <SubClassOf>
1841    <Class IRI="#noDelivs"/>
1842    <Class IRI="#termsOfUse"/>
1843    </SubClassOf>
1844    <SubClassOf>
1845    <Class IRI="#nonCommercial"/>
1846    <Class IRI="#termsOfUse"/>
1847    </SubClassOf>
1848    <SubClassOf>
1849    <Class IRI="#peerReviewer"/>
1850    <Class IRI="http://www.w3.org/ns/prov#Person"/>
1851    </SubClassOf>
1852    <SubClassOf>
1853    <Class IRI="#persistentIdentifier"/>
1854    <Class IRI="#metadata"/>
1855    </SubClassOf>
1856    <SubClassOf>
1857    <Class IRI="#policy"/>
1858    <Class IRI="http://www.w3.org/ns/prov#Entity"/>
1859    </SubClassOf>
1860    <SubClassOf>
1861    <Class IRI="#preservationPolicy"/>
1862    <Class IRI="#policy"/>
1863    </SubClassOf>
1864    <SubClassOf>
1865    <Class IRI="#quality"/>
1866    <Class IRI="#selectionPolicy"/>
1867    </SubClassOf>
1868    <SubClassOf>
1869    <Class IRI="#reporting"/>
1870    <Class IRI="#termsOfUse"/>
1871    </SubClassOf>
1872    <SubClassOf>
1873    <Class IRI="#repositorySystem"/>
1874    <Class IRI="http://www.w3.org/ns/prov#SoftwareAgent"/>
1875    </SubClassOf>
1876    <SubClassOf>
1877    <Class IRI="#researchData"/>
1878    <Class IRI="http://www.w3.org/ns/prov#Entity"/>
1879    </SubClassOf>
1880    <SubClassOf>
1881    <Class IRI="#retrievalMetadata"/>
1882    <Class IRI="#metadata"/>
1883    </SubClassOf>
1884    <SubClassOf>
1885    <Class IRI="#reusability"/>
1886    <Class IRI="#selectionPolicy"/>
1887    </SubClassOf>
1888    <SubClassOf>
1889    <Class IRI="#secondaryUseProhibited"/>
1890    <Class IRI="#termsOfUse"/>
1891    </SubClassOf>
1892    <SubClassOf>
1893    <Class IRI="#selectionPolicy"/>
1894    <Class IRI="#policy"/>
1895    </SubClassOf>
1896    <SubClassOf>
1897    <Class IRI="#selectionResults"/>
1898    <Class IRI="#curationRecord"/>
1899    </SubClassOf>
1900    <SubClassOf>
1901    <Class IRI="#size"/>
1902    <Class IRI="#selectionPolicy"/>
1903    </SubClassOf>
1904    <SubClassOf>
1905    <Class IRI="#softwareRegistry"/>
1906    <Class IRI="http://www.w3.org/ns/prov#Entity"/>
1907    </SubClassOf>
```

```
1908   <SubClassOf>
1909   <Class IRI="#sourceCode"/>
1910   <Class IRI="http://www.w3.org/ns/prov#Entity"/>
1911   </SubClassOf>
1912   <SubClassOf>
1913   <Class IRI="#sourceCode"/>
1914   <ObjectSomeValuesFrom>
1915   <ObjectProperty IRI="http://xmlns.com/foaf/0.1/primaryTopic"/>
1916   <Class IRI="#researchData"/>
1917   </ObjectSomeValuesFrom>
1918   </SubClassOf>
1919   <SubClassOf>
1920   <Class IRI="#submittedData"/>
1921   <Class IRI="#researchData"/>
1922   </SubClassOf>
1923   <SubClassOf>
1924   <Class IRI="#successionPlan"/>
1925   <Class IRI="#policy"/>
1926   </SubClassOf>
1927   <SubClassOf>
1928   <Class IRI="#systemAdministrator"/>
1929   <Class IRI="http://www.w3.org/ns/prov#Person"/>
1930   </SubClassOf>
1931   <SubClassOf>
1932   <Class IRI="#technicalInformation"/>
1933   <Class IRI="#metadata"/>
1934   </SubClassOf>
1935   <SubClassOf>
1936   <Class IRI="#termsOfUse"/>
1937   <Class IRI="http://www.w3.org/ns/prov#Entity"/>
1938   </SubClassOf>
1939   <SubClassOf>
1940   <Class IRI="#typeOfResource"/>
1941   <Class IRI="#selectionPolicy"/>
1942   </SubClassOf>
1943   <SubClassOf>
1944   <Class IRI="#usageResults"/>
1945   <Class IRI="#curationRecord"/>
1946   </SubClassOf>
1947   <SubClassOf>
1948   <Class IRI="#useLatestVersion"/>
1949   <Class IRI="#termsOfUse"/>
1950   </SubClassOf>
1951   <SubClassOf>
1952   <Class IRI="#versionInformation"/>
1953   <Class IRI="#metadata"/>
1954   </SubClassOf>
1955   <SubClassOf>
1956   <Class IRI="#visualizedData"/>
1957   <Class IRI="#researchData"/>
1958   </SubClassOf>
1959   <SubClassOf>
1960   <Class IRI="#visualizedData"/>
1961   <ObjectMinCardinality cardinality="1">
1962   <ObjectProperty IRI="http://www.w3.org/ns/prov#wasDerivedFrom"/>
1963   <ObjectUnionOf>
1964   <Class IRI="#curatedData"/>
1965   <Class IRI="#submittedData"/>
1966   </ObjectUnionOf>
1967   </ObjectMinCardinality>
1968   </SubClassOf>
1969   <SubClassOf>
1970   <Class IRI="#waiver"/>
1971   <Class IRI="#termsOfUse"/>
1972   </SubClassOf>
1973   <DisjointClasses>
1974   <Class IRI="#DataCurationActivities"/>
1975   <Class IRI="#DataPreservationActivities"/>
1976   </DisjointClasses>
1977   <DisjointClasses>
1978   <Class IRI="#RiskManagement"/>
1979   <Class IRI="#Selection"/>
1980   </DisjointClasses>
1981   <AnnotationAssertion>
```

```
1982   <Annotation>
1983   <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
1984   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
1985   </Annotation>
1986   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
1987   <IRI>#ActivatingMetadataBrokerage</IRI>
1988   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Active dissemination of a data set's metadata to search and discovery
       services (e.g., article databases, catalogs, web-based indexes) for federated search and discovery.</Literal>
1989   </AnnotationAssertion>
1990   <AnnotationAssertion>
1991   <Annotation>
1992   <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
1993   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
1994   </Annotation>
1995   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
1996   <IRI>#AllowingFileDownload</IRI>
1997   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Allow access to the data materials by authorized third parties.</Lite
       ral>
1998   </AnnotationAssertion>
1999   <AnnotationAssertion>
2000   <Annotation>
2001   <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2002   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://www.dcc.ac.uk/guidance/how-guides/appraise-select-data</Literal>
2003   </Annotation>
2004   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2005   <IRI>#Appraisal</IRI>
2006   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Appraisal is the process whereby some records are selected for retent
       ion, others (the great majority) are deemed of insufficient value to justify permanent retention.</Literal>
2007   </AnnotationAssertion>
2008   <AnnotationAssertion>
2009   <Annotation>
2010   <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2011   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2012   </Annotation>
2013   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2014   <IRI>#ArrangementAndDescription</IRI>
2015   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">The re-organization of files (e.g., new folder directory structure) i
       n a dataset that may also involve the creation of new file names, file descriptions, and the recording of technical metadata inherent to
       the files (e.g., date last modified).</Literal>
2016   </AnnotationAssertion>
2017   <AnnotationAssertion>
2018   <Annotation>
2019   <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2020   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2021   </Annotation>
2022   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2023   <IRI>#Authentication</IRI>
2024   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">The process of confirming the identity of a person, generally the dep
       ositor, who is contributing data to the data repository. (e.g., password authentication or authorization via digital signature). Used for
       tracking provenance of the data files.</Literal>
2025   </AnnotationAssertion>
2026   <AnnotationAssertion>
2027   <Annotation>
2028   <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2029   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2030   </Annotation>
2031   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2032   <IRI>#CeaseDataCuration</IRI>
2033   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Plan for any contingencies that will ultimately terminate access to t
       he data. For example, providing tombstones or metadata records for data that have been deselected and removed from stewardship.</Literal>
2034   </AnnotationAssertion>
2035   <AnnotationAssertion>
2036   <Annotation>
2037   <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2038   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2039   </Annotation>
2040   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2041   <IRI>#ChainOfCustody</IRI>
2042   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Intentional recording of provenance metadata of the files (e.g., meta
       data about who created the file, when it was last edited, etc.) in order to preserve file authenticity when data are transferred to third
       -parties.</Literal>
2043   </AnnotationAssertion>
2044   <AnnotationAssertion>
2045   <Annotation>
```

```
2046   <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2047   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2048   </Annotation>
2049   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2050   <IRI>#CodeReview</IRI>
2051   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Run and validate computer code (e.g., look for missing files and/or e
       rrors) in order to find mistakes overlooked in the initial development phase, improving the overall quality of software.</Literal>
2052   </AnnotationAssertion>
2053   <AnnotationAssertion>
2054   <Annotation>
2055   <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2056   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2057   </Annotation>
2058   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2059   <IRI>#ConnectingDiscoveryServices</IRI>
2060   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">To connect services that incorporate machine-based search and retriev
       al functionality that help users identify what data exist, where the data are located, and how can they be accessed (e.g., full-text inde
       xing or web optimization).</Literal>
2061   </AnnotationAssertion>
2062   <AnnotationAssertion>
2063   <Annotation>
2064   <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2065   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2066   </Annotation>
2067   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2068   <IRI>#Contextualization</IRI>
2069   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Use metadata to link the data set to related publications, dissertati
       ons, and/or projects that provide added context to how the data were generated and why.</Literal>
2070   </AnnotationAssertion>
2071   <AnnotationAssertion>
2072   <Annotation>
2073   <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2074   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">http://doi.org/10.5334/dsj-2020-053</Literal>
2075   </Annotation>
2076   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2077   <IRI>#Contracts</IRI>
2078   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">An agreement with a research partner, contractor, etc. that restricts
       the data publishing in joint research or contract research.</Literal>
2079   </AnnotationAssertion>
2080   <AnnotationAssertion>
2081   <Annotation>
2082   <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2083   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2084   </Annotation>
2085   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2086   <IRI>#Conversion</IRI>
2087   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">In effort to increase the usability of a data set, the information is
       transferred into digital file formats (e.g., analog data keyed into a database).</Literal>
2088   </AnnotationAssertion>
2089   <AnnotationAssertion>
2090   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2091   <IRI>#CreatingLandingPage</IRI>
2092   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">A web page associated with some dataset or identifier.</Literal>
2093   </AnnotationAssertion>
2094   <AnnotationAssertion>
2095   <Annotation>
2096   <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2097   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2098   </Annotation>
2099   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2100   <IRI>#DataCleaning</IRI>
2101   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">A process used to improve data quality by detecting and correcting (o
       r removing) defects &amp; errors in data.</Literal>
2102   </AnnotationAssertion>
2103   <AnnotationAssertion>
2104   <Annotation>
2105   <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2106   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2107   </Annotation>
2108   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2109   <IRI>#DataCurationActivities</IRI>
2110   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">The encompassing work and actions taken by curators of a data reposit
       ory in order to provide meaningful and enduring access to data.</Literal>
2111   </AnnotationAssertion>
```

```
2112    <AnnotationAssertion>
2113    <Annotation>
2114    <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2115    <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://casrai.org/term/evaluation/</Literal>
2116    </Annotation>
2117    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2118    <IRI>#DataEvaluation</IRI>
2119    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Evaluation is a decision about significance, value, or quality of som
        ething, based on careful study of its good and bad features.</Literal>
2120    </AnnotationAssertion>
2121    <AnnotationAssertion>
2122    <Annotation>
2123    <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2124    <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">http://doi.org/10.5334/dsj-2020-053</Literal>
2125    </Annotation>
2126    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2127    <IRI>#DataPolicy</IRI>
2128    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Where the research funder has a policy on limited data sharing for th
        e research to be funded, or where a strategic business decision restricts the data publishing relating to pending industrial property rig
        hts or research data where the commercialization of the research results is envisaged.</Literal>
2129    </AnnotationAssertion>
2130    <AnnotationAssertion>
2131    <Annotation>
2132    <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2133    <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://casrai.org/term/data-processing/</Literal>
2134    </Annotation>
2135    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2136    <IRI>#DataProcessing</IRI>
2137    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">A generic concept referring to all kinds of procedures being executed
        on data at any point in the data life cycle.</Literal>
2138    </AnnotationAssertion>
2139    <AnnotationAssertion>
2140    <Annotation>
2141    <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2142    <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://casrai.org/term/data-publication/</Literal>
2143    </Annotation>
2144    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2145    <IRI>#DataPublishing</IRI>
2146    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">The release of research data, associated metadata, accompanying docum
        entation, and software code (in cases where the raw data have been processed or manipulated) for re-use and analysis in such a manner tha
        t they can be discovered on the Web and referred to in a unique and persistent way. Data publishing occurs via dedicated data repositorie
        s and/or (data) journals which ensure that the published research objects are well documented, curated, archived for the long term, inter
        operable, citable, quality assured and discoverable – all aspects of data publishing that are important for future reuse of data by third
        party end-users.</Literal>
2147    </AnnotationAssertion>
2148    <AnnotationAssertion>
2149    <Annotation>
2150    <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2151    <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2152    </Annotation>
2153    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2154    <IRI>#DataVisualization</IRI>
2155    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">The presentation of pictorial and/or graphical representations of a d
        ata set used to identify patterns, detect errors, and/or demonstrate the extent of a data set to third party users.</Literal>
2156    </AnnotationAssertion>
2157    <AnnotationAssertion>
2158    <Annotation>
2159    <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2160    <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2161    </Annotation>
2162    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2163    <IRI>#Deidentification</IRI>
2164    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Redacting or removing personally identifiable or protected informatio
        n (e.g., sensitive geographic locations) from a dataset prior to sharing with third-parties.</Literal>
2165    </AnnotationAssertion>
2166    <AnnotationAssertion>
2167    <Annotation>
2168    <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2169    <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2170    </Annotation>
2171    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2172    <IRI>#DepositAgreement</IRI>
2173    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">The certification by the data author (or depositor) that the data con
        form to all policies and conditions (e.g., do not violate any legal restrictions placed on the data) and are fit for deposit into the rep
```

```
        ository. A deposit agreement may also include rights transfer to the repository for ongoing stewardship.</Literal>
2174    </AnnotationAssertion>
2175    <AnnotationAssertion>
2176    <Annotation>
2177    <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2178    <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">http://doi.org/10.5334/dsj-2020-053</Literal>
2179    </Annotation>
2180    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2181    <IRI>#DiplomaticOrNationalSecurity</IRI>
2182    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Research data pertaining to national security. Data related to the de
        velopment of weapons of mass destruction, etc. (as defined in the Foreign Exchange and Foreign Trade Act) and defense secrets (the Self-D
        efense Forces). law), important data that may affect national life (e.g., domestic energy (e.g., location of resources, blueprints for cr
        itical equipment, etc.).</Literal>
2183    </AnnotationAssertion>
2184    <AnnotationAssertion>
2185    <Annotation>
2186    <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2187    <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">http://doi.org/10.5334/dsj-2020-053</Literal>
2188    </Annotation>
2189    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2190    <IRI>#DisciplinaryCustoms</IRI>
2191    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Practices and standards in a specific discipline or research communit
        y that limit the data publishing. In some cases this is stated as an international treaty, but in others it is not always explicitly stat
        ed.</Literal>
2192    </AnnotationAssertion>
2193    <AnnotationAssertion>
2194    <Annotation>
2195    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2196    <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2197    </Annotation>
2198    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2199    <IRI>#DisplayingDataCitation</IRI>
2200    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Display of a recommended bibliographic citation for a dataset to enab
        le appropriate attribution by third-party users in order to formally incorporate data reuse as part of the scholarly ecosystem.</Literal>
2201    </AnnotationAssertion>
2202    <AnnotationAssertion>
2203    <Annotation>
2204    <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2205    <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2206    </Annotation>
2207    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2208    <IRI>#Documentation</IRI>
2209    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">To describe Information about any necessary information to use and un
        derstand the data. Documentation may be structured (e.g., a code book) or unstructured (e.g., a plain text "Readme" file).</Literal>
2210    </AnnotationAssertion>
2211    <AnnotationAssertion>
2212    <Annotation>
2213    <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2214    <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2215    </Annotation>
2216    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2217    <IRI>#Embargo</IRI>
2218    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">To restrict or mediate access to a data set, usually for a set period
        of time. In some cases an embargo may be used to protect not only access, but any knowledge that the data exist.</Literal>
2219    </AnnotationAssertion>
2220    <AnnotationAssertion>
2221    <Annotation>
2222    <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2223    <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2224    </Annotation>
2225    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2226    <IRI>#Emulation</IRI>
2227    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Provide legacy system configurations in modern equipment in order to
        ensure long-term usability of data. (E.g., arcade games emulated on modern web-browsers)</Literal>
2228    </AnnotationAssertion>
2229    <AnnotationAssertion>
2230    <Annotation>
2231    <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2232    <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2233    </Annotation>
2234    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2235    <IRI>#FileAuditing</IRI>
2236    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Periodic review of the digital integrity of the data files and taking
        action when needed to protect data from digital erosion (e.g., bitrot) and/or hardware failure.</Literal>
```

```
2237   </AnnotationAssertion>
2238   <AnnotationAssertion>
2239   <Annotation>
2240   <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2241   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2242   </Annotation>
2243   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2244   <IRI>#FileFormatTransformation</IRI>
2245   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Transform files into open, non-proprietary file formats that broaden
       the potential for long-term reuse and ensure that additional preservation actions might be taken in the future.</Literal>
2246   </AnnotationAssertion>
2247   <AnnotationAssertion>
2248   <Annotation>
2249   <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2250   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2251   </Annotation>
2252   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2253   <IRI>#FileInventoryOrManifest</IRI>
2254   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">The data files are inspected periodically and the number, file types
       (extensions), and file sizes of the data are understood and documented. Any missing, duplicate, or corrupt (e.g., unable to open) files a
       re discovered.</Literal>
2255   </AnnotationAssertion>
2256   <AnnotationAssertion>
2257   <Annotation>
2258   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2259   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2260   </Annotation>
2261   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2262   <IRI>#FileRenaming</IRI>
2263   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">To rename files in a dataset, often to standardize and/or reflect imp
       ortant metadata.</Literal>
2264   </AnnotationAssertion>
2265   <AnnotationAssertion>
2266   <Annotation>
2267   <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2268   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2269   </Annotation>
2270   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2271   <IRI>#FileValidation</IRI>
2272   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">A computational process to ensure that the intended data transfer to
       a repository was perfect and complete using means such as generating and validating file checksums (e.g., test if a digital file has chan
       ged at the bit level) and format validation to ensure that file types match their extensions.</Literal>
2273   </AnnotationAssertion>
2274   <AnnotationAssertion>
2275   <Annotation>
2276   <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2277   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2278   </Annotation>
2279   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2280   <IRI>#FormulateSuccessionPlanning</IRI>
2281   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Planning for contingency, and/or escrow arrangements, in the case tha
       t the repository (or other entity responsible) ceases to operate or the institution substantially changes its scope.</Literal>
2282   </AnnotationAssertion>
2283   <AnnotationAssertion>
2284   <Annotation>
2285   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2286   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2287   </Annotation>
2288   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2289   <IRI>#GeneratingFulltextIndexing</IRI>
2290   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Enhance the data for discovery purposes by generating search-engine-o
       ptimized formats of the text inherent to the data.</Literal>
2291   </AnnotationAssertion>
2292   <AnnotationAssertion>
2293   <Annotation>
2294   <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2295   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2296   </Annotation>
2297   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2298   <IRI>#Indexing</IRI>
2299   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Verify all metadata provided by the author and crosswalk to descripti
       ve and administrative metadata compliant with a standard format for repository interoperability.</Literal>
2300   </AnnotationAssertion>
2301   <AnnotationAssertion>
```

```
2302    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2303    <IRI>#Ingest</IRI>
2304    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">The process of turning a Submission Information Package (SIP) into an
        Archival Information Package (AIP), i.e. putting data into a digital archive (OAIS term)</Literal>
2305    </AnnotationAssertion>
2306    <AnnotationAssertion>
2307    <Annotation>
2308    <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2309    <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2310    </Annotation>
2311    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2312    <IRI>#Interoperability</IRI>
2313    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Formatting the data using a disciplinary standard for better integrat
        ion with other datasets and/or systems.</Literal>
2314    </AnnotationAssertion>
2315    <AnnotationAssertion>
2316    <Annotation>
2317    <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2318    <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2319    </Annotation>
2320    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2321    <IRI>#MaintainingContactInformation</IRI>
2322    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Keep up-to-date contact information for the data authors and/or the c
        ontact persons in order to facilitate connection with third-party users. Often involves managing ephemeral information that will change o
        ver time.</Literal>
2323    </AnnotationAssertion>
2324    <AnnotationAssertion>
2325    <Annotation>
2326    <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2327    <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2328    </Annotation>
2329    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2330    <IRI>#MetadataGeneration</IRI>
2331    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">A generation process of metadata which is the Information about a dat
        a set that is structured (often in machine-readable format) for purposes of search and retrieval. Metadata elements may include basic inf
        ormation (e.g. title, author, date created, etc.) and/or specific elements inherent to datasets (e.g., spatial coverage, time periods).</
        Literal>
2332    </AnnotationAssertion>
2333    <AnnotationAssertion>
2334    <Annotation>
2335    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2336    <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2337    </Annotation>
2338    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2339    <IRI>#Migration</IRI>
2340    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Monitor and anticipate file format obsolescence and, as needed, trans
        form obsolete file formats to new formats as standards and use dictate.</Literal>
2341    </AnnotationAssertion>
2342    <AnnotationAssertion>
2343    <Annotation>
2344    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2345    <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2346    </Annotation>
2347    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2348    <IRI>#MintingPersistentIdentifier</IRI>
2349    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">To mint a permanent URL (or Uniform Resource Locator) that is monitor
        ed by an authority to ensure a stable web location for consistent citation and long-term discoverability. Provides redirection when neces
        sary. E.g., a Digital Object Identifier or DOI.</Literal>
2350    </AnnotationAssertion>
2351    <AnnotationAssertion>
2352    <Annotation>
2353    <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2354    <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2355    </Annotation>
2356    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2357    <IRI>#PeerReview</IRI>
2358    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">The review of a data set by an expert with similar credentials and su
        bject knowledge as the data creator for the purposes of validating the soundness and trustworthiness of the file contents.</Literal>
2359    </AnnotationAssertion>
2360    <AnnotationAssertion>
2361    <Annotation>
2362    <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2363    <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">http://doi.org/10.5334/dsj-2020-053</Literal>
2364    </Annotation>
```

```
2365  <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2366  <IRI>#PersonalInformation</IRI>
2367  <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">It stipulates the handling of data that can identify individuals. It
      includes guidelines that define individual policies on anonymization and information disclosure.</Literal>
2368  </AnnotationAssertion>
2369  <AnnotationAssertion>
2370  <Annotation>
2371  <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2372  <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2373  </Annotation>
2374  <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2375  <IRI>#ProvidingRestrictedAccess</IRI>
2376  <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">In order to maintain the privacy of research subjects without losing
      integral components of the data, some data access will be protected and/or mediated to individuals that meet predefined criteria.</Litera
      l>
2377  </AnnotationAssertion>
2378  <AnnotationAssertion>
2379  <Annotation>
2380  <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2381  <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2382  </Annotation>
2383  <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2384  <IRI>#QualityAssurance</IRI>
2385  <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Ensure that all documentation and metadata are comprehensive and comp
      lete. Example actions might include: open and run the data files; inspect the contents in order to validate, clean, and/or enhance data f
      or future use; look for missing documentation about codes used, the significance of "null" and "blank" values, or unclear acronyms.</Lite
      ral>
2386  </AnnotationAssertion>
2387  <AnnotationAssertion>
2388  <Annotation>
2389  <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2390  <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2391  </Annotation>
2392  <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2393  <IRI>#RegisteringSoftware</IRI>
2394  <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Maintain copies of modern and obsolete versions of software (and any
      relevant code libraries) so that data may be opened/used overtime.</Literal>
2395  </AnnotationAssertion>
2396  <AnnotationAssertion>
2397  <Annotation>
2398  <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2399  <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2400  </Annotation>
2401  <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2402  <IRI>#Restructure</IRI>
2403  <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Organize and/or reformat poorly structured data files to clarify the
      ir meaning and importance.</Literal>
2404  </AnnotationAssertion>
2405  <AnnotationAssertion>
2406  <Annotation>
2407  <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2408  <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2409  </Annotation>
2410  <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2411  <IRI>#RightsManagement</IRI>
2412  <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">The process of tracking and managing ownership and copyright inherent
      to a data set as well as monitoring conditions and policies for access and reuse (e.g., licenses and data use agreements).</Literal>
2413  </AnnotationAssertion>
2414  <AnnotationAssertion>
2415  <Annotation>
2416  <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2417  <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2418  </Annotation>
2419  <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2420  <IRI>#RiskManagement</IRI>
2421  <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">The process of reviewing data for known risks such as confidentiality
      issues inherent to human subjects data, sensitive information (e.g., sexual histories, credit card information) or data regulated by law
      (e.g. HIPAA, FERPA) and taking actions to reject or facilitate remediation (e.g., de-identification services) when necessary.</Literal>
2422  </AnnotationAssertion>
2423  <AnnotationAssertion>
2424  <Annotation>
2425  <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2426  <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2427  </Annotation>
```

```
2428   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2429   <IRI>#SecuringStorage</IRI>
2430   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Data files are properly stored in a well-configured (in terms of hard
       ware and software) storage environment that is routinely backed-up and physically protected. Perform routine fixity checks (to detect deg
       radation or loss) and provide recovery services as needed.</Literal>
2431   </AnnotationAssertion>
2432   <AnnotationAssertion>
2433   <Annotation>
2434   <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2435   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2436   </Annotation>
2437   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2438   <IRI>#Selection</IRI>
2439   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">The process of a successful appraisal. The data are determined approp
       riate for acceptance and ingest into the repository according to local collection policy and practice.</Literal>
2440   </AnnotationAssertion>
2441   <AnnotationAssertion>
2442   <Annotation>
2443   <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2444   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2445   </Annotation>
2446   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2447   <IRI>#SettingTermsOfUse</IRI>
2448   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Information provided to end users of a data set that outline the requ
       irements or conditions for use (e.g., a Creative Commons License).</Literal>
2449   </AnnotationAssertion>
2450   <AnnotationAssertion>
2451   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2452   <IRI>#SubmitData</IRI>
2453   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">To submit data to start data curation activities.</Literal>
2454   </AnnotationAssertion>
2455   <AnnotationAssertion>
2456   <Annotation>
2457   <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2458   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2459   </Annotation>
2460   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2461   <IRI>#TechnologyMonitoringAndRefreshing</IRI>
2462   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Formal, periodic review and assessment to ensure responsiveness to te
       chnological developments and evolving requirements of the digital infrastructure and hardware storing the data.</Literal>
2463   </AnnotationAssertion>
2464   <AnnotationAssertion>
2465   <Annotation>
2466   <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2467   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2468   </Annotation>
2469   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2470   <IRI>#TrackingUseAnalytics</IRI>
2471   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Monitor and record how often data are viewed, requested, and/or downl
       oaded. Track and report reuse metrics, such as data citations and impact measures for the data over time.</Literal>
2472   </AnnotationAssertion>
2473   <AnnotationAssertion>
2474   <Annotation>
2475   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2476   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2477   </Annotation>
2478   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2479   <IRI>#Transcoding</IRI>
2480   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">With audio and video files, detect technical metadata (min resolutio
       n, audio/video codec) and encode files in ways that optimize reuse and long-term preservation actions. (E.g, Convert QuickTime files to M
       PEG4).</Literal>
2481   </AnnotationAssertion>
2482   <AnnotationAssertion>
2483   <Annotation>
2484   <AnnotationProperty abbreviatedIRI="rdfs:isDefinedBy"/>
2485   <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">https://hdl.handle.net/11299/188638</Literal>
2486   </Annotation>
2487   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2488   <IRI>#Versioning</IRI>
2489   <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Provide mechanisms to ingest new versions of the data overtime that i
       ncludes metadata describing the version history and any changes made for each version.</Literal>
2490   </AnnotationAssertion>
2491   <AnnotationAssertion>
2492   <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
```

```
2493    <IRI>#creditOnTheResults</IRI>
2494    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">This condition of use requires displaying the creator's name and data
        URL information</Literal>
2495    </AnnotationAssertion>
2496    <AnnotationAssertion>
2497    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2498    <IRI>#feeForUse</IRI>
2499    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">This condition of use requires some fee for data use.</Literal>
2500    </AnnotationAssertion>
2501    <AnnotationAssertion>
2502    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2503    <IRI>#imposeTheSameConditions</IRI>
2504    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">This condition of use requires the same condition of use under redist
        ribution and combination of multiple pieces of data. Unlike "Attribution," it may prevent them from combining the data with other sets th
        at have an incompatible license.</Literal>
2505    </AnnotationAssertion>
2506    <AnnotationAssertion>
2507    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2508    <IRI>#improperUse</IRI>
2509    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">This condition of use prohibits "Improper use" of data in the data pr
        ocessing phase. Note that the definition of inappropriate use will probably depend on the conventions of the field.</Literal>
2510    </AnnotationAssertion>
2511    <AnnotationAssertion>
2512    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2513    <IRI>#noDelivs</IRI>
2514    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">This condition of use prohibits data publishing after any modificatio
        n. The case in which the prohibition of modification is effective is presumed to largely depend on the type of data and the manner of use
        (e.g., image data that are practically a work of art).</Literal>
2515    </AnnotationAssertion>
2516    <AnnotationAssertion>
2517    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2518    <IRI>#nonCommercial</IRI>
2519    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">This condition of use requires the "non-commercial" use of data. Als
        o, the limitation is on redistribution and data processing (e.g., selling visualizations derived from the data).</Literal>
2520    </AnnotationAssertion>
2521    <AnnotationAssertion>
2522    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2523    <IRI>#reporting</IRI>
2524    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">This condition of use requires post-reuse reporting, suggesting that
        the objective is to know detailed usage practices rather than mechanical access statistics.</Literal>
2525    </AnnotationAssertion>
2526    <AnnotationAssertion>
2527    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2528    <IRI>#secondaryUseProhibited</IRI>
2529    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">This condition of use prohibits the secondary use of data. This licen
        se clearly prohibits data redistribution, translation, or adaption and is intended for one-on-one use of its original form.</Literal>
2530    </AnnotationAssertion>
2531    <AnnotationAssertion>
2532    <Annotation>
2533    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2534    <Literal datatypeIRI="http://www.w3.org/2001/XMLSchema#anyURI">http://doi.org/10.5334/dsj-2020-053</Literal>
2535    </Annotation>
2536    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2537    <IRI>#termsOfUse</IRI>
2538    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">"Condition of use" generally means use permission/prohibition, obliga
        tions, and constraints based on the types of research data.</Literal>
2539    </AnnotationAssertion>
2540    <AnnotationAssertion>
2541    <AnnotationProperty IRI="http://www.w3.org/2004/02/skos/core#altLabel"/>
2542    <IRI>#termsOfUse</IRI>
2543    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Conditions of Use</Literal>
2544    </AnnotationAssertion>
2545    <AnnotationAssertion>
2546    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2547    <IRI>#useLatestVersion</IRI>
2548    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">This condition of use is used to limit the use of data to the latest
        one. Data in the past cannot be reused.</Literal>
2549    </AnnotationAssertion>
2550    <AnnotationAssertion>
2551    <AnnotationProperty abbreviatedIRI="rdfs:comment"/>
2552    <IRI>#waiver</IRI>
2553    <Literal datatypeIRI="http://www.w3.org/2000/01/rdf-schema#Literal">Waives copyright and all other related rights</Literal>
2554    </AnnotationAssertion>
2555    </Ontology>
```

```
2556
2557
2558
2559    <!-- Generated by the OWL API (version 4.5.9.2019-02-01T07:24:44Z) https://github.com/owlcs/owlapi -->
2560
```

# Appendix 4. Topic guide of "Basic research on licensing and legal interoperability of research data Questions"

We used a topic guide in Section 5.3.1 to share our interview outline with the interviewee. We set three major sections within the topic guide: "Sharing and publishing of research data," "Regarding granting licenses to research data," and "Licensing of research data and promotion of legal interoperability." Details of each section are as follows. Note that the topic guide sent to interviewees was prepared in Japanese.

**List of sections**

1. Sharing and publishing of research data
- Overview and characteristics of your data
    e.g., ownership of data, public or commercial value, existence of personal information

- Current status of data sharing and publication in your research, your institution, and your research community

- Difficulties in sharing and publishing research data

- What you would like to ask/prohibit users from doing with the research data you have made public

    e.g., prohibition or suspension of use and penalties for license violations

2. Regarding granting licenses to research data
- The type of license used, the stipulations of the license granted, and information/guidelines referred to at the time of granting the license

- Difficulties in granting licenses to research data

3. Licensing of research data and promotion of legal interoperability
- The needs for licensing and legal interoperability of research data

- Personal views on existing licenses, guidelines, and existing discussions

- The degree of protection of rights to the research data sought and the basis for the request for protection of rights (and why such a request is made by the user)

- Expectations regarding licensing, rights management, data release support, etc.

4. Other comments

# Appendix 5. Guidelines for specifying conditions of use in research data publishing

In Section 5.5 we introduced the "Guidelines for specifying conditions of use in research data publishing" (Minamiyama et al., 2020) to help researchers and stakeholders understand and make appropriate publication decisions. These guidelines provide necessary information and examples that should be considered when publishing research data with an interdisciplinary perspective. The original document was written in Japanese, and a tentative English translation of the document is shown below.

December 25, 2019

# Guidelines for specifying conditions of use
# in research data publishing

Research Data License Subcommittee

under the Research Data Utilization Forum (RDUF)

## Index

Introduction. Five Questions on the Research Data Publishing and the Specification of Conditions of Use

<u>Purpose and objectives</u>

These guidelines are for the research data publishing and the specification of conditions of use developed by the Research Data License Subcommittee under the Research Data Utilization Forum (RDUF) [1]. The basic policy of open science in Japan is to expand the utilization of research results funded by public research funds as much as possible[2], but some types of research data are exceptions. Therefore, these guidelines aim to enable data providers to publish research data under appropriate conditions of use by organizing information and examples that generally require attention when publishing research data, along with the decision-making process. It also expects to be used as a tool for data reusers to easily understand the background of the conditions of use required by the data provider.

(1) When the data provider specifies the conditions of use

    To enable those who wish to publish research data (individual researchers, teams, and repository managers) to inform third parties of their conditions of use concisely. Possible scenarios are as follows: a) when publishing research data underlying a research paper and b) when publishing research data itself as research results. These guidelines prevent unauthorized reproduction, plagiarism, inappropriate processing, and trouble with interested parties in research data publishing.

(2) When the data reuser checks the existing conditions of use

    When the researcher acquires and reuses published research data, they can easily understand the conditions of use required by the data provider.
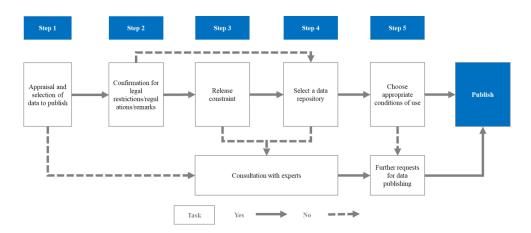
<u>How to use these guidelines</u>

First, select your research data to be published in Q1. Next, confirm any external constraints as listed in Q2. In Q3, you can confirm the processes for enabling research data publishing for the external constraints identified in Step 2. Even if external constraints exist, it may be

---

[1] Research Data Utilization Forum (RDUF) was established in May 2016. The principal mission of the RDUF is to encourage the utilization of research data based on open science. The RDUF is willing to support the communication on research data management and open science in various inter-disciplinary and inter-sectional views.

[2] 5th Science and Technology Basic Plan.

https://www8.cao.go.jp/cstp/kihonkeikaku/index5.html

possible to publish your research data by setting a certain embargo period, as described in Q3. If your research data can be made public, refer to Q4 to select a data repository. Finally, specify the conditions of use for your research data in Q5.

Publishing flow and licensing scenarios for research data

<u>Usage Notes</u>

- These guidelines organize information and examples that generally need to be kept in mind when publishing research data and the decision-making process. We do not expect to cover all academic fields; When publishing research data, please check general guidelines and policies related to research ethics and the handling of research data in your own academic field.

- When the researchers share research data, even if one wants a particular publication method or conditions of use, the other researchers may not have the same intention. Please confirm their intent well in advance.

- The conditions of use recommended by these guidelines are proposed regarding the Creative Commons 4.0 International License. Note that we do not consider the compatibility with different versions.

- If you wish to specify conditions of use for databases and/or repositories, we recommend seeking advice from experts in the relevant legal systems. The same advice should be given when there is more than one interested party, such as in the case of research conducted through industry-academia-government collaboration. Also, if you wish to state your conditions of use for research data as an organization or institution, please consider developing your institution's data policy.

<u>Conditions of use for this document/Disclaimers</u>

- Copyright of the text and figures in these guidelines belongs to the Research Data License Subcommittee under the Research Data Utilization Forum (RDUF). Except where otherwise noted, all materials are available under the Creative Commons Attribution 4.0 International License (CC-BY) terms. When using the materials, please clearly indicate the source information and any modifications regarding the following license notice:

  Source: "Guidelines for specifying conditions of use in research data publishing ver.1.0". Research Data License Subcommittee under the Research Data Utilization Forum (RDUF), 2019, 32p. https://doi.org/10.11502/rduf_license_guideline, (accessed YYYY-MM-DD).

- We do not guarantee the accuracy, certainty, fitness for purpose, or other quality of the statements in your context. The responsibility for all actions using these guidelines rests entirely with the user. The users themselves should make decisions based on the information obtained.

4

Q1. Appraisal and selection of data to publish

First, select the research data to be published in the whole study. The scope of the "research data" term varies from field to field, so these guidelines limit the scope of the term to the extent that it can be managed by electronic means. In other words, the term "research data" does not include physical materials such as samples (specimens, samples) or recording media (paper, disks, etc.) in these guidelines.

- Definition of "research data" (*Subject of these guidelines)
  Digital data used as a source information for scientific research. It includes a variety of formats, such as numerical, textual, image, audio, and video. Various designations may be used depending on the context in which the data is used; e.g., evidence data, source data, and derived data.

  Example:
  1) Evidence data
     Data underlying a research paper or research results is called "evidence data." Raw data may be published as the evidence data, or sometimes processed data are selected to publish.
  2) Source data
     The original data newly collected from the observation is called "source data (or primary data)." Researchers may observe and create the source data themselves, or it may be held by a third party (e.g., another researcher, a company, or a local government).
  3) Derived data
     Data created by derivation from source data is called "derived data." If the creator of the source data and the creator of the derived data are different, the data citation method and rights attribution tends to be complicated. It is necessary to pay attention to the source data version information.

(cf: Data not included in the "research data") (Outside the scope of these guidelines)
- Non research data
  Descriptions to explain the outline and status of the research data, as well as physical objects such as research notes, diaries, samples, etc.
  - Ex. Descriptions about research data (including metadata or meta-information)
  - Other records or logs (research notes, samples, and other physical objects)

5

- Copyrighted work

  For example:
  - ➢ Research papers
  - ➢ Books, derivative works by the authors, secondary works
  - ➢ Derivative works by third parties

- Research environment

  For example:
  - ➢ Databases
  - ➢ Software (e.g., Analysis and visualization programs, estimation models, machine learning algorithms)
  - ➢ Other source codes

## When publishing research data is required

Publishing research data may be required by your funding agency, publisher, or institution to promote research data reuse.

[Ex. 1] by funding agency
JST Policy on Open Access to Research Publications and Research Data Management (April 1, 2017)
https://www.jst.go.jp/pr/intro/openscience/guideline_openscience.pdf

[Ex. 2] by publisher
Elsevier. Research Data Guidelines
https://www.elsevier.com/authors/author-resources/research-data/data-guidelines

[Ex. 3] by institution (project data)
JAXA. ISAS Data Policy (March 14, 2018)
http://www.isas.jaxa.jp/researchers/data-policy/

[Ex.4] by institution (evidence data)
National Institute of Polar Research. National Institute of Polar Research Open Access Policy (November 24, 2017)
https://www.nipr.ac.jp/outline/activity/oap.html

6

Considerations for handling data not included in "research data"

- Metadata plays an essential role in informing the existence of research data. Even if the source or derived data is kept private, the metadata should be widely available for the public to search, view, and retrieve.

- Data not included in "research data" should be segregated and managed with the research data. If you treat past research notes or diaries as source data, be careful what you publish in digital form.

- The requirement of copyrighted works such as articles, papers, posters, slide materials, and projection materials differs from the research data to be published. For these works, consider publishing them in your institutional repository and applying standard licensing tools such as the Creative Commons License.

- The research environment and research data must be treated separately in the licensing context. You may refer to the following if you wish to specify a database license or a source code license.

[Ref. 1] Open Data Commons. Open Data Common Open Database License (ODbL).
https://opendatacommons.org/licenses/odbl/index.html

[Ref. 2] Choosealicense.com. Licenses.
https://choosealicense.com/licenses/

Q2. Confirmation for legal restrictions/regulations/remarks

There may be restrictions on publishing research data due to the sensitive content contained in the data (e.g., privacy information) or the research participant's request. Please confirm if your data falls into the descriptions contained in the following categories.

In cases of disciplinary customs restriction, including international treaty

- Individual disciplines and research communities may have conventions or standards regarding data release restrictions. The provisions of international treaties are indicated in some cases[3], but the provisions are not always explicitly stated[4].

  [Ex. 1] Convention on International Trade in Endangered Species of Wild Fauna and Flora (CITES)
  https://www.cites.org/eng

  [Ex. 2] Convention on the Means of Prohibiting and Preventing the Illicit Import, Export and Transfer of Ownership of Cultural Property
  http://portal.unesco.org/en/ev.php-URL_ID=13039&URL_DO=DO_TOPIC&URL_SECTION=201.html

In cases of containing personal information

- The laws of each country regulate the handling of personal information. In Japan, if the research entity is private, the Act on the Protection of Personal Information applies. If the research entity is an Incorporated Administrative Agency, including a National Research and Development Agency, the Act on the Protection of Personal Information Held by Incorporated Administrative Agencies, etc.
- As for discipline-specific regulations, for example, guidelines may be formulated for each field, with separate policies for anonymization and information disclosure.

---

[3] Other examples of publishing research data restrictions are mentioned within the guidelines prepared by RDA/CODATA, such as the protection of endangered species, cultural resources, sovereign genetic resources, and traditional knowledge.
Legal Interoperability of Research Data: Principles and Implementation Guidelines
https://doi.org/10.5281/zenodo.162241
[4] For example, materials may be withheld from the public due to the Bereaved family's request in literary research.

8

[Ref. 1] Personal Information Protection Commission, Government of Japan. "Laws and guidelines" (only in Japanese)
https://www.ppc.go.jp/personalinfo/legal/

[Ref. 2] Japan External Trade Organization(JETRO). "About General Data Protection Regulation (GDPR)" (only in Japanese)
https://www.jetro.go.jp/world/europe/eu/gdpr/

[Ref. 3] Ministry of Health, Labor and Welfare (Japan). "About research guidelines" (only in Japanese)
https://www.mhlw.go.jp/stf/seisakunitsuite/bunya/hokabunya/kenkyujigyou/i-kenkyu/index.html

In cases of Diplomatic / National security restriction

- Publishing research data on national security is regulated by law and includes data related to the weapons development of mass destruction (as defined in the Foreign Exchange and Foreign Trade Law) and defense secrets (defined in the Self Defense Forces Law).
- In other cases, there are special legal measures for data that may affect the lives of the public (e.g., location of domestic energy resources, blueprints of important facilities, etc.).

[Ref. 1] Japan Society for Intellectual Production. "Security Trade Control Guidelines for Researchers in universities and other institutions of higher education. Revised 2nd ed"
http://j-sip.org/info/pdf/anzenhosho1-1_2.pdf

In cases of keeping agreements, contracts, Intellectual Property rights

- In joint or contract research, it is necessary to comply with agreements with research partners, contractors, etc. regarding publishing research data.
- If you have some agreements or contracts with a private company (including a commercial publisher) that restricts publishing your data, you must comply with the terms of the agreements or contracts.

9

[Ref. 1] Ministry of Economy, Trade and Industry (Japan). "Operation guidelines for data management in contract research and development" (only in Japanese)
https://www.meti.go.jp/press/2017/12/20171227001/20171227001-1.pdf

[Ref. 2] Ministry of Economy, Trade and Industry (Japan). "Contract Guidelines on Utilization of AI and Data. Data Section"
https://www.meti.go.jp/press/2019/12/20191209001/20191209001.html

In cases of complying data policy

- Your institution may have a restricted data sharing policy. If your institution has an intellectual property policy or data policy[5], you must confirm the scope or embargo period within the policy.
- In some cases, publishing research data on industrial property rights pending application or research data expected to be commercialized may be restricted as a management strategy decision[6]. You also need to confirm the target data attribution.

[Ex. 1] National Institute for Environmental Studies. "NIES Data Policy" (only in Japanese)
https://www.nies.go.jp/kihon/kitei/kt_datapolicy.pdf

---

[5] In Japan, all national research and development agencies must have a data policy by 2020. Cabinet Office, Japan. "Integrated Innovation Strategy"
https://www8.cao.go.jp/cstp/tougosenryaku/index.html
Cabinet Office, Japan. "Guidelines for the Development of Data Policies in National Research and Development Agency"
https://www8.cao.go.jp/cstp/stsonota/datapolicy/datapolicy.html
[6] The "Open and Closed Strategy" is "a strategic choice to increase the company's profits by adopting an open model IP strategy of disclosing or licensing to other companies, in addition to a closed model IP strategy of keeping technologies and other information secret or implementing exclusive rights such as patent rights." The "Guidelines for the Formulation of Data Policies for National Research and Development Institutions" also calls for the formulation of policies based on this concept.
Japan Patent Office. "Open and Closed Strategy"
https://faq.inpit.go.jp/content/tradesecret/files/100578260.pdf

[Ex. 2] Teikyo University. "Intellectual Property policy in Teikyo University" (only in Japanese)
https://www.teikyo-u.ac.jp/affiliate/laboratory/tttc_center/policy.html

[Ex. 3] Japan Agency for Medical Research and Development. "Data sharing policy for realization of genomic medicine" (only in Japanese)
https://www.amed.go.jp/content/000023353.pdf
https://humandbs.biosciencedbc.jp/files/DAC/4th_meeting/2_ref_AMED_DSP.pdf

Q3. Release constraint

In most cases, even if there are some restrictions on publishing research data as described in Q2, it is possible to publish the data by applying appropriate data processing or allowing a certain embargo period[7]. Set the necessary conditions of use based on the following information[8] and show them to the data reuser:

In cases of disciplinary customs restriction, including international treaty

> You need to confirm the data publishing procedures with the corresponding national law if the disciplinary customs restriction is explicitly stated in an international treaty. In the absence of a specified period for restricted publication, you must set an appropriate period that considers the disciplinary practice and/or the treaty's purpose. In addition to checking with the IP department of your institution, you can consult with an expert if necessary.
>
> [Ref. 1] Ministry of Foreign Affairs of Japan. "Treaty Data Search"
> https://www3.mofa.go.jp/mofaj/gaiko/treaty/
>
> [Ref. 2 ] FAIRSharing. "FAIRSharing policies"
> https://fairsharing.org/policies/
> * Collection of field-specific policies/guidelines

---

[7] In principle, the legal protection of a copyrighted work expires 70 years after the author's death. Still, there is no corresponding provision for research data, so it is necessary to be more careful in setting an embargo period. From the viewpoint of protecting research papers, it is generally 12 months for science, engineering, and medicine, and typically 24 to 36 months for humanities and social sciences. However, in recent years, the evidence data tend to publish immediately.
[Ref.] U.S. Department of Health & Human Services. "NIH Public Access Policy Details"
https://publicaccess.nih.gov/policy.htm
[Ref.] Wiley. "Wiley's Self-Archiving Policy" https://authorservices.wiley.com/author-resources/Journal-Authors/licensing/self-archiving.html
[Ref.] SHERPA/RoMEO. "Publisher copyright policies & self-archiving"
http://sherpa.ac.uk/romeo/index.php
[8] There are some options for setting an embargo period, such as a timer, date/time-specified, and user-only limitation methods. These options will be used in some combination.

## In cases of containing personal information

Even if your research data contains some personal information, it can be published when you anonymize your data in an appropriate method.

[Ref. 1] Personal Information Protection Commission, Government of Japan. "Guidelines for the Act on the Protection of Personal Information. Anonymized Information section"
   https://www.ppc.go.jp/files/pdf/guidelines04.pdf

You can also refer to these particular guidelines to specify anonymizing methods in your fields:

[Ref. 2] Personal Information Protection Commission, Government of Japan. "Guidelines for Specific Fields"
https://www.ppc.go.jp/personalinfo/legal/guidelines/

[Ref. 3] Japan Pharmaceutical Manufacturers Association. "Institutional Review Board (IRB) overview and date of event"
http://www.jpma.or.jp/medicine/shinyaku/tiken/allotment/leaflet/009.html

## In cases of Diplomatic / National security restriction

If your research data is restricted by export control, military diversion, or other treaties, you must follow established procedures to consider whether or not to disclose the data. Please consult with the department in charge and follow the disclosure procedures.

[Ref. 1] Ministry of Finance, Japan. "Overview of Foreign Exchange and Foreign Trade Act"
https://www.mof.go.jp/international_policy/gaitame_kawase/gaitame/index.html

[Ref. 2] Ministry of Economy, Trade and Industry, Japan. "Security Export Control System in Japan"
https://www.meti.go.jp/policy/anpo/

If your research data is categorized in public records, it is subject to the Public Records and Archives Management Act. You must be considered for disclosure following the enforcement order of the Act. As above, please consult with the department in charge

13

and follow the disclosure procedures.

[Ref. 1] Public Records and Archives Management Act
https://elaws.e-gov.go.jp/search/elawsSearch/elaws_search/lsg0500/detail?lawId=421AC0000000066

[Ref. 2] Cabinet Office, Government of Japan. "Documents subject to management under the Public Records and Archives Management Act"
https://www8.cao.go.jp/chosei/koubun/about/bunsho/bunsho.html

[Ref. 3] Cabinet Secretariat, Japan. "Standard document retention period"
https://www.cas.go.jp/jp/koukai/hyoujunbunsho/anzenhosho.html

## In cases of keeping agreements, contracts, Intellectual Property rights

Based on the agreement or contract, set the publication date and time after confirming and coordinating the cases in which publication is possible. In addition to the statements related to the publication of research results, check the agreements after the joint research period has ended. Even in cases where there is no explicit agreement or contract, you should check with your IP department before publishing your research data.

[Ref. 1] Joint Research and Development Agreement (in Japanese)

第〇条（研究成果の公表等）
甲又は乙は、本契約の有効期間中及び契約終了後〇年間は、本共同研究によって得られた研究成果を公表又は第三者に開示しようとする場合には、その内容、時期、方法等について、書面により事前に相手方の承諾を受けるものとする。

Source: Ministry of Economy, Trade and Industry, Japan. "Handbook for Protection of Confidential Information - Toward Enhancing Corporate Value"
https://www.meti.go.jp/policy/economy/chizai/chiteki/trade-secret.html#handbook

[Ref. 2] Ministry of Education, Culture, Sports, Science and Technology, Japan. "Sakura tool"
https://www.mext.go.jp/a_menu/shinkou/sangaku/1383777.htm
* Provides contract templates that can be used for joint research. The link provides the consortium type of contract and revised materials for individual versions.

14

## In cases of complying data policy

Based on the applicable data policy, set the publication date and time to the end of the embargo period. If the policy does not specify the embargo period, it should be decided in consultation with the department in charge.

## In cases of not available to be published

Even if your research data cannot be published at this time, it is necessary to leave a trail of evidence that the research data exists to support future research activities. Please document the decision-making process up to this point and store the research data in appropriate storage. Also, please publish it as metadata to the institution's platform if possible.

## Q4. Select a data repository

When you have finished checking the constraints, you need to select an appropriate repository. You can publish your research data using file-sharing services or data management software. However, we recommend using data repositories in related fields or institutional repositories from the viewpoints of organization and preservation, data permanence assurance, user recognition, and security management[9]. Some examples of data repositories in Japan are as below. If you want to search for international data repositories, please see the "List of International Data Repositories."

## Disciplinary data repositories

- Social sciences
  - SSJDA (https://ssjda.iss.u-tokyo.ac.jp/Direct/)
    - Acceptable data type: Micro data with questionnaires used in various social and statistical surveys
  - RUDA (https://ruda.rikkyo.ac.jp/dspace/)
    - Acceptable data type: Social survey data (Economics, Business Administration, Sociology, Social Psychology, Political Science, Political Psychology, Law, Sociology of Law, Education, Sociology of Education, etc.)
- Life sciences
  - DDBJ (https://www.ddbj.nig.ac.jp/index.html)
    - Acceptable data type: Annotated/assembled sequences, Sequencing and alignment data from next-generation sequencing platforms, Functional genomics data, Research project, Biological sample, Human data requiring controlled-access
      * Further information:
        https://www.ddbj.nig.ac.jp/data-categories.html
  - NBDC human database (https://humandbs.biosciencedbc.jp/)
    - Acceptable data type: Human data produced from publicly funded research

---

[9] Even when research data is published as an Appendix or Supplement to a research paper, registering it in repositories and databases in related fields will further increase its discoverability and make it more likely to be used. It also facilitates maintenance of broken links.

16

- jPOSTrepo (https://repository.jpostdb.org/)
  - Acceptable data type: ProteOme data in Japan
- GlyTouCan (https://glytoucan.org/)
  - Acceptable data type: Glycan structures data
- Life Science Database Archive (https://dbarchive.biosciencedbc.jp/index.html)
  - Acceptable data type: Datasets generated by domestic life science researchers
- Earth science
  - DIAS (http://www.diasjp.net/)
    - Acceptable data type: Earth and Environmental data
  - IUGONET (http://search.iugonet.org/list.jsp)
    - Acceptable data type: Solar-Terrestrial Science Observations data
  - Global Environment Database (http://db.cger.nies.go.jp/portal/)
    - Acceptable data type: Global Environmental Research Data
- Biological science
  - Biological Information System for Marine Life (BISMaL) (https://www.godac.jamstec.go.jp/bismal/j/)
    - Acceptable data type: Data on marine biotic occurrence records
  - Global Biodiversity Information Facility Japan Node (JBIF) (http://www.gbif.jp/v2/)
    - Acceptable data type: World's biodiversity data
    - Contact: http://www.gbif.jp/v2/regist/index.html
- Synchrotron radiation science
  - SPring-8 case studies & reports cross research (http://www.spring8.or.jp/ja/science/customsearch/)
    - Acceptable data type: Data on polymers, organic thin films, and green energy fields
  - SPring-8 BL14B2 XAFS Standard Sample Database (https://support.spring8.or.jp/xafs/standardDB/standardDB.html)
    - Acceptable data type: XAFS (X-ray absorption fine structure) data

- Other fields
  - List of international data repositories (https://www.re3data.org/)

Institutional repositories

- List of Japanese institutional repositories (https://www.nii.ac.jp/irp/list/)

<u>Considerations</u>

- In the event of data leakage or unauthorized use, you will warn the data reuser and request an injunction against the data reuse. Depending on the circumstances, you may consider injunctive relief, damages, or criminal legal proceedings. Protection may be available under the Copyright Act, or relief may be obtained under the Unfair Competition Prevention Act. You should consult with the department in charge and follow the appropriate procedures.
- When selecting a data repository, information such as whether it conforms to international standards certification and which country's laws it complies with may be helpful.

[Ref. 1] FAIR principles
- FAIR principles as a standard for data sharing (https://doi.org/10.18908/a.2018041901)
- Is the repository listed on the "FAIRsharing" website? (https://fairsharing.org/)
- Is the repository listed on the "Repository Finder," which complies with FAIR principles? (https://repositoryfinder.datacite.org/)

[Ref. 2] International Standards Certification
- From the re3data.org search page (https://www.re3data.org/search), select "Certificates" in the Filter to check.
- A list of recommended repositories may be provided by the publisher.
  Ex. Nature "Scientific Data. Recommended Data Repositories"
  https://www.nature.com/sdata/policies/repositories

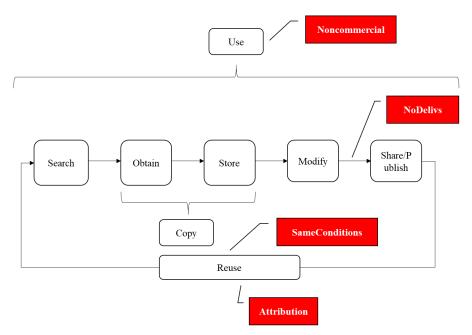## List of Legal Protection of Data in Japan

| | Requirements | | Civil Affairs | | Criminal Affairs | Comparison with Limited Provision Data |
|---|---|---|---|---|---|---|
| | Data to be protected | Misconduct | Demand an Injunction | Claim for Damages | Imprisonment /Fines | |
| Database works (Copyright Act Article 12-2 (1)) | A database that, by reason of the selection or systematic construction of information contained therein, constitutes a creation | Reproduction or any other acts without the permission of the right holder (regardless of the maliciousness of the form) | ○ | | ○ | Data that is not creative (e.g., factory operating data) is not protected |
| Patented invention (Patent Act Article 2 (1), Article 29) | 1) the highly advanced creation of technical ideas utilizing the laws of nature 2) an invention for which a patent has been granted | Implementation or any other acts without the permission of the right holder (regardless of the maliciousness of the form) | ○ | | ○ | |
| Trade secrets (Unfair Competition Prevention Act Article 2 (1) (iv) to (x)) | 1) confidentiality 2) non-public 3) value | Unauthorized acquisition, unauthorized use, etc. (list malicious acts) | ○ | | ○ | Data widely available to the public is not protected |
| Disclosing shared data (Unfair Competition Prevention Act Article 2 (1) (xi) to (xvi)) | 1) shared data with limited access 2) managed by electronic or magnetic means 3) accumulated to a significant extent | Unauthorized acquisition, unauthorized use, etc. (list malicious acts) | ○ | | × | - |
| Torts (Civil Code Article 709) | Data in general | intentionally or negligently infringed the rights or legally protected interests | × (Except in Moral rights infringement) | ○ | × | No injunction is available (in principle) |
| Contract (Non-Performance) (Civil Code Article 415) | Data in general (depends on contract) | Violation of contract | ○ (only contracting parties) | | × | Cannot be applied to other than contracting parties |

Ref. Ministry of Economy, Trade and Industry, Japan. Major legal systems against data misuse

https://www.meti.go.jp/policy/economy/chizai/chiteki/H30nen_fukyohoshosai.pdf

## Q5. Choose appropriate conditions of use

Once you have decided where to publish the data, you need to specify the conditions of use for your research data. When you have multiple datasets or are publishing them together with derived data, it is convenient to specify their conditions of use together. Specified conditions of use should be appropriately described as metadata[10] when registering to the repository.

Flow of research data reuse and actions subject to recommended conditions of use specification



- Please clearly indicate the following four points: attribution, same conditions, NoDelivs, and Noncommercial.
- "NoDelivs" means that the release of modified data is prohibited.

---

[10] To ensure transparency, we strongly recommend that a description of how data is acquired and modified be created in the metadata according to the conditions of use or a link be provided to a report or data paper describing the procedure.

Recommended conditions of use

In specifying the conditions of use, you can combine some conditions such as attribution and/or permission for modification. You can also choose to waive your rights. The description in square brackets for each conditional specification can be used when describing the conditions of use.

List of conditions of use combinations

| Conditions of use | Description |
|---|---|
| Waiver | Freely available |
| Attribution | Clearly indicate the data source and credit information |
| Attribution - Noncommercial | Clearly indicate the data source and credit information; Commercial use prohibited |
| Attribution - NoDelivs | Clearly indicate the data source and credit information; the release of modified data prohibited |
| Attribution - Noncommercial - NoDelivs | Clearly indicate the data source and credit information; Commercial use prohibited; the release of modified data prohibited |
| Attribution - SameConditions | Clearly indicate the data source and credit information; the release of modified data granted different conditions of use prohibited |
| Attribution - SameConditions - Noncommercial | Clearly indicate the data source and credit information; the release of modified data granted different conditions of use prohibited; Commercial use prohibited |
| Other | Individual restrictions by contract (e.g., limited sharing) |

- Although these guidelines aim to set appropriate conditions of use for a non-copyrighted data, there are many cases in which it is difficult to determine whether or not a copyrighted work. The conditions of use recommended in this section are compatible with the Creative Commons License (https://creativecommons.jp/licenses/) International 4.0, and you can use these conditions of use regardless of their copyrightability.
- Note that when research data are based on derived data (see p. 5), it is impossible to grant fewer conditions of use than the source data. For example, you cannot grant

"Waiver" if the original conditions of use are "Attribution," even if the research data was published by yourself.

- "Waiver" can be understood as a declaration by the data provider that they waive the right to take legal action against copyright infringement. However, even if the data provider declares a "waiver," moral rights, privacy rights, and the right to prevent unfair competition remain. Trademark and patent rights are not also waived. Therefore, legal action against rights infringement other than copyright can be considered by the data provider, the institution to which the data belongs, or the repository manager. Cf. https://creativecommons.jp/sciencecommons/aboutcc0/

- "NoDelivs" prohibits the sharing or publication of modified data, so it does not prohibit private data reuse. Also, even under conditions of use that do not grant "NoDelivs," the Unfair Competition Prevention Act may be applied in case of data falsification.

1）Waiver

- The data are freely available for commercial or non-commercial purposes. It is unnecessary to indicate the data source information or the modification methods.

---

[Description (Japanese)]

- 本データの利用に当たり、原則として、何らの制約はありません。

※本データが著作物である場合は、CC0（権利放棄）が付与されます。

---

[Notice]

- If you select a "waiver," you may not revoke or change your choice in the future. Please check carefully with your institution or professional before selection.
- This conditions of use will be interpreted that the data reuser does not require third parties to indicate the data source information or the same conditions of use as the original data.
- Even if the data provider has waived their rights, it may be necessary to indicate the data source information according to the someone's policy, such as journal policy. You need to cite data appropriately, taking into account research ethics regulations.

2) Attribution

• The data are freely available as long as the data source and credit information are clearly indicated.

[Description (Japanese)]
- ➤ 本データの公開に当たっては、出所を明示してください。
- ➤ 本データを改変した場合には、その手順を何らかの手段で明記してください。

※本データに著作権が発生する場合、クリエイティブ・コモンズ 表示 4.0 国際ライセンス（CC-BY）の条件で利用することが可能です。著作権が発生しない場合でも、出所の明示を条件に利用することが可能です。

[Notice]
• In displaying attribution for your research data, clearly indicate the credit information, including the version and date/time information on the landing page.
• There are several ways to specify the means of data modification: 1) mentioning when the data source is referred, 2) including it in the metadata, and 3) writing a report or data paper that describes the procedure in more detail. Choose a method appropriate to the degree of alteration.

3) Attribution - Noncommercial

- For noncommercial purposes, the data are freely available as long as the data source and credit information are clearly indicated.

[Description (Japanese)]
- ➢ 本データの公開に当たっては、出所を明示してください。
- ➢ 本データは、営利目的で利用することができません。
- ➢ 本データを改変した場合には、その手順を何らかの手段で明記してください。

※本データに著作権が発生する場合、クリエイティブ・コモンズ 表示-非営利 4.0 国際ライセンス（CC-BY-NC）の条件で利用することが可能です。著作権が発生しない場合でも、出所の明示及び非営利目的での利用を条件に利用することが可能です。

4) Attribution - NoDelivs

* The data are freely available as long as the data source and credit information are clearly indicated; The release of modified data is prohibited.

[Description (Japanese)]
> 本データの公開に当たっては、出所を明示してください。
> 本データを改変した場合、改変されたデータを公開することはできません。

※本データに著作権が発生する場合、クリエイティブ・コモンズ 表示-改変禁止 4.0 国際ライセンス（CC-BY-ND）の条件で利用することが可能です。著作権が発生しない場合でも、出所の明示及び改変されたデータを公開しないことを条件に利用することが可能です。

[Notice]
* It is common for data acquired by third parties to be modified in the reuse process, except when the data is only for observation, viewing, or browsing. If you have a particular modification method that you want to prohibit, clearly state it.
* The term "Modified" includes partially altered from the source data and derived data. For example, this applies when tabular data are prepared based on individual data or estimates are organized based on observed data.

5) Attribution - Noncommercial - NoDelivs

• For noncommercial purposes, the data are freely available as long as the data source and credit information are clearly indicated; The release of modified data is prohibited.

[Description (Japanese)]
  ➢ 本データの公開に当たっては、出所を明示してください。
  ➢ 本データは、営利目的で利用することができません。
  ➢ 本データを改変した場合、改変されたデータを公開することはできません。

※本データに著作権が発生する場合、クリエイティブ・コモンズ 表示-非営利-改変禁止 4.0 国際ライセンス（CC-BY-NC-ND）の条件で利用することが可能です。著作権が発生しない場合でも、出所の明示、非営利目的での利用及び改変されたデータを公開しないことを条件に利用することが可能です。

6) Attribution - SameConditions

- The data are freely available as long as the data source and credit information are clearly indicated; The release of modified data granted different conditions of use is prohibited.

[Description (Japanese)]
  ➢ 本データの公開に当たっては、出所を明示してください。
  ➢ 本データを改変した場合には、本データと同じ利用条件で公開し、かつその手順を何らかの手段で明記してください。

※本データに著作権が発生する場合、クリエイティブ・コモンズ 表示-継承 4.0 国際ライセンス（CC-BY-SA）の条件で利用することが可能です。著作権が発生しない場合でも、出所の明示及び元データと同じ利用条件要素を付与することを条件に利用することが可能です。

7）Attribution - SameConditions - Noncommercial

- For noncommercial purposes, the data are freely available as long as the data source and credit information are clearly indicated; The release of modified data granted different conditions of use is prohibited.

[Description (Japanese)]
   ➢ 本データの公開に当たっては、出所を明示してください。
   ➢ 本データは、営利目的で利用することができません。
   ➢ 本データを改変した場合には、本データと同じ利用条件で公開し、かつその手順を何らかの手段で明記してください。

※本データに著作権が発生する場合、クリエイティブ・コモンズ 表示-継承-非営利 4.0 国際ライセンス（CC-BY-SA-NC）の条件で利用することが可能です。著作権が発生しない場合でも、出所の明示、非営利目的での利用及び元データと同じ利用条件要素を付与することを条件に利用することが可能です。

8） Additional conditions of use

In cases of some additional conditions of use caused by an institution's data policy or individual contracts, the detailed description will be more readily understood when displayed along with the usage notes. We recommend preparing a concise user guide with references to these policies or contracts. Note that any new conditions granted to the data will no longer make it compatible with the Creative Commons License.

## Appendix. Terms of use statement

Since the research data assumed by these guidelines are not protected by copyright law, it is necessary to indicate the specified terms of use in the metadata and establish more detailed terms in advance to legally guarantee the specified conditions of use. Referring to the following sample format, please check your description, such as the data source and credit information, an example of how to modify the data, and a disclaimer. If there is missing information in the landing page, add the necessary information to the metadata and consider changing the data repository if necessary.

- Sample format: in cases for "Attribution" (Japanese)

　本データ及び付録資料に収録された情報(以下「本データ等」といいます)に関する一切の権利は、原則として、本データ等の作成に関与した研究者、研究機関又は当該データの提供者 (以下「情報提供者」といいます) に帰属します。本データ等に関する権利は、我が国国内法及び国際条約により保護されており、情報提供者が指定する利用規約又はライセンス表示に従う場合を除いて、本データ等を無断で利用することはできません（使用、複製、頒布、上映、公衆送信、上演、出版、送信可能化、翻案、改変及び商用利用を含みますが、これらに限られません）。本データ等の利用に当たっては、情報提供者が指定する利用規約又はライセンス表示に同意したものとみなします。

　（利用条件）
- 本データの公開に当たっては、出所を明示してください。
- 本データを改変した場合には、その手順を何らかの手段で明記してください。

　※本データ等に著作権が発生する場合、クリエイティブ・コモンズ 表示 4.0 国際ライセンス（CC-BY）の条件で利用することが可能です。著作権が発生しない場合でも、出所の明示を条件に利用することが可能です。

29

（例 1）出所：「本データ等の名称」（本データ等の作者名）（本デー
タ等の URL）（バージョン表記などの日時情報）

（例 2）出所：「本データ等の名称」（本データ等の作者名）（本デー
タ等の URL）をもとに（利用者名）が加工して作成

　なお、本データ等に関しては、万全を期してはおりますが、正確性、確実性、目的適合性その他の品質を保証するものではありません。本データ等を用いて行うすべての行為に関して、その責任はすべて利用者自身に帰属します。

　万が一、本データ等を用いたことによって利用者が何らかの損害を被った場合、その損害に関して情報提供者は一切の責任を負うものではありません。得られた情報に基づく決定は、本データ等の利用者ご自身でご判断いただきますようお願い申し上げます。

　また、情報提供者は本データ等からアクセス可能な、第三者が権利を有する情報の正確性、信頼性、安全性を何ら保証するものではなく、第三者が権利を有する情報の利用により生じたいかなる損害に関しても、情報提供者は一切の責任を負うものではありません。

　本データ等は、予告なく追加、変更、削除されることがありますので、あらかじめご了承ください。

（氏名）

（文書の公開年月日）

References

1) Shu Higashi. "Considering open data licenses (6)" (translated from Japanese) Open Data Commons, 2012.
http://okfn.jp/2012/10/24/opendata_license06/, (accessed 2019-12-20).

2) Masafumi Ono, Toshio Koike, Ryosuke Shibasaki. Survey for research data sharing in earth environmental information domain: Realities in research community. Journal of Information Processing and Management. 2016, Vol. 59, No. 8, p. 514-525.
https://doi.org/10.1241/johokanri.59.514, (accessed 2019-12-20).

3) Cabinet Office, Japan. "Government of Japan Standard Terms of Use (Version 2.0)". 2015. https://www.kantei.go.jp/jp/singi/it2/densi/kettei/gl2_betten_1.pdf, (accessed 2019-12-20).

4) Liaison Committee of Relevant Ministries and Agencies and Working-Level Consultative Meeting on Collaboration of Digital Archives. "Guidelines for Building, Sharing, and Using Digital Archives" (translated from Japanese). 2017.
https://www.kantei.go.jp/jp/singi/titeki2/digitalarchive_kyougikai/guideline.pdf, (accessed 2019-12-20).

5) Shinnosuke Fukuoka and Hidetoshi Matsumura. Data laws and contracts. SHOJIHOMU. 2019, 440p.

6) Alex Ball. How to License Research Data. 2011. http://www.dcc.ac.uk/resources/how-guides/license-research-data, (accessed 2019-12-20).

7) FORCE11: Data Citation Synthesis Group. Joint Declaration of Data Citation Principles. 2014. https://doi.org/10.25490/a97f-egyk, (accessed 2019-12-20).
（Research Data Utilization Forum (RDUF): Research Data Citation Subcommittee, Trans. https://doi.org/10.11502/rduf_rdc_jddcp_ja, (accessed 2019-12-20).)

8) RDA-CODATA: Legal Interoperability Interest Group. Legal Interoperability of Research Data. Principles and Implementation Guidelines. Zenodo, 2016.
https://doi.org/10.5281/zenodo.162241, (accessed 2019-12-20).

* All URLs in the text and footnotes are as of December 20, 2019.

Members (ver.1, as of December 2019)

Yasuyuki Minamiyama (National Institute of Informatics)

Ui Ikeuchi (Bunkyo University)

Kunihiko Ueshima (Japan Data Exchange, Inc.)

Misaki Suto (Formerly Mitsubishi UFJ Research and Consulting Co., Ltd.)

Nobuya Okayama (Hitachi Consulting)

Issaku Yamada (The Noguchi Institute)

Ken Ebisawa (Japan Aerospace Exploration Agency, Institute of Space and Astronautical Science)

Hodaka Nakanishi (Teikyo University)

Yui Kumazaki (Japan Atomic Energy Agency)

# Glossary of "Guideline for specifying conditions of use

# in research data publishing"

| No. | Term | Definition | Source/Reference |
|---|---|---|---|
| 1 | Falsification | Manipulating research materials, equipment, or processes to change data or results obtained from research activities. | Ministry of Education, Culture, Sports, Science and Technology, Japan. "Guidelines for Responding to Misconduct in Research" https://www.mext.go.jp/b_menu/shingi/ gijyutu/gijyutu12/houkoku/attach/13346 60.htm, (accessed 2019-12-25). |
| 2 | Creative Commons | The name of a project or a non-profit organization that promotes the smooth distribution and reuse of copyrighted works by providing copyright holders with a means of indicating the conditions of use of their works with a simple mark. By declaring the conditions of use of texts, photos, videos, sounds, etc. on websites, etc., using the Creative Commons-defined marks in advance, the copyright holder can save users from having to go through the licensing procedure. | Japan Society of Library and Information Science, Dictionary of Terms Editorial Committee. Dictionary of Library and Information Science Terms. 4th edition. 2014. |
| 3 | Credit | Formal recognition of the contributions made by an individual or group to the research outputs. | RDA-CODATA: Legal Interoperability Interest Group. "Legal Interoperability of Research Data. Principles and Implementation Guidelines". Zenodo, 2016. https://doi.org/10.5281/zenodo.162241, (accessed 2019-12-25). |
| 4 | Research data | Digital data used as a source of information for scientific research. It includes a variety of formats, such as numerical, textual, image, audio, and video. In these guidelines, it does not include physical materials such as samples (specimens, samples) or recording media (paper, disks, etc.). | 1) Cabinet Office, Japan. "Report of the Working Group on Research Data Infrastructure and Global Outreach - Strategies for the Development of Research Data Infrastructure and Global Outreach". 2019. https://www8.cao.go.jp/cstp/tyousak |

| | | | ai/kokusaiopen/houkokusho.pdf, (accessed 2019-12-25).<br>2) Japan Science and Technology Agency. "JST Policy on Open Access to Research Publications and Research Data Management". 2017. https://www.jst.go.jp/pr/intro/openscience/policy_openscience.pdf, (accessed 2019-12-25).<br>3) OECD. OECD principles and guidelines for access to research data from public funding. 2007. https://doi.org/10.1787/9789264034020-en-fr, (accessed 2019-12-25). |
|---|---|---|---|
| 5 | Industrial Property rights | Four of the intellectual property rights, patent rights, utility model rights, design rights, and trademark rights, are referred to as industrial property rights. The purpose of a system of industrial property rights is to encourage and to motivate inventors of inventions and creators of designs, to protect their rights, and to instill confidence in the maintenance of business activities related to trademarks. | Japan Patent Office. "System of Industrial Property Rights". https://www.jpo.go.jp/system/patent/gaiyo/seidogaiyo/chizai01.html, (accessed 2019-12-25). |
| 6 | Sample | A physical object that has substance, such as a specimen to be used for research. | Science Council of Japan. "Response: Improving Soundness in Scientific Research" (in Japanese). 2015. http://www.scj.go.jp/ja/info/kohyo/pdf/kohyo-23-k150306.pdf, (accessed 2019-12-25). |
| 7 | Information disclosure | A person with information shows it to another. This document distinguishes between "access to information," which is done by government agencies and academic institutions for the public. In other words, information disclosure may have some restrictions, such as limiting the scope and subject matter of disclosure, confidentiality | 1) Ministry of Internal Affairs and Communications, Japan. "Access to Information System" (in Japanese). https://www.soumu.go.jp/main_sosiki/gyoukan/kanri/jyohokokai/index.html, (accessed 2019-12-25).<br>2) Strike Co., Ltd. "Glossary of M&A" |

| | | obligations, and payment of compensation. When the transfer, lease, or assignment of rights to that information is involved, it is called "provision of information." | https://www.strike.co.jp/maword/0331.html, (accessed 2019-12-25).<br>3) Cambridge Dictionary. https://dictionary.cambridge.org/ja/dictionary/english/disclosure, (accessed 2019-12-25). |
|---|---|---|---|
| 8 | Storage | External memory, one of the main devices that make up a computer, is used to store data for an unspecified period. This term is used to maintain its contents even when electricity is not supplied, such as hard disks, optical disks (CDs and DVDs), flash memory storage devices (USB memory sticks and memory cards), magnetic tapes, and so on. | 1) Online Dictionary for Library and Information Science, https://www.abc-clio.com/ODLIS/odlis_s.aspx, (accessed 2019-12-25).<br>2) IDC Frontier Inc. "Glossary of Cloud / Data Center". https://www.idcf.jp/words/storage.html, (accessed 2019-12-25). |
| 9 | Data sharing | Providing or disclosing data only to a limited number of subjects. "Data sharing" term is sometimes used in the sense of data publishing; these guidelines distinguish it from data publishing. | Cabinet Office, Japan. "Report of the Working Group on Research Data Infrastructure and Global Outreach - Strategies for the Development of Research Data Infrastructure and Global Outreach". 2019. https://www8.cao.go.jp/cstp/tyousakai/kokusaiopen/houkokusho.pdf, (accessed 2019-12-25). |
| 10 | Data publishing | Publishing data to third parties on websites, repositories, or supplements to research papers and is accessible via the Internet. | 1) Cabinet Office, Japan. "Report of the Working Group on Research Data Infrastructure and Global Outreach - Strategies for the Development of Research Data Infrastructure and Global Outreach". 2019. https://www8.cao.go.jp/cstp/tyousakai/kokusaiopen/houkokusho.pdf, (accessed 2019-12-25).<br>2) G7 Science and Technology Ministers' Meeting. "Tsukuba Communiqué". 2016. |

| | | | |
|---|---|---|---|
| | | | https://www8.cao.go.jp/cstp/kokusai teki/g7_2016/2016communique.html , (accessed 2019-12-25). |
| 11 | Data paper | A paper that describes the content, acquisition method, data format, access information, etc., regarding publicly available data such as observation data, measurement data, analysis data, and calculation simulation results. It does not include analysis, interpretation, or scientific conclusions. | 1) Vishwas Chavan, Lyubomir Penev. The data paper: a mechanism to incentivize data publishing in biodiversity science. BMC Bioinformatics Vol.12, S2, 2011. https://doi.org/10.1186/1471-2105-12-S15-S2, (accessed 2019-12-25). <br> 2) Introduction of "Data Paper": New Category for JAMSTEC-R Article, JAMSTEC Report of Research and Development, 2017, Vol. 24, P. 21-22. 2017, https://doi.org/10.5918/jamstecr.24. 21, (accessed 2019-12-25). |
| 12 | Data policy | A statement of the data and information management processes that the organization has designed to support and protect the organization's research data assets. It is a set of high-level principles that establishes a guiding framework for data management. Data policies can be used to address strategic aspects such as data access, relevant legal matters, data management issues and storage operations, data acquisition, and other issues. | Research Data Canada and CASRAI. Trans-Disciplinary Glossary for Research Data Management. https://dictionary.casrai.org/Data_policy, (accessed 2019-12-25). |
| 13 | Anonymizati on | A manipulation that reduces the risk of personal identification by processing information that could directly/indirectly identify an individual. It is distinguished from "pseudonymization"; Pseudonymization is an operation that removes or separates information that can directly identify an individual by itself (e.g., name, mug shot, fingerprints, driver's license number, etc.) from other information. Anonymization involves | 1) Personal Information Protection Commission, Government of Japan. "Guidelines for the Act on the Protection of Personal Information. Anonymized Information section". 2017. https://www.ppc.go.jp/files/pdf/rep ort_office.pdf, (accessed 2019-12-25). |

4

| | | disambiguating or replacing age, gender, occupation, behavior logs, etc., so normal methods cannot recover the original information. | 2) Jun Sakuma. Privacy Protection in Data Analysis: Machine Learning Professional Series. 2016. |
|---|---|---|---|
| 14 | Plagiarism | The appropriation of another person's ideas, processes, results, or words without giving appropriate credit. | U.S. Department of Health and Human Service, Office of Research Integrity. Definition of Research Misconduct. 2000. https://ori.hhs.gov/definition-misconduct, (accessed 2019-12-25). |
| 15 | Metadata | Data describing the characteristics of an information resource in order to effectively identify, describe, and explore it. Every academic community has its unique metadata tied to its information use practices, which are numerous and varied. | Japan Society of Library and Information Science, Dictionary of Terms Editorial Committee. Dictionary of Library and Information Science Terms. 4th edition. 2014. |
| 16 | License | To grant official permission for any action, use or possession. It also means an official document that sets forth the terms and conditions for a patented invention or the right to use intellectual property such as software. The Creative Commons License and the MIT License are well known in the academic community. In this document, the basic declaration of intent by the rights holder is called a "license indication" and is distinguished from the detailed "conditions of use" set forth in contracts or terms of use. | Cambridge Dictionary. https://dictionary.cambridge.org/dictionary/english/licence, (accessed 2019-12-25). |
| 17 | Landing page | In a broad sense, it refers to the first page that a visitor sees through a website through a link or advertisement. It is distinguished from the top page, which is the entrance to a website. In the academic community, it is a page that describes the metadata and access methods for digital content in a repository. | 1) NTTCom Online Marketing Solutions Corporation. Visionalist, "Glossary of Web Marketing". https://www.visionalist.com/glossaries/10_ra_001.html, (accessed 2019-12-25). 2) Basic Inc. ferret, "Web Marketing Dictionary". https://ferret-plus.com/words/1048, (accessed 2019-12-25). 3) Cambridge Dictionary. |

| | | | | |
|---|---|---|---|---|
| | | | | https://dictionary.cambridge.org/dictionary/english/landing-page, (accessed 2019-12-25). |
| | | | 4) | Japan Link Center Joint Steering Committee. "Guidelines for Registering DOIs for Research Data". https://doi.org/10.11502/rd_guideline_ja, (accessed 2019-12-25). |
| 18 | Repository (Digital repository) | Information systems that capture, store, manage, preserve, and provide access to digital content. Depending on the managing entity, there are different names for institutional repositories, discipline-specific repositories, government repositories, etc. The definition category also differs depending on the purpose of repository construction. For example, an open access repository is defined as a collection of full-text documents available in an online database on the Internet and characterized by free and immediate access. | 1) | Iris Xie, Krystyna K. Matusiak, Chapter 1 - Introduction to digital libraries, Discover Digital Libraries. 2016. https://doi.org/10.1016/B978-0-12-417112-1.00001-6, (accessed 2019-12-25). |
| | | | 2) | Iris Xie, Krystyna K. Matusiak, Chapter 9 - Digital preservation, Discover Digital Libraries. 2016. https://doi.org/10.1016/B978-0-12-417112-1.00009-0, (accessed 2019-12-25). |