

Documentation and metadata.
James Bukosia Wafula - KWTRP

June 25, 2019

Course content

- Documentation for reproducible research.
- Categories of documentation in RDM.
- When & what to document.
- Data management process flow documentation.
- Metadata: What is this?
- Functions of metadata.
- Types/categories of metadata.
- Metadata standards.
- Exposing metadata via harvesting.

Documentation

- Documentation ensures that anyone re-using the data can understand it and interpret it correctly.
- Explains how the data was created, context for the data, structure of the data and its contents, and any manipulations that have been done to the data. E.g: Segmentation manipulations on DICOM files, manifold learning techniques employed for NLP and other pattern recognition applications.
- Documentation ensures that data can be searched for and retrieved efficiently by users of data centres and repositories.

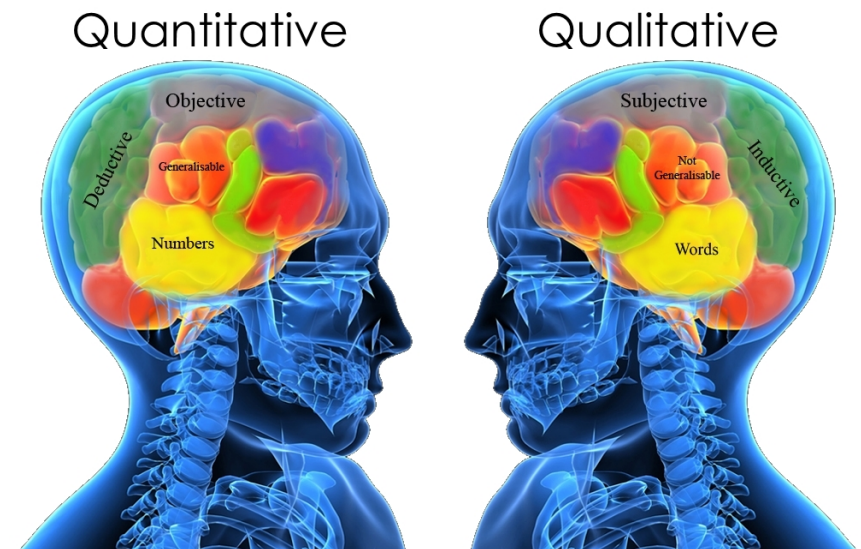
Documentation for reproducible research

- Documentation is not an option especially now that reproducible research is a key requirement from funders.
- Other researchers will be interested in reproducing results with minimal effort.
- Hence data documentation should be shared so that someone else can follow and replicate your results. E.g: Set random seeds for simulated data.



Documentation: Qualitative data

- Qualitative data is usually not generalizable.
- Documentation will make it easier for other qualitative researchers to understand your data.
- Makes research output visible and increases re-use potential, even when research results are not reproducible.



Categories of documentation in RDM

Documentation can be categorized into two:

- **Internal (embedded) documentation:** Which is written as program comments. It includes:
 - Code, field and label descriptions.
 - Descriptive headers or summaries.
 - Transcripts.
- **External (supporting) documentation:** Written for people who need to use the software, in a separate file. This can further be split into:
 - **Library documentation:** Which describes tools that a programmer can use.
 - **User documentation:** Which is intended for users of an application e.g. codebooks, questionnaires or interview guides.

Data documentation in RDM life cycle

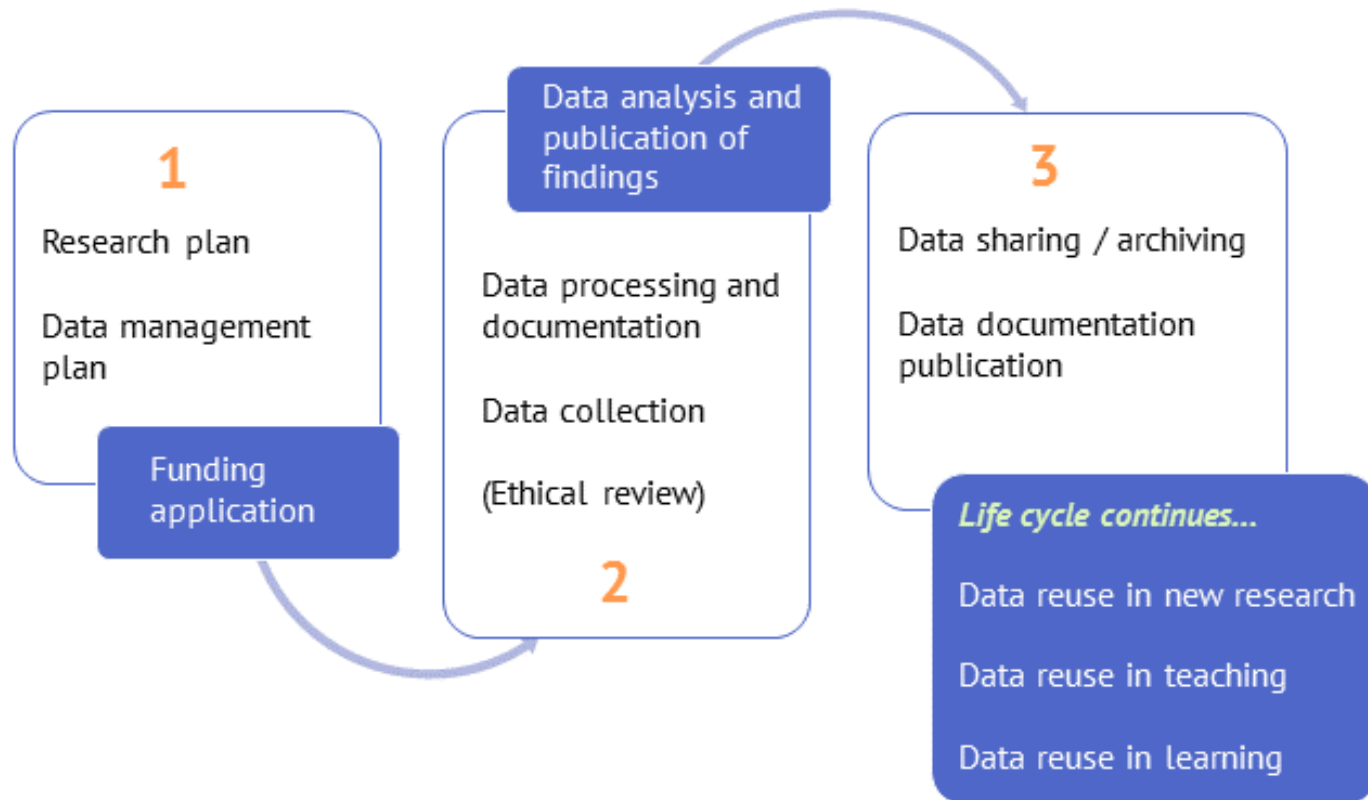


Figure: Data life cycle

What to document

It is vital to document the study for which the data has been collected and the data itself. Therefore, we can split documentation into:

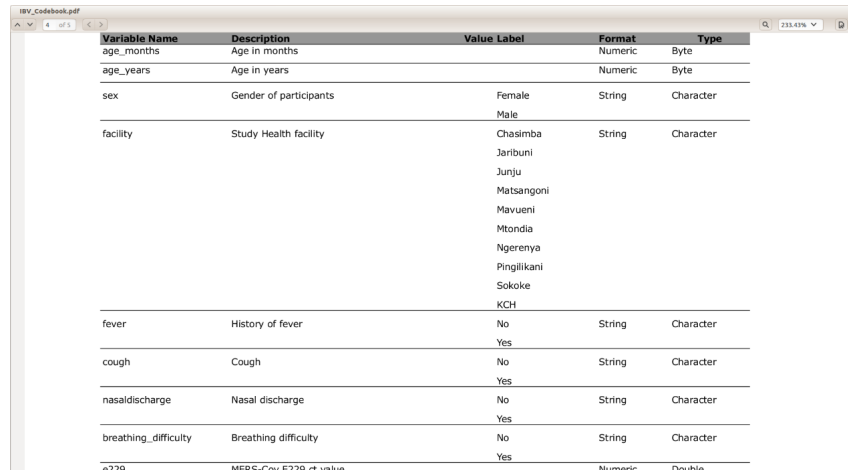
- Project-level documentation.
- Data-level documentation.

Project-level documentation

- Purpose for which data was collected.
- Contents of the dataset.
- Data collection method employed.
- Who collected the data.
- Data processing techniques.
- Manipulations/ transformations done on data.
- Quality assurance.
- Data access.

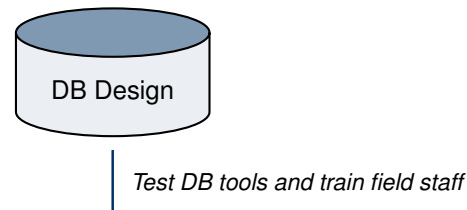
Data-level (object-level) documentation

- Provides information at the level of individual objects e.g. interview transcripts or variables in a database.
- Can be embedded in datafiles e.g. descriptive information about an interview at the beginning of each file; or variable names embedded in each file for quantitative data.

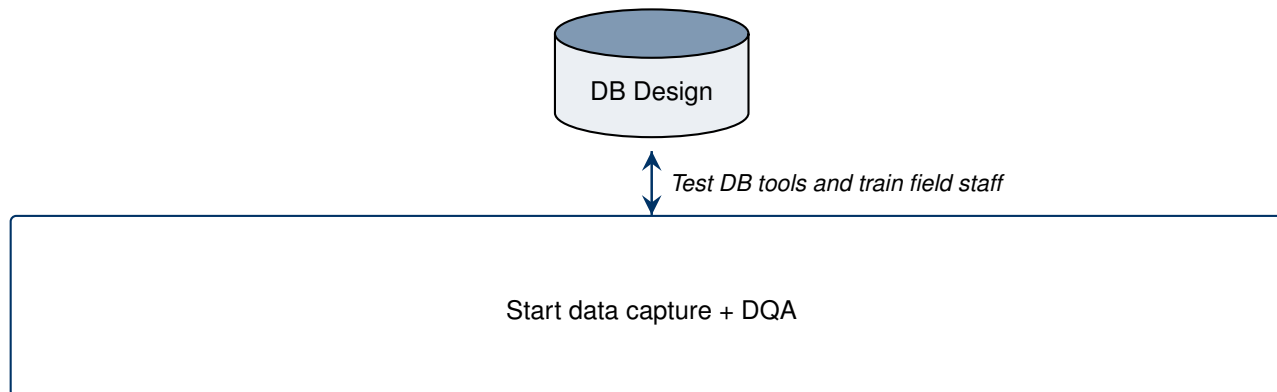


Variable Name	Description	Value Label	Format	Type
age_months	Age in months		Numeric	Byte
age_years	Age in years		Numeric	Byte
sex	Gender of participants	Female Male	String	Character
facility	Study Health facility	Chasimba Jaribuni Junju Matsangoni Mavueni Mtondia Ngerenya Pingilikani Soko KCH	String	Character
fever	History of fever	No Yes	String	Character
cough	Cough	No Yes	String	Character
nasaldischarge	Nasal discharge	No Yes	String	Character
breathing_difficulty	Breathing difficulty	No Yes	String	Character
e229	MERS-Cov E229 ct value		Numeric	Double

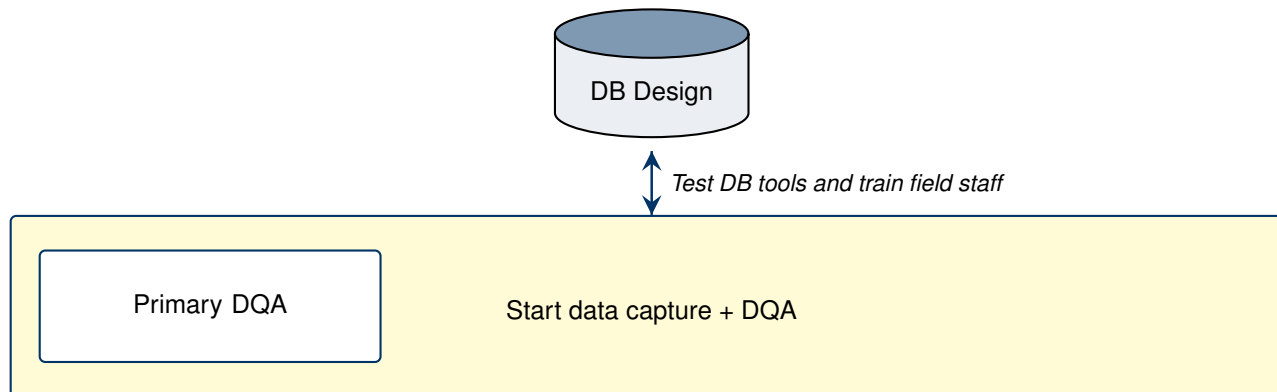
Data management process flow documentation example



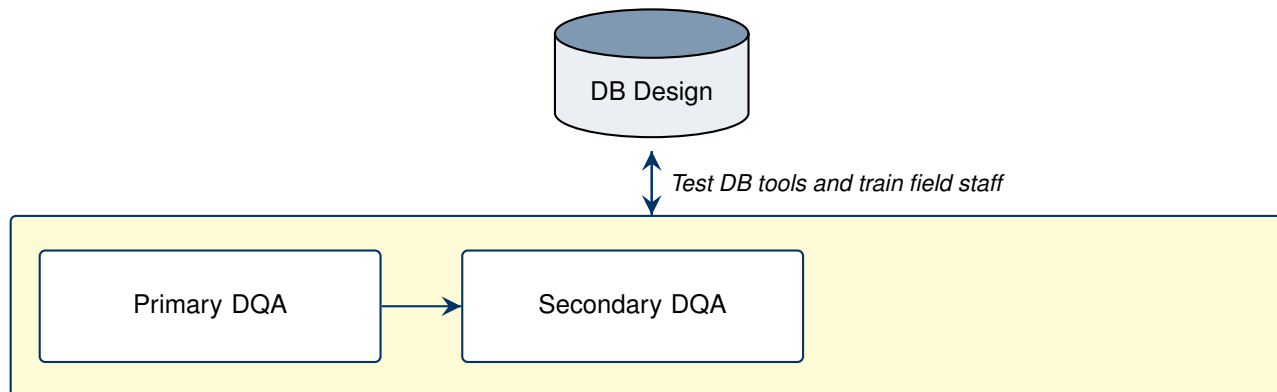
Data management process flow documentation example



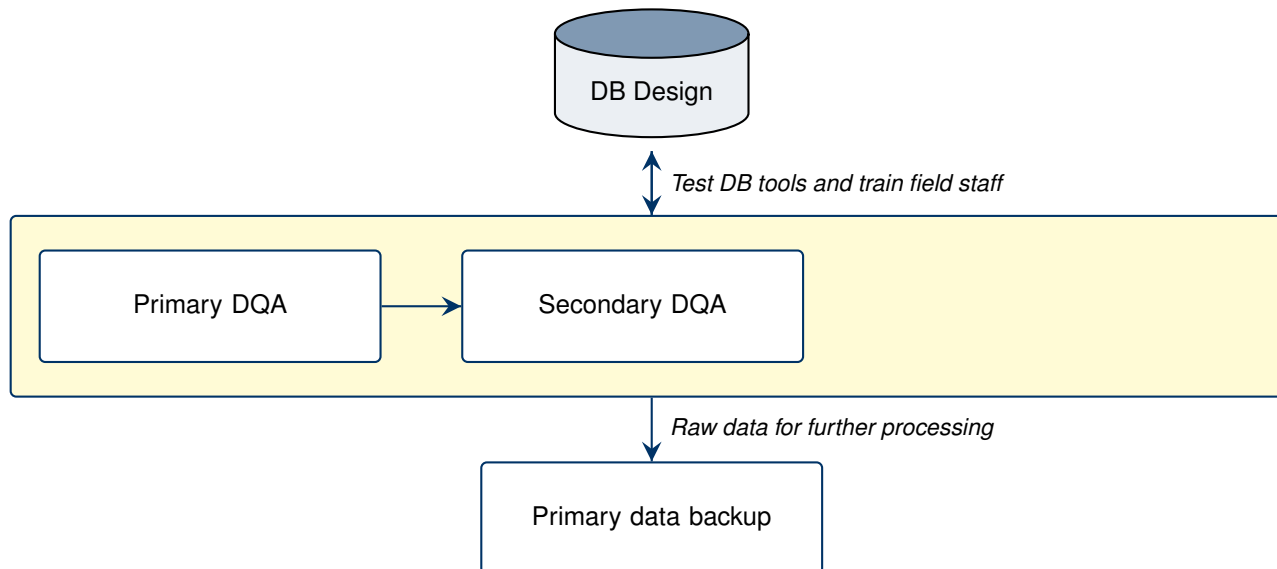
Data management process flow documentation example



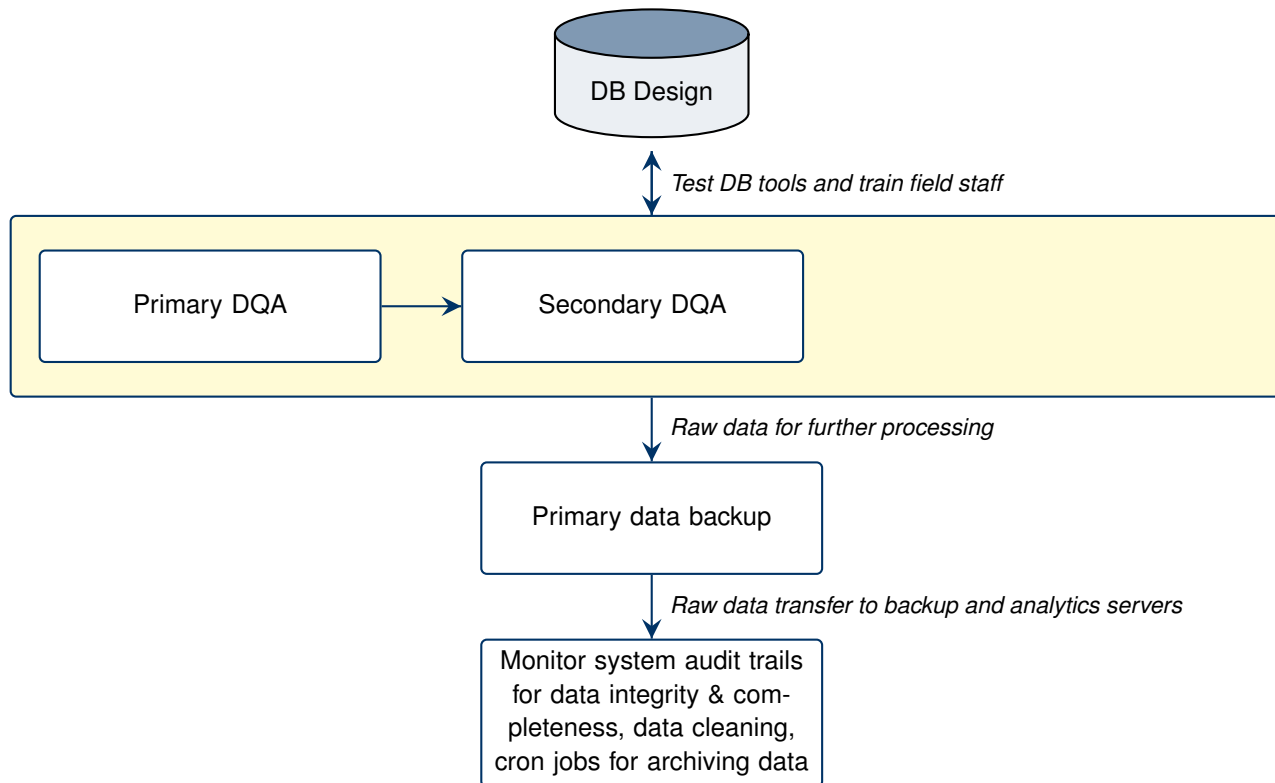
Data management process flow documentation example



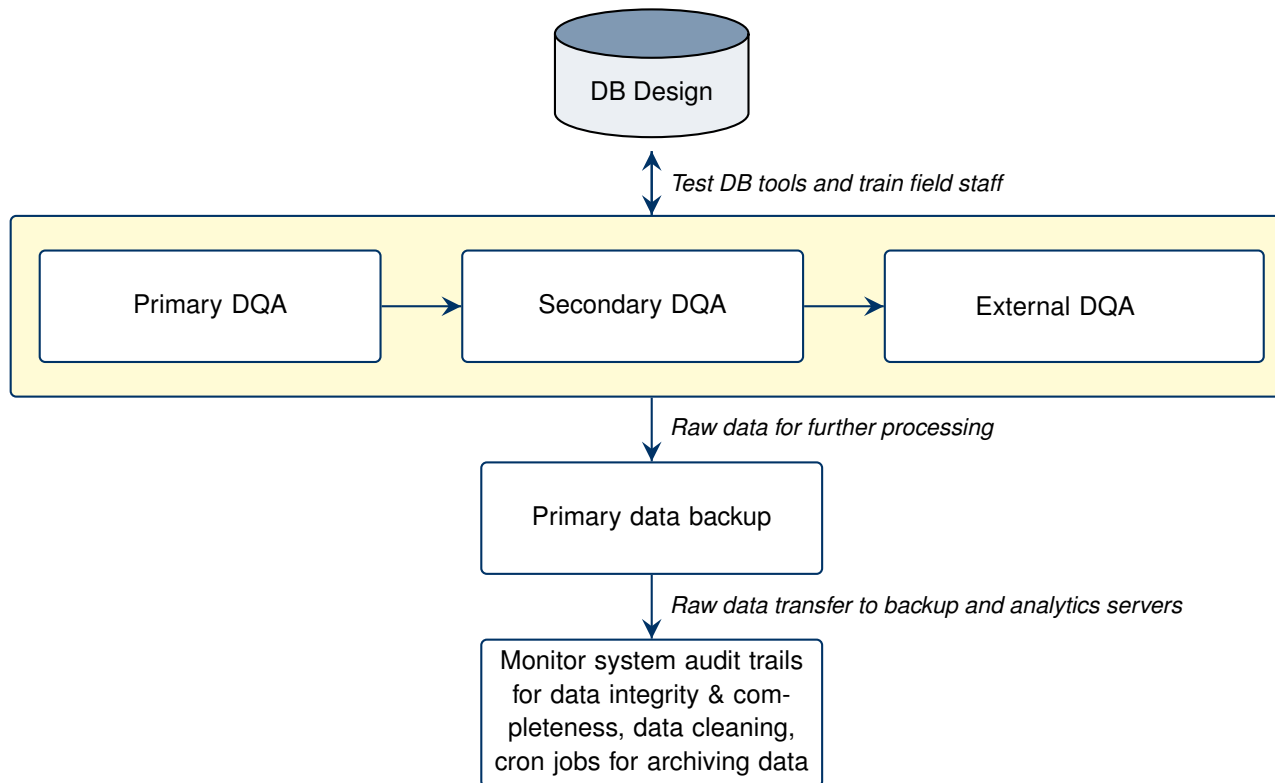
Data management process flow documentation example



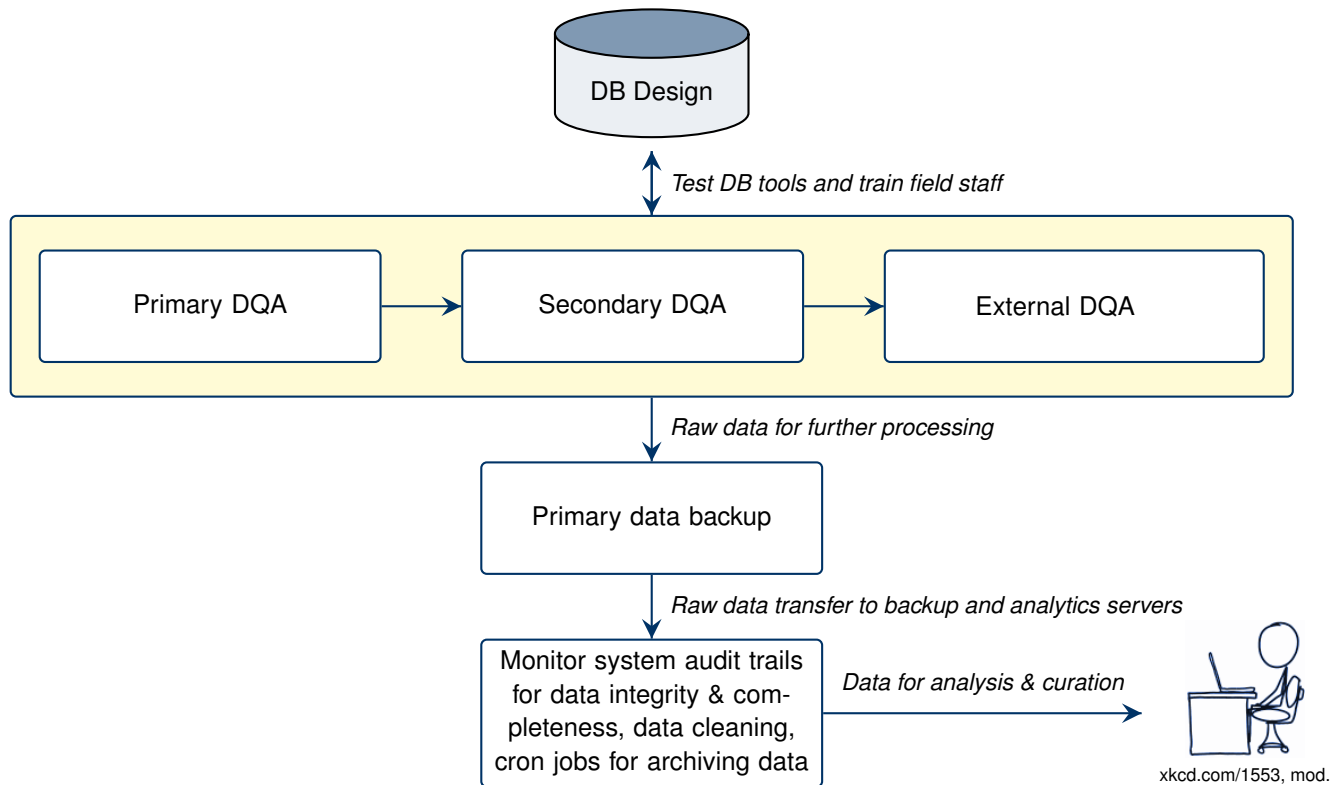
Data management process flow documentation example



Data management process flow documentation example



Data management process flow documentation example



Data capture documentation (SOP) example

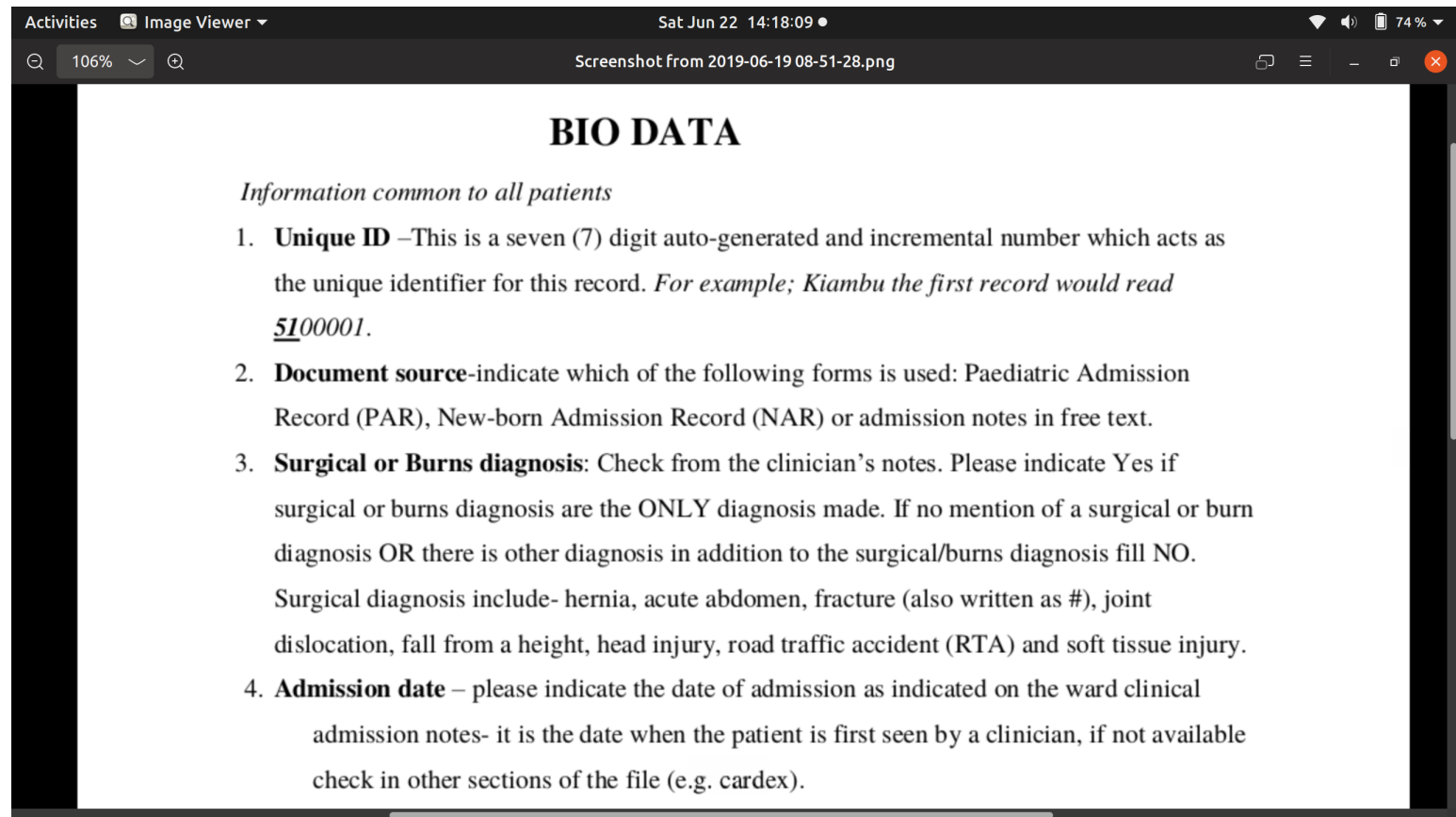


Figure: Eg.SOP

Example post-project (publication) documentation

See example documentation files @
<https://doi.org/10.7910/DVN/PP6QRQ>

- variables codebook.
- readme file.

METADATA

Metadata

What is this?



TOP DEFINITION

metadata

A fancy word for "information" invented by tech folks to make their jobs sound harder than they really are.

My metadata is flaring up today.

#information #data #words #seo #marketing

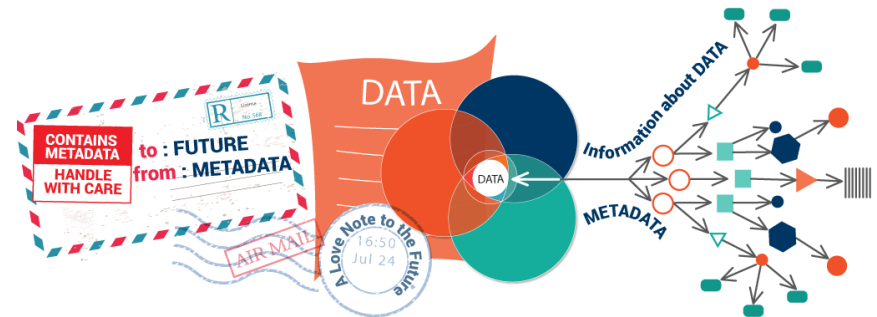
by **clayton rattlesnake** January 25, 2012

 18  9

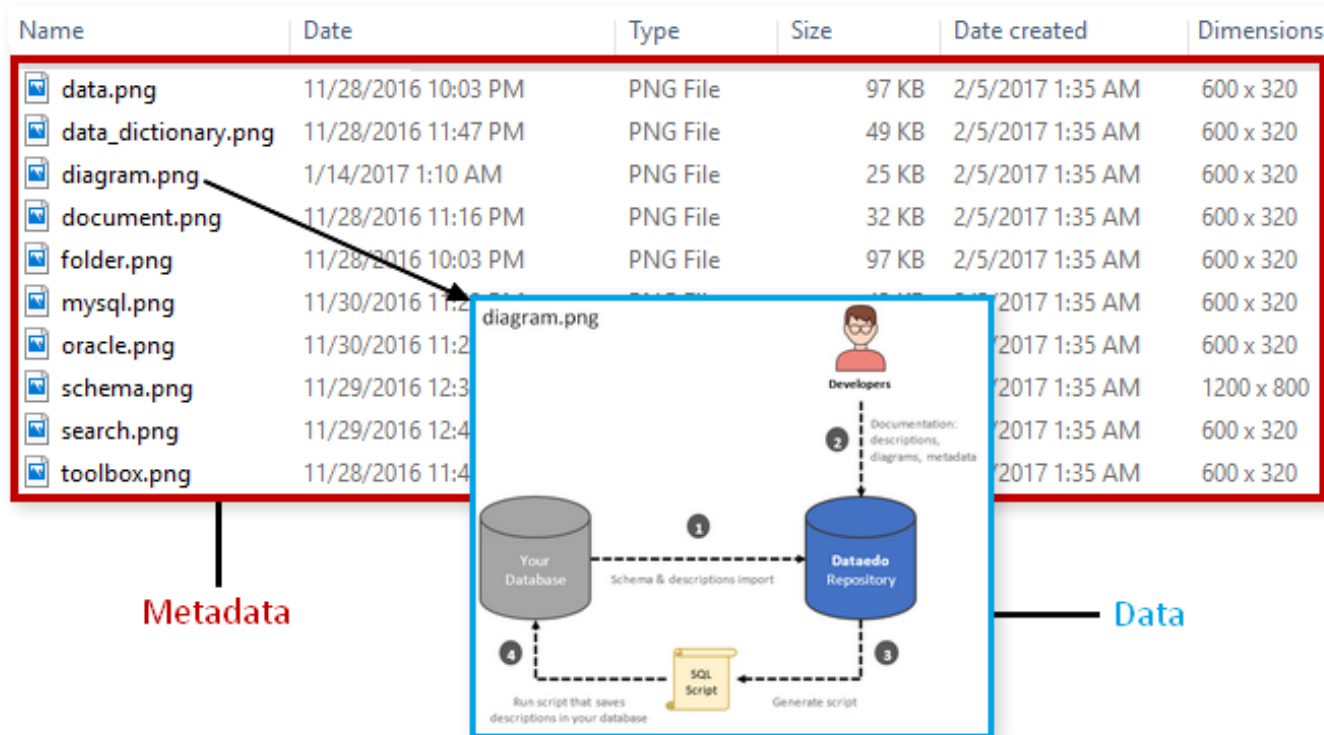


Metadata made simple

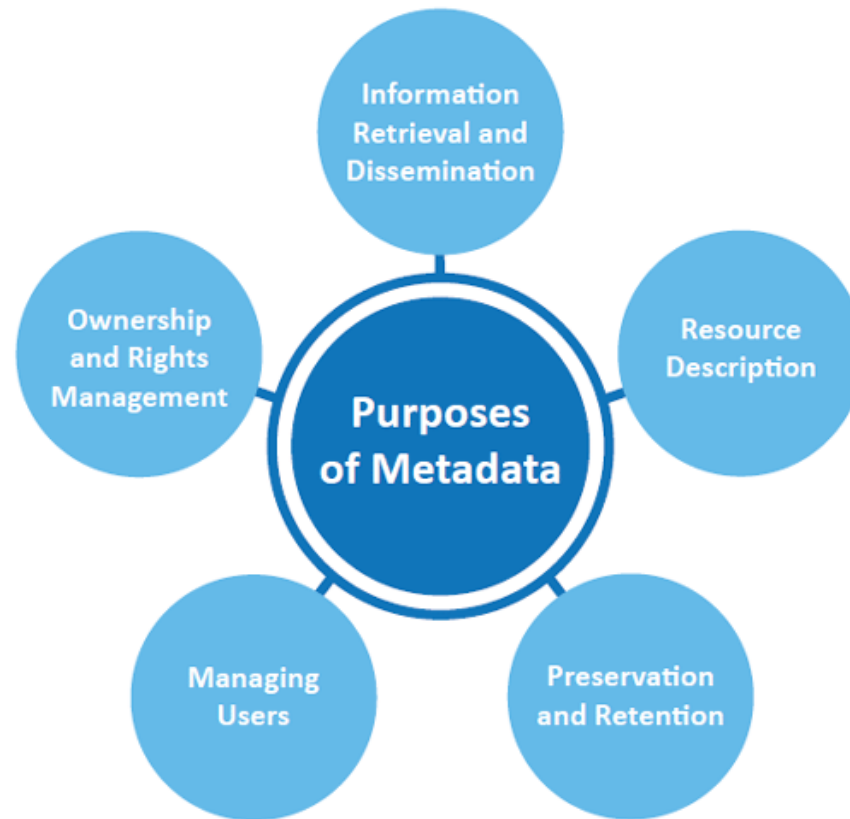
- Data about data.
- It explains the origin, purpose, time, geographic location, creator, access, and terms of use of the data.
- Important for purposes of retrieving and indexing data in a repository or archives, and for citation.



Quick example of everyday metadata



Functions of metadata



Functions of metadata

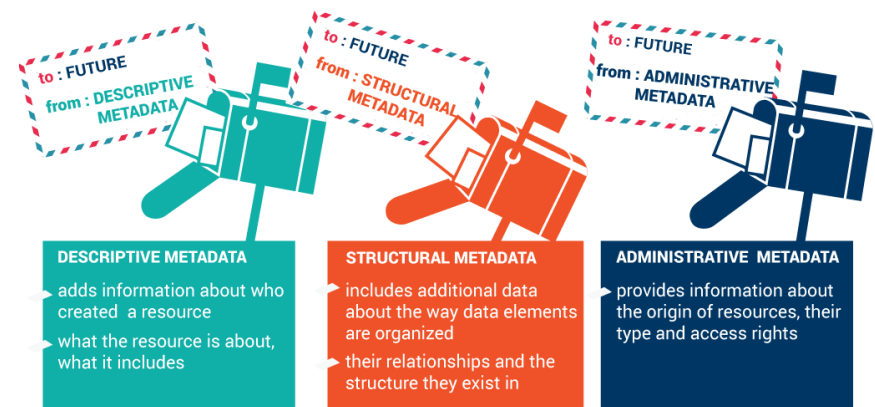
- **Resource discovery:** Allowing resources to be found by relevant criteria; Identifying resources; Bringing similar resources together; Distinguishing dissimilar resources; Giving location information.
- **Digital identification:** e.g. using digital object identifiers (DOI).
- **Organizing e-resources:** Organizing links to resources based on audience or topic; Building these pages dynamically from metadata stored in databases.

Functions of metadata

- **Facilitating interoperability:** Using defined metadata schemes, shared transfer protocols, and crosswalks between schemes, resources across the network can be searched more seamlessly. Cross-system search, e.g., using Z39.50 protocol; Metadata harvesting, e.g., OAI protocol.
- **Archiving and preservation:** requires tracking of the lineage of a digital object, and documentation of its behavior in order to emulate it in future technologies.

Types of metadata

- **Descriptive:** for discovery and identification.
- **Structural:** indicates how compound objects are put together e.g. how pages are ordered to form chapters.
- **Administrative:** Rights management and preservation.



Categories of metadata

- **Application Metadata:** Data created by the application specific to the electronically stored information (ESI) being addressed, embedded in the file, and moved with the file when copied. Copying may alter application metadata.
- **Document Metadata:** Properties about the file stored in the file, as opposed to stored in the document content. Often this data is not immediately viewable in the software application used to create/edit the document, but can generally be accessed via a "Properties" view. Examples include document author and company, and creation or revision dates.

Categories of metadata

- **File System Metadata:** Metadata generated by the system to track the demographics (name, size, location, usage, etc.) of the ESI that are stored externally from, rather than embedded within, the ESI.
- **Embedded Metadata:** Generally hidden, but an integral part of ESI, such as "track changes" or "comments" in a word processing file or "notes" in a presentation file. While some metadata is routinely extracted during processing and conversion for e-discovery, embedded data may not be. Therefore, it may only be available in the original, native file.

Metadata example 1

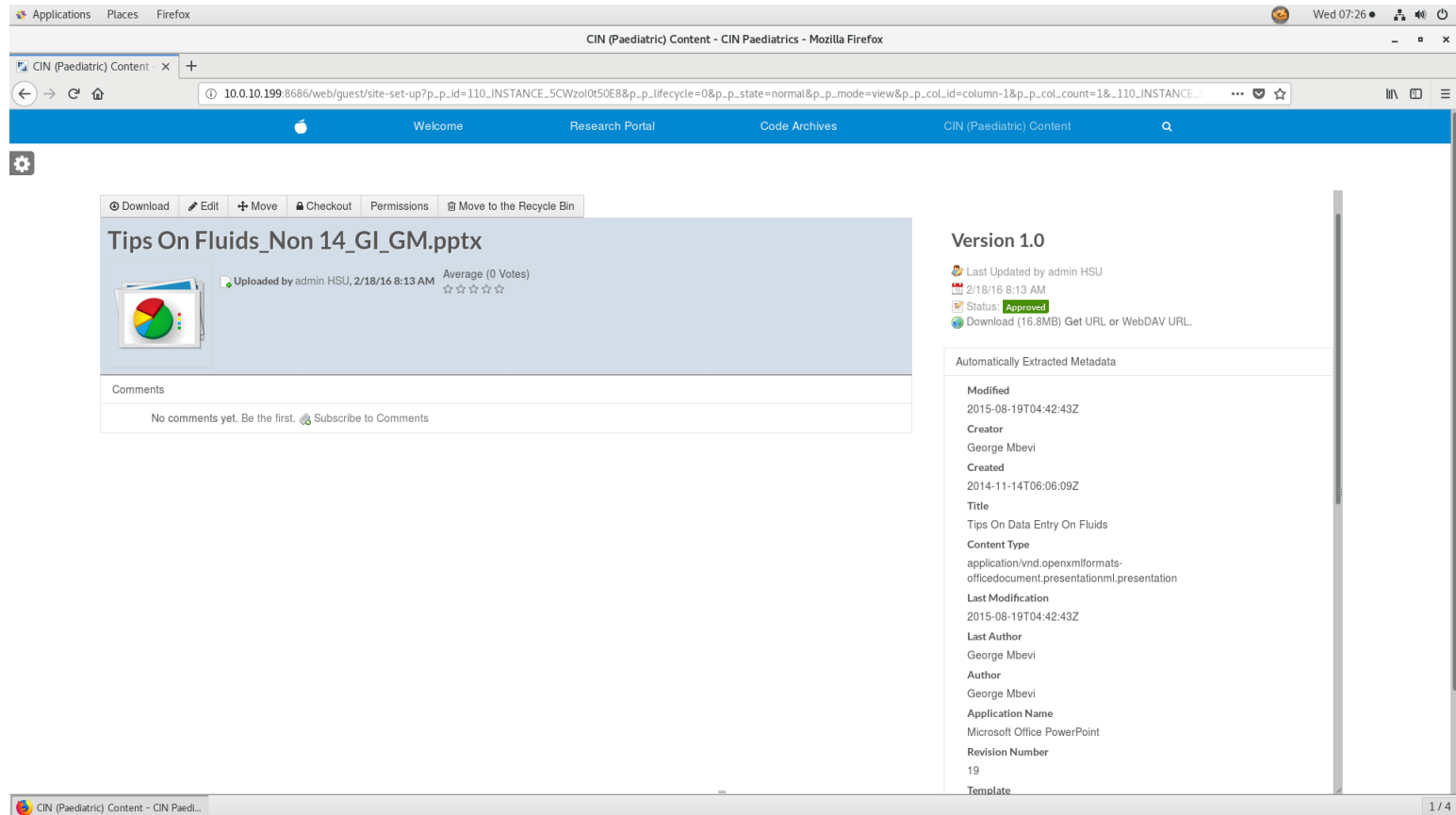


Figure: Eg.1 Metadata

Metadata example 2

Activities LibreOffice Calc Wed Jun 19 07:41:36

AntibioticConsumptionSurvey201_DataDictionary_2017-07-07.csv - LibreOffice Calc

File Edit View Insert Format Styles Sheet Data Tools Window Help

Liberation Sans: 10

A1 fx Σ = Variable / Field Name

	A	B	C	D	E
	Variable / Field Name	Form Name	Section Header	Field Type	Field Label
1	unique_id	biodata		text	Unique ID
2	ward_code	biodata		text	Ward code
3	activity	biodata		dropdown	Activity
4	ip_no	biodata		text	Patient's IP Number
5	age_recorded	biodata		radio	Age Recorded?
6	age_less1mnth	biodata		yesno	Age less than 1 month
7	age_days	biodata		text	Age (Days)
8	age_years	biodata		text	Age (years)
9	age_mths	biodata		text	Age (months)
10	weight	biodata		text	Weight (kgs)
11	gender	biodata		radio	Gender
12	date_today	biodata		text	Today's Date
13	date_adm	biodata		text	Admission Date
14	readm_hosp	biodata		radio	Re-admission to this hospital?
15	num_of_antibios_given	treatment_culture		text	Number of antibiotics given
16	antibio_1	treatment_culture		sql	Antibiotic 1
17	date_started_1	treatment_culture		text	Start date
18	single_unit_dose_1	treatment_culture		text	Single unit dose
19	dose_unit_1	treatment_culture		radio	<i>Unit<i>
20	dose_freq_1	treatment_culture		dropdown	<i>Doses/day<i>
21	route_1	treatment_culture		radio	<i>Route<i>
22	dx_1	treatment_culture		sql	Diagnosis
23	indication_type_1	treatment_culture		sql	Type of indication
24	reason_in_notes_1	treatment_culture		yesno	Reason in notes?
25	guideline_compliance_1	treatment_culture		radio	Guideline compliance
26	review_date_doc_1	treatment_culture		yesno	Step/revision date as duration documented?

Sheet 1 of 1 | Default | English (USA) | Average: ; Sum: 0 | 100%

Figure: Eg.2 Metadata

Metadata example 3

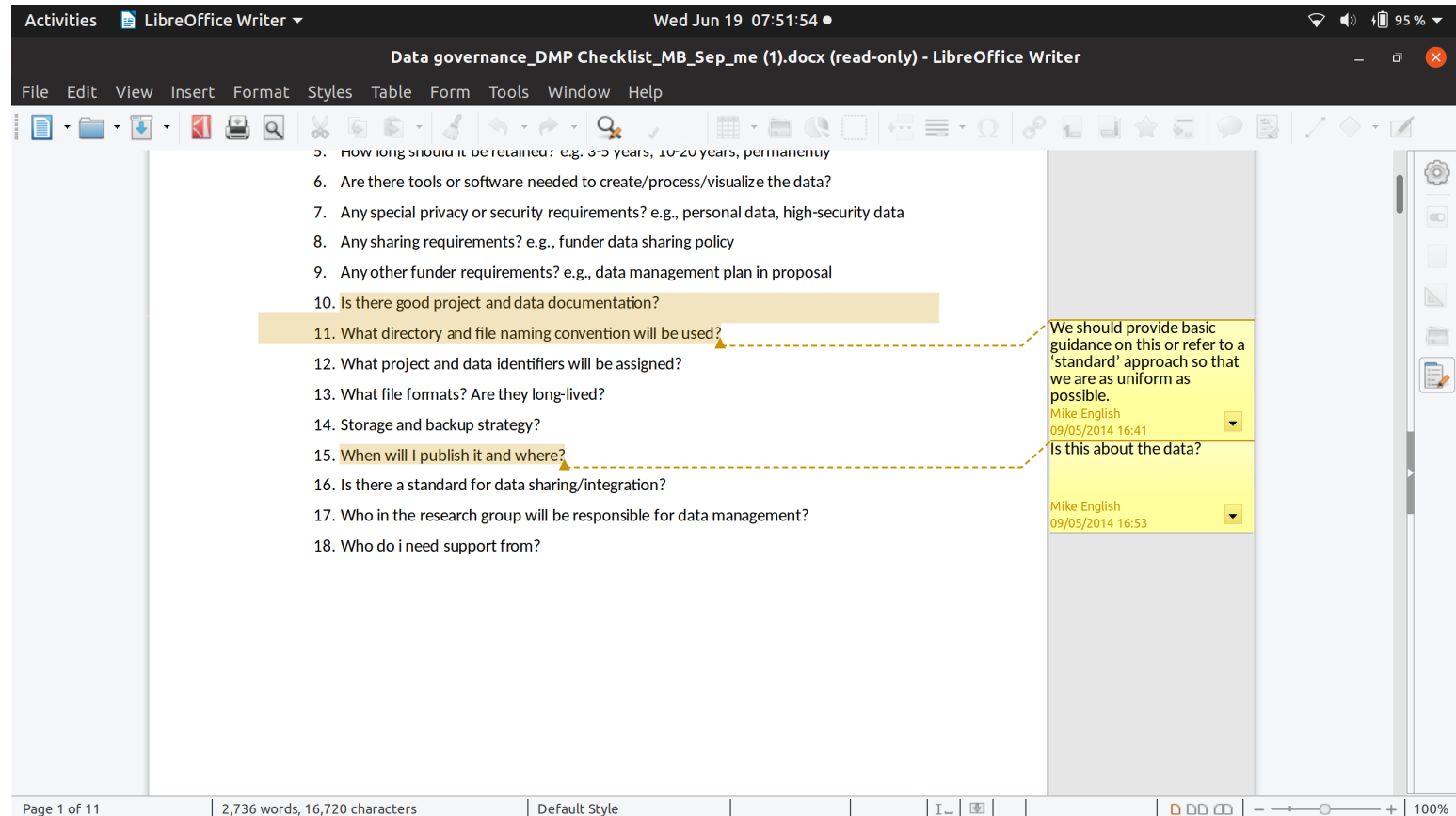
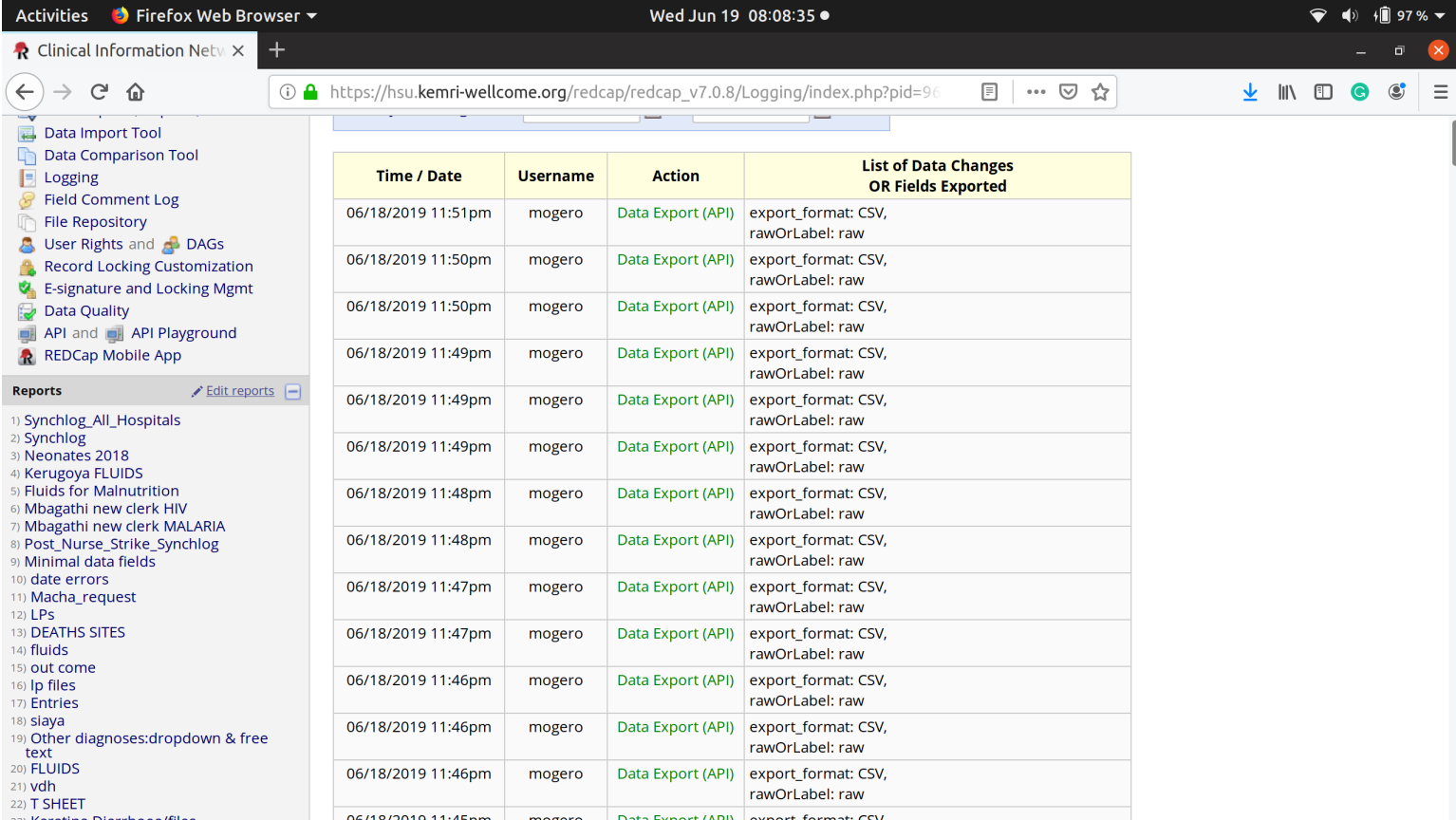


Figure: Eg.3 Metadata

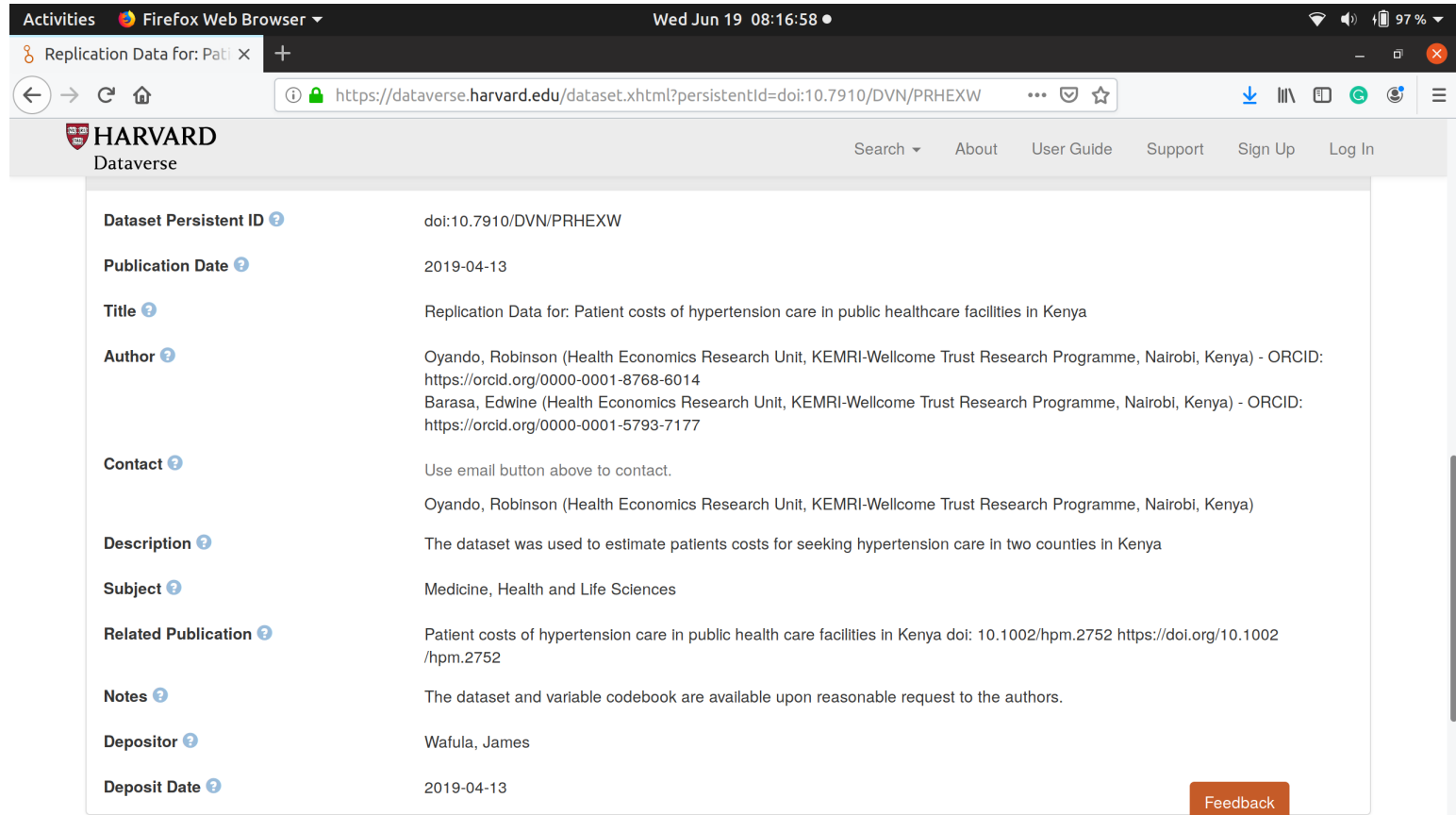
Metadata example 4



Time / Date	Username	Action	List of Data Changes OR Fields Exported
06/18/2019 11:51pm	mogero	Data Export (API)	export_format: CSV, rawOrLabel: raw
06/18/2019 11:50pm	mogero	Data Export (API)	export_format: CSV, rawOrLabel: raw
06/18/2019 11:50pm	mogero	Data Export (API)	export_format: CSV, rawOrLabel: raw
06/18/2019 11:49pm	mogero	Data Export (API)	export_format: CSV, rawOrLabel: raw
06/18/2019 11:49pm	mogero	Data Export (API)	export_format: CSV, rawOrLabel: raw
06/18/2019 11:49pm	mogero	Data Export (API)	export_format: CSV, rawOrLabel: raw
06/18/2019 11:48pm	mogero	Data Export (API)	export_format: CSV, rawOrLabel: raw
06/18/2019 11:48pm	mogero	Data Export (API)	export_format: CSV, rawOrLabel: raw
06/18/2019 11:47pm	mogero	Data Export (API)	export_format: CSV, rawOrLabel: raw
06/18/2019 11:47pm	mogero	Data Export (API)	export_format: CSV, rawOrLabel: raw
06/18/2019 11:46pm	mogero	Data Export (API)	export_format: CSV, rawOrLabel: raw
06/18/2019 11:46pm	mogero	Data Export (API)	export_format: CSV, rawOrLabel: raw
06/18/2019 11:46pm	mogero	Data Export (API)	export_format: CSV, rawOrLabel: raw
06/18/2019 11:45pm	mogero	Data Export (API)	export_format: CSV, rawOrLabel: raw

Figure: Eg.4 Metadata

Metadata example 5



Activities Firefox Web Browser Wed Jun 19 08:16:58

Replication Data for: Pat: X

https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/PRHEXW

HARVARD
Dataverse

Search About User Guide Support Sign Up Log In

Dataset Persistent ID	doi:10.7910/DVN/PRHEXW
Publication Date	2019-04-13
Title	Replication Data for: Patient costs of hypertension care in public healthcare facilities in Kenya
Author	Oyando, Robinson (Health Economics Research Unit, KEMRI-Wellcome Trust Research Programme, Nairobi, Kenya) - ORCID: https://orcid.org/0000-0001-8768-6014 Barasa, Edwine (Health Economics Research Unit, KEMRI-Wellcome Trust Research Programme, Nairobi, Kenya) - ORCID: https://orcid.org/0000-0001-5793-7177
Contact	Use email button above to contact. Oyando, Robinson (Health Economics Research Unit, KEMRI-Wellcome Trust Research Programme, Nairobi, Kenya)
Description	The dataset was used to estimate patients costs for seeking hypertension care in two counties in Kenya
Subject	Medicine, Health and Life Sciences
Related Publication	Patient costs of hypertension care in public health care facilities in Kenya doi: 10.1002/hpm.2752 https://doi.org/10.1002/hpm.2752
Notes	The dataset and variable codebook are available upon reasonable request to the authors.
Depositor	Wafula, James
Deposit Date	2019-04-13

Feedback

Figure: Eg.5 Metadata

Metadata example 6

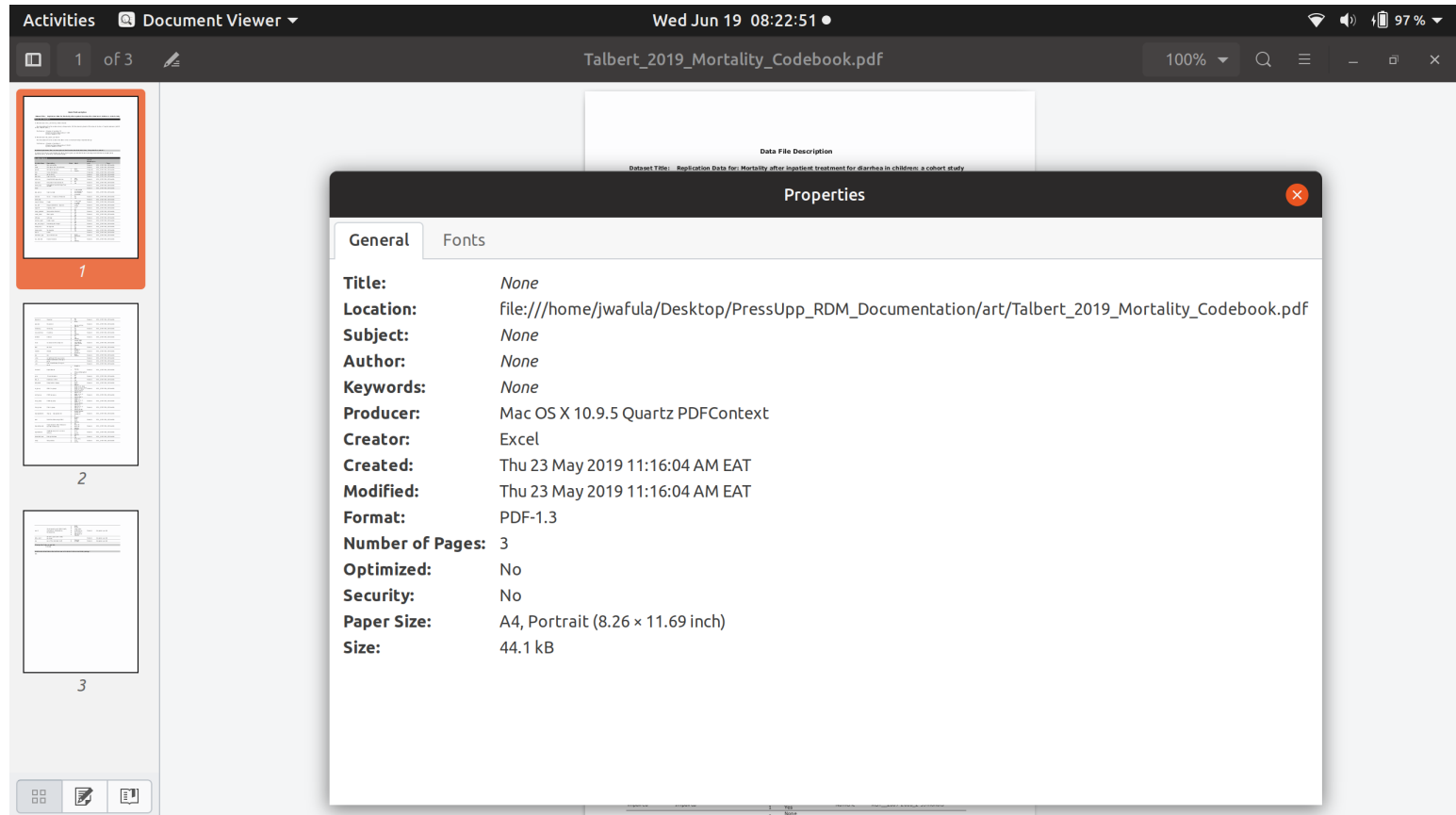


Figure: Eg.6 Metadata

Metadata standards

- Metadata takes the form of a structured set of elements which describe the information resource and assists in the identification, location and retrieval of it by users.
- Aim of standardisation is to make metadata support a number of defined functions.
- Often starts as schemas developed by a particular user community to enable the best possible description of a resource type for their needs.

Examples of metadata standards

Refer to:

<http://www.dcc.ac.uk/resources/subject-areas/biology>

<http://rd-alliance.github.io/metadata-directory/standards/>

- OAI Dublin Core - minimal requirement for all OAI providers.
- MARCXML - utilizes marcxml as the metadata schema name.
- e-GMS (e-Government Metadata Standard).
- ISO 19115: 2003(E) — Geographic Information: Metadata.
- PREMIS: Data Dictionary for Metadata Preservation.
- Metadata Object Description Schema (MODS): for library applications.

Metadata standards supported by Harvard Dataverse

- Citation Metadata: compliant with DDI Lite, DDI 2.5 Codebook, DataCite 3.1, and Dublin Core's DCMI Metadata Terms.
- Geospatial Metadata: compliant with DDI Lite, DDI 2.5 Codebook, DataCite, and Dublin Core.
- Social Science Humanities Metadata: compliant with DDI Lite, DDI 2.5 Codebook, and Dublin Core.
- Astronomy and Astrophysics Metadata : Based on Virtual Observatory (VO) Discovery and Provenance Metadata.
- Life Sciences Metadata: based on ISA-Tab Specification.

Exposing metadata via harvesting

- There is increasing need for open access to published work and reproducibility in research.
- Metadata harvesting is the automated collection of metadata descriptions from different sources to create useful aggregations of metadata.
- Metadata can be harvested for data sharing through the Open Archives Initiative Protocol for Metadata Harvesting (OAI-PMH). OAI-PMH is a low-barrier mechanism for repository interoperability.
- OAI-PMH has two levels of participants:
 - Data providers - who administer the system.
 - Service providers - who use the metadata harvested to build their digital collection.

Exposing metadata via harvesting

- A harvester is used by a service provider as a way to collect metadata from a repository.
- A repository is a network-accessible server that is able to process OAI-PMH requests. Dataverse is an example of an OAI server.
- Only the published, unrestricted datasets in Dataverse can be made harvestable.
- A repo is managed by the data provider to allow harvesters access to its metadata.

Metadata harvesting

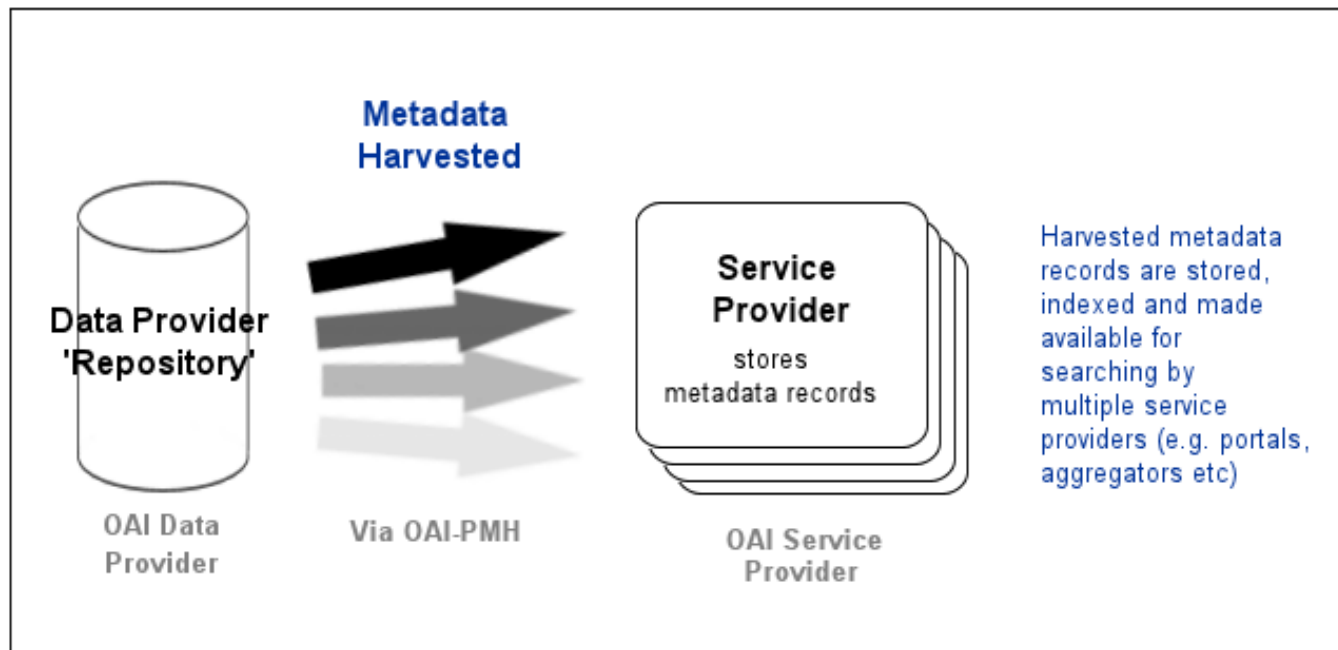


Figure 2. Exposing Metadata via Harvesting

Metadata harvesting

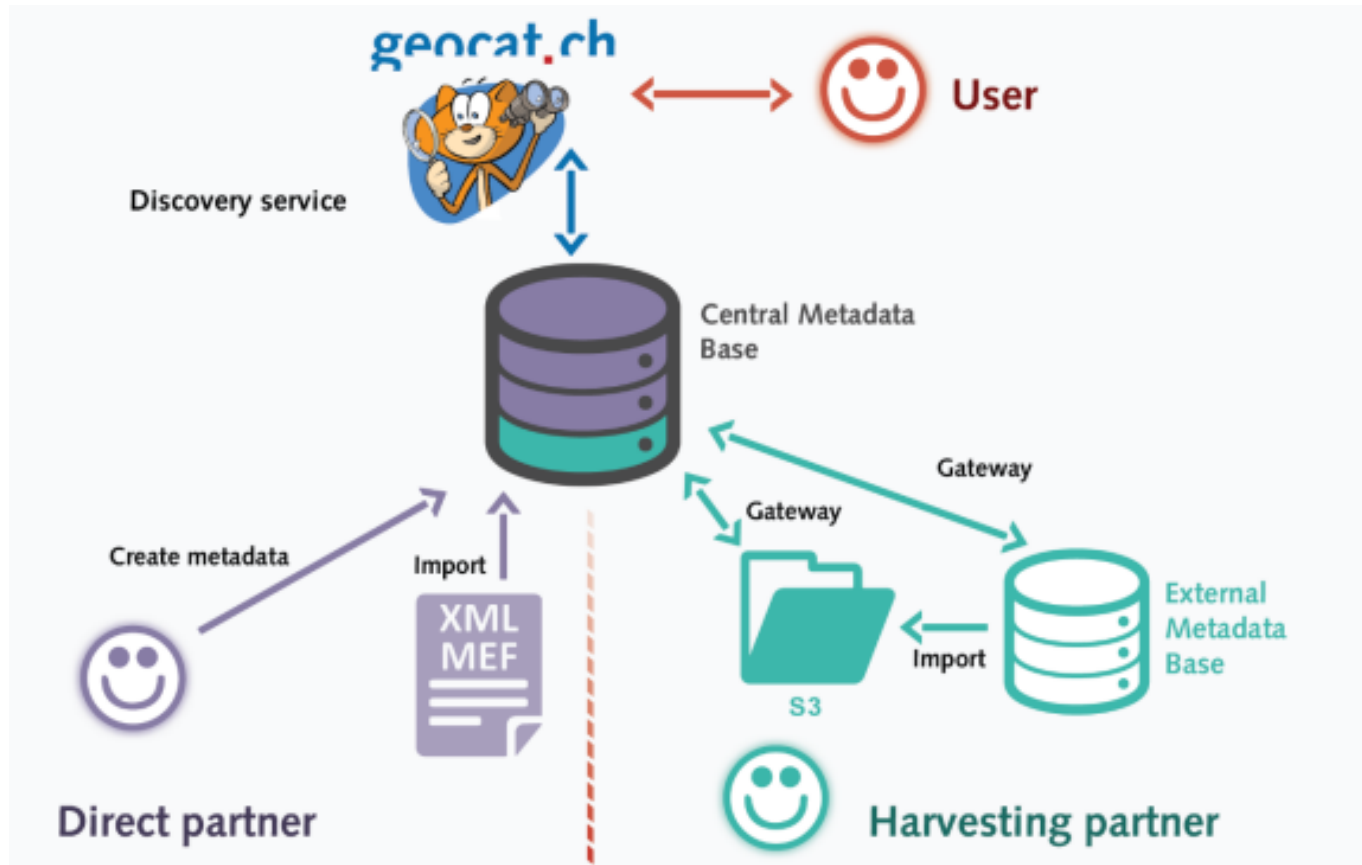


Figure: <https://www.geocat.admin.ch/en/information/partner-model.html>

Acknowledgements

