



# FOUNDATIONS of DATA CURATION

Allen H. Renear, Cheryl A Thompson, Katrina S Fenlon, Myrna Morales



School of Information Sciences



University of Illinois at Urbana-Champaign



# METADATA





4

HOW DOES METADATA  
SUPPORT DATA CURATION?

# How Does Metadata Support Data Curation?

Relating metadata to curation *objectives, activities, actions*

# Recall the data curation objective

*Data curation is concerned with all aspects of the management of data in order to efficiently and reliably support the analysis of data, and enable reuse over time*

.



# ***Recall: Areas of curatorial activities***

<b>Collection:</b>	Support the collection and acquisition of data
<b>Organization:</b>	Employ an appropriate data model and use appropriate standards
<b>Storage:</b>	Support reliable and effective storage
<b>Preservation:</b>	Ensure that data will be understandable and useable in the future
<b>Discoverability:</b>	Support the ability to search for and locate relevant data
<b>Access:</b>	Support the ability to retrieve and distribute data
<b>Workflow:</b>	Support the ability to systematize data workflows
<b>Identification:</b>	Support the ability to identify, authenticate, and validate data
<b>Integration:</b>	Support integration of data from different sources using different data models
<b>Reformatting:</b>	Support reformatting for use by different tools or to match new format standards
<b>Reproducibility:</b>	Support ability to reproduce results, ensuring scientific validity
<b>Sharing:</b>	Support sharing data between researchers, teams, and institutions.
<b>Communication:</b>	Support representation, publishing, and visualizations that provide insight
<b>Provenance:</b>	Support identifying what inputs and calculations are responsible for data values
<b>Modification:</b>	Support management of corrections and updates
<b>Compliance:</b>	Ensure compliance to legal, regulatory, and local policy requirements
<b>Security:</b>	Ensure that data is secure from tampering or inappropriate access and distribution

# *Recall:* Methods of curatorial action

## Analysis

To determine needs, and develop relevant data models and *metadata*, and reformat, correct, or update data.

## → Documentation

To record essential information (typically via *metadata*)

## System design and implementation

To support all data curatorial activities

To support the generation and use of data documentation and processing documentation

## Policy

To specify objectives, procedures, practices, and formats.

## Process

To ensure success and efficiency by managing the development of appropriate organizational units and roles, providing training, advocating for change, and managing curatorial activities.

→ Metadata is rigorously defined machine-processable documentation

# Areas of curatorial activities (Part I: the core)

Area	Description	Supported by metadata that documents. . .
<b>Collection</b>	Support the collection and acquisition of data	Method, location, time, instruments, settings, calibration. . .
<b>Organization</b>	Employ an appropriate data model and use appropriate standards	Schemas and schema documentation for semantics, syntax, and encoding.
<b>Storage</b>	Support reliable and effective storage	Authoritative and alternative copies, physical locations, redundancy, compression, reduction, backups
<b>Preservation</b>	Ensure that data will be understandable and useable in the future	[see Organization, Identification, Storage]
<b>Discoverability</b>	Support the ability to search for and locate relevant data	Topic, coverage, formats, availability, currency.
<b>Access</b>	Support the ability to retrieve and distribute data	[see Organization, Security], licensing, owner, location.
<b>Workflow</b>	Support the ability to systematize data workflows	[See Organization], scripts, processes, transformations, inputs.
<b>Identification</b>	Support the ability to identify, authenticate, and validate data	[See Organization], identifiers, version data, integrity checks, authentication.
<b>Integration</b>	Support integration of data from different sources using different data models	[See Organization, Identification, Access, Discoverability]



# Areas Of Curatorial Activities (Part II: Partial Dependencies)

Activity	Description	Metadata documenting . . .
<b>Modification</b>	Support management of corrections and updates	[See Organization, Workflow, Provenance, Identification]
<b>Reformatting</b>	Support reformatting for use by different tools or to match new format standards	[See Organization, Identification, Workflow]
<b>Provenance</b>	Support identifying what inputs and calculations are responsible for data values	[See Workflow, Identification, Reproducibility]
<b>Reproducibility</b>	Support ability to reproduce results, ensuring scientific validity	[See Organization, Workflow, Provenance, Identification]
<b>Preservation</b>	Ensure that data will be understandable and useable in the future	[see Organization, Identification, Storage]
<b>Compliance</b>	Ensure compliance to legal, regulatory, and local policy requirements	[see Organization, Provenance, Workflow, Discoverability]. Certification.
<b>Security</b>	Ensure that data is secure from tampering or inappropriate access and distribution	[see Organization, Provenance, Workflow, Discoverability]. Encryption, Certification.
<b>Communication</b>	Support representation, publishing, and visualizations that provide insight	[See Organization, Identification, Reproducibility, Compliance]
<b>Sharing</b>	Support sharing data between researchers, teams, and institutions.	[See Discoverability, Organization, Workflow, Provenance, Identification]

# *Frictions* In Creating And Managing Metadata

Metadata Friction Categories	Specific frictions
Standardization	<ul style="list-style-type: none"><li>• Unfinished and multiple metadata versions</li><li>• Proliferation of metadata standards</li></ul>
Temporal	<ul style="list-style-type: none"><li>• Metadata in the life cycle</li><li>• Timeframe for metadata use</li></ul>
Data Sharing	<ul style="list-style-type: none"><li>• Audience for metadata</li><li>• Individual vs. group knowledge</li></ul>
Human Support	<ul style="list-style-type: none"><li>• Metadata is nobody's job</li><li>• Proliferation of metadata tools</li></ul>

*Metadata frictions* specific to the “problems and conflicts that must be addressed to make metadata useful”





# FOUNDATIONS OF DATA CURATION (IS531)

**Allen H. Renear, Cheryl A Thompson, Katrina S Fenlon, Myrna Morales**

**School of Information Sciences**

**University of Illinois at Urbana-Champaign**

**Includes material adapted from work by Carole Palmer, Melissa Cragin,  
David Dubin, Karen Wickett, Bertram Ludæscher, Ruth Duerr and Simone Sacchi.**

**Comments and corrections to: [renear@illinois.edu](mailto:renear@illinois.edu).**