



Norwegian
Meteorological
Institute

Data Management Handbook template for MET and partners in S-ENDA written in asciidoc

Nina E. Larsgård, Elodie Fernandez, Morten W. Hansen, ...

Table of Contents

1. Introduction	3
1.1. The principles of data management for geophysical data	3
1.1.1. External data management requirements and forcing mechanisms	4
1.1.2. The geophysical value chain	4
1.1.3. A data management model based on the FAIR principles	5
1.1.4. Dataset	6
1.1.5. Metadata	6
1.2. Human roles in data management	8
1.2.1. Data consumer	8
1.2.2. Data provider	8
1.2.3. Data Management Roles	8
1.3. Summary of data management requirements	8
2. Global attributes that should be added to NetCDF-CF files	11
Glossary of Terms and Names	15

Abstract

Abstract will come here..

Revision history

Version	Date	Comment	Responsible
2.0	2021-??-??	New version based on original MET DMH	Nina E. Larsgård, Elodie Fernandez, Morten W. Hansen, ...

1. Introduction

The purpose of the Data Management Handbook (DMH) is threefold:

1. to provide an overview of the principles for data management to be employed.
2. to help personnel identify their roles and responsibilities for good data management.
3. to provide personnel with practical guidelines for carrying out good data management.

Data management is the term used to describe the handling of data in a systematic and cost-effective manner. The data management regime should be continuously evolving, to reflect the evolving nature of data collection. Therefore this DMH is a living document that will be revised and updated from time to time in order to maintain its relevance.

The DMH is a strategic governing document and should be used as part of the quality framework the organisation is using.

This DMH focuses on the management of geophysical data. Types of data and the usage of these data are described in the organisation specific information in Section X.X.X. This introduction (Chapter 1) lays forth the background and principles for the data management regime. Chapters 2-5 describe the implementation of the main building blocks: structuring and documenting data (Chapter 2), data services (Chapter 3), portals and documentation aimed at users (Chapter 4) and governance issues (Chapter 5). Each chapter starts with a brief statement of its purpose, followed by a description of what is implemented at the organisation at present, as well as the planned developments for the short-term (<2 years) and expected developments for the longer term (2-5 years). Practical guidelines for carrying out good data management are addressed in Chapter 6 and especially in the Quick Manual for Data Providers. The Quick Manual is a concise, informal HOW-TO for data management practitioners.

The intended audience for this DMH is any co-worker who is involved in any part of the value chain for geophysical data ([Figure 1](#)).

The handbook can be used in three ways:

1. Read the Introduction (Chapter 1) to find out the background and principles of data management;
2. Read Chapters 2-5 to learn about how data management is currently implemented and how it is expected to evolve in the next few years;
3. Read the Quick Manual for Data Providers for guidelines on what to do in real-life cases; alternatively read Chapter 6 for typical workflow examples.

1.1. The principles of data management for geophysical data

One of the main motivations for implementing a unified data management model is to better serve the users of the data. Primarily, this can be approached by making user needs and requirements the guide for determining what data we provide and how. For example, it will be described below how

the specification of datasets should be determined by user needs. By implementing the data management practices described here, it is expected that users will benefit from:

- the ease of discovering, viewing and accessing all the datasets that are offered by the institute;
- standardised ways of downloading data, which reduces the need for special solutions on the user side;
- reducing their own data storage needs, by downloading just what they need;
- easy and standardised access to remote datasets and catalogues, when using their own visualisation/analysis tools;
- the ability to compare and combine data from internal and external sources (through metadata catalogues);
- the ability to apply common data transformations, like spatial, temporal and variables subsetting and reprojection, before downloading anything; *the ease of building specific metadata catalogues and data portals that include data from the institute and can target a specific user community;
- the access to datasets which can be integrated in their internal and external workflows through standardised web services.

The principles of standardised data documentation, publication, sharing and preservation have been formalised in the The [FAIR](#) Guiding Principles for scientific data management and stewardship [RD3] through a process facilitated by FORCE11.

FAIR - findability, accessibility, interoperability and reusability

1.1.1. External data management requirements and forcing mechanisms

Any organisation that strives to implement [FAIR](#) data management model has to relate to external forcing mechanisms concerning data management at several levels. At the national level, the organisation must comply with national regulations as decided by the government. Some of these are indications of expected behaviour (e.g. OECD regulations) and some are implemented through a legal framework. The Norwegian government has over time promoted free and open sharing of public data. Mechanisms for how to do this are governed by the [Geodataloven](#) (implemented as [Geonorge](#)), which is a national implementation of the European INSPIRE directive (to be amended in 2019). INSPIRE defines a federated multinational Spatial Data Infrastructure (SDI) for the European Union, similar to NSDI in the USA or UNSDI under the United Nations. The goal is to provide a standardised access to data and provide the necessary tools to be able to work with the data in a unified manner. In short, these legal frameworks require standardised documentation (at discovery and use level; these concepts are described later) and access (through specified protocols) to the data identified.

Other external requirements and forcing mechanisms that are organisation-specific are listed under section X.X.X

1.1.2. The geophysical value chain

An example of a geophysical value chain is presented in Figure 1. Typically, data from a wide variety of providers are used in the value chain. Traditionally, the data used have been transmitted

on request from one data centre to another, and used in the specific processing chains that requested the data. The focus on reuse of data in various contexts has been missing.

[Value chain] | *value_chain.png*

Figure 1. Value chain for geophysical data

At the end of the data management value chain are the users of the data (aka. data consumers, see Section 1.2.2.1), who may be either external or internal to the institute.

1.1.3. A data management model based on the FAIR principles

This model is based on the model of the Arctic Data Centre, which adheres to the FAIR principles. The model's basic functions fall into three main categories:

1. **Documentation of data** using discovery and use metadata (metadata are further described below). The documentation identifies who, what, when, where, and how, and shall make it easy for consumers to find and understand data. This requires application of information containers and utilisation of controlled vocabularies and ontologies where textual representation is required. It also covers the topic of data provenance which is used to describe the origin and all actions done on a dataset. Data provenance is closely linked with workflow management. Furthermore, it covers the relationship between datasets. Application of ontologies in data documentation is closely linked to the concept of linked data.
2. **Publication and sharing of data** focuses on making data accessible to consumers internally and externally. Application of standardised approaches is vital, along with cost efficient solutions that are sustainable. Direct integration of data in applications for analysis through data streaming minimises the complexity and overhead in dissemination solutions. This category also covers persistent identifiers for data.
3. **Preservation of data** includes short and long term management of data, which secures access and availability throughout the lifespan of the data. Good solutions in this area depend on expected and actual usage of the data. Preservation of data includes the concept of data life cycle, i.e., the documented flow of data from initial storage through to obsolescence and permanent archiving (or deletion).

For its implementation, the data management model is built upon the following principles:

- **Standardisation** – compliance with established international standards;
- **Interoperability** – enabling machine-to-machine interfaces and standardised documentation and encoding of data;
- **Integrity** – ensuring that data and access to them can be maintained over time, ensuring the user receives the same data each time;
- **Traceability** – documentation of the provenance of a dataset, i.e., all actions taken to produce and maintain the dataset and the usage of the data in downstream systems;
- **Modularisation** – enabling replacement of one component of the system without necessitating other changes. In this data management model there are two terms that are fundamental and may be addressed immediately: **dataset** and **metadata**.

1.1.4. Dataset

A dataset is a collection of data. In the context of the data management model, the storage mode of the dataset is irrelevant, since access mechanisms can be decoupled from the storage layer as experienced by a data consumer. Typically, a dataset represents a number of variables in time and space. A more detailed definition is provided in the [Glossary of Terms](#). In order to best serve the data through the web services developed, the following guidance is given for defining datasets:

1. A dataset can be a collection of variables stored in, for example, a relational database or as flat files.
2. A dataset is defined as a number of spatial and/or temporal variables.
3. A dataset should be defined by the information content and not the production method. This implies that the output of, for example, a numerical model may be divided into several datasets that are related. This is also important in order to efficiently serve the data through [web services](#). For instance, model variables defined on different vertical coordinates should be separated as [linked datasets](#), since some OGC [//Link here//](#) services (e.g. WMS) are unable to handle mixed coordinates in the same dataset.
4. A good dataset does not mix feature types, e.g. do not combine trajectories and gridded data in one dataset.

Most importantly, a dataset should be defined to meet a consumer need. This means that the specification of a dataset should follow not only the content guidelines just listed, but also address the user needs for delivery, security and preservation.

1.1.5. Metadata

Metadata is a broad concept. In our data management model the term “metadata” is used in several contexts, specifically the five categories that are briefly described in [Table 1](#).

Table 1. Brief introduction to different types of metadata.

Type	Purpose	Description	Examples
Discovery metadata	Used to find relevant data	Discovery metadata are also called index metadata and are a digital version of the library index card. They describe who did what, where and when, how to access data and potential constraints on the data. They shall also link to further information on the data like site metadata. Discovery metadata are thus WIS metadata.	ISO 19115 GCMD DIF

Type	Purpose	Description	Examples
Use metadata	Used to understand data found	Use metadata describe the actual content of a dataset and how it is encoded. The purpose is to enable the user to understand the data without any further communication. They describe the content of variables using standardised vocabularies, units of variable, encoding of missing values, map projections, etc.	Climate and Forecast (CF) Convention BUFR GRIB
Configuration metadata	Used to tune portal services for datasets for users	Configuration metadata are used to improve the services offered through a portal to the user community. This can be e.g. how to best visualise a product.	
Site metadata	Used to understand data found	Site metadata are used to describe the context of observational data. They describe the location of an observation, the instrumentation, procedures, etc. To a certain extent they overlap with discovery metadata, but also extend the discovery metadata. Site metadata can be used for observation network design. Site metadata can be considered a type of use metadata.	WIGOS OGC O&M StInfoSys

Type	Purpose	Description	Examples
System metadata	Used to understand the technical structure of the data management system and track changes in it	System metadata covers e.g. technical details of the storage system (e.g. Lustre metadata), web services, their purpose and how they interact with other components of the data management system, available and consumed storage, number of users and other KPI elements etc.	SysDok

The tools and facilities used to manage the information contained in the metadata are further described in Chapter 2.

1.2. Human roles in data management

1.2.1. Data consumer

The Data Consumer may be a scientist or student, employee of a governmental agency, consultant or some other person with a professional or personal interest in the data provided. Data consumers may be internal or external to the entities providing and managing the data.

1.2.2. Data provider

The Data Provider is generating datasets managed by the data management system described in this document. Data providers can be internal or external to the system. They should be able to maintain the datasets they have committed.

1.2.3. Data Management Roles

Between the data providers and data consumers are the processes that manage and deliver the datasets (cf. [img-value-chain](#)). A number of human roles may be defined with responsibilities that, together, ensure that these processes are carried out in accordance with the data management requirements of the organisation. The definition and filling of these roles depend heavily on the particular organisation, and each organisation must devise its own best solution.

1.3. Summary of data management requirements

The data management regime described in this DMH follows the Arctic Data centre model and shall ensure that:

1. There are relevant metadata for all datasets, and both data and metadata are available in a form and in such a way that they can be utilised by both humans and machines.

- a. There are sufficient metadata for each dataset for both discovery and use purposes.
 - b. Discovery metadata are indexed and can be retrieved from available services in a standard way and with standard protocols.
 - c. There are interfaces for discovery, visualisation and download, as well as portals for human access, that operate seamlessly across institutions.
 - d. The data are described in a relevant, standardised and managed vocabulary that supports machine-machine interfaces.
 - e. Datasets have attached a unique and permanent identifier, i.e., UUID, that enables traceability.
 - f. Datasets have licensing that ensures free use and reuse wherever possible.
 - g. Datasets are available for download in a standard form according to the FAIR guiding principles (NetCDF/CF and equivalent) and through standard protocols (OPeNDAP, OGC WCS, secure direct download, etc.) that are accepted and utilised in the user environment.
 - h. There are authentication and authorisation mechanisms that ensure access control to data with restrictions, and that are compatible with and coupled to relevant public authentication solutions (FEIDE, eduGAIN, Google, etc.).
2. There is an organisation that provides for the management of each dataset throughout its lifetime (life cycle management).
 - a. There is documentation that describes physical storage, lifetime of each dataset, degree of storage redundancy, metadata consistency methods, how dataset versioning is implemented and unique IDs to ensure traceability. The organisation provides seamless access to data from distributed data centres through various portals.
 - b. The above and a business model at dataset level are described in a Data Management Plan (DMP)
 3. There are services or tools that provide the following functionalities on the datasets:
 - a. Transformations
 - i. subsetting
 - ii. slicing of gridded data sets to points, sections, profiles
 - iii. reprojection
 - iv. resampling
 - v. reformatting
 - b. Visualisation (time series, mapping services, etc.)
 - c. Aggregation
 - d. Upload of new datasets (including enabling and configuring data access services)

2. Global attributes that should be added to NetCDF-CF files

In order to add netCDF-CF datasets to the discovery metadata catalog, the data producer should populate the file with certain global attributes mainly described in the Attribute Convention for Data Discovery (ACDD). For a complete description of the ACDD elements, please refer to http://wiki.esipfed.org/index.php/Attribute_Convention_for_Data_Discovery.

The ACDD recommendations should always be followed in order to properly document your netCDF-CF files. However, the below tables summarize the attributes that are needed to properly populate a discovery metadata catalog which fulfills the requirements of international standards (e.g., GCMD/DIF, the WMO profile of ISO19115, etc.).

The following ACDD elements are required:

ACDD Attribute	Repetition allowed	Separator	Default	MMD equivalent
id	no			metadata_identifier
date_created	yes	,		last_metadata_update>update>datetime
date_metadata_modified	yes	,		last_metadata_update>update>datetime
title	yes	;		title>title
summary	yes	;		abstract>abstract
geospatial_lat_max	no		90	geographic_extent>rectangle>north
geospatial_lat_min	no		-90	geographic_extent>rectangle>south
geospatial_lon_max	no		180	geographic_extent>rectangle>east
geospatial_lon_min	no		-180	geographic_extent>rectangle>west
keywords_vocabulary	yes	,		keywords>vocabulary
keywords	yes	,		keywords>keyword

The following ACDD elements are recommended:

ACDD Attribute	Repetition allowed	Separator	Default	MMD equivalent
time_coverage_start	yes	,	1850-01-01T00:00:00Z	temporal_extent>start_date

time_coverage_end	yes	,		temporal_extent>end_date
processing_level	no			operational_status
license	no			use_constraint>identifier
['creator_role', 'contributor_role']	yes	,	unknown	personnel>role
['creator_name', 'contributor_name']	yes	,	unknown	personnel>name
creator_email	yes	,	unknown	personnel>email
creator_institution	yes	,	unknown	personnel>organisation
institution	yes	,		data_center>data_center_name>short_name
institution	yes	,		data_center>data_center_name>long_name
publisher_url	yes	,		data_center>data_center_url
project	yes	;		project>short_name
project	yes	;		project>long_name
platform	yes	,		platform>short_name
platform	yes	,		platform>long_name
platform_vocabulary	yes	,		platform>resource
instrument	yes	,		platform>instrument>short_name
instrument	yes	,		platform>instrument>long_name
instrument_vocabulary	yes	,		platform>instrument>resource
source	yes	;		activity_type
creator_name	yes	,		dataset_citation>author
date_created	yes	,		dataset_citation>publication_date

title	yes	,		dataset_citation>title
publisher_name	yes	,		dataset_citation>publisher
metadata_link	yes	,		dataset_citation>url
references	yes	,		dataset_citation>other

In addition, some global attributes are useful for the discovery metadata catalog but do not exist in ACDD. Please refer to the documentation of [MMD](https://htmlpreview.github.io/?https://github.com/metno/mmd/blob/master/doc/mmd-specification.html) [https://htmlpreview.github.io/?https://github.com/metno/mmd/blob/master/doc/mmd-specification.html] for a description of these elements:

Extra Attribute	Repetition allowed	Separator	Default	MMD equivalent
date_created_type	yes	,	Created	last_metadata_update>update>type
collection	yes	,		collection
title_lang	yes	,	en	title>lang
abstract_lang	yes	,	en	abstract>lang
dataset_production_status	no			dataset_production_status
license_resource	no			use_constraint>resource
contributor_email	yes	,	unknown	personnel>email
contributor_organisation	yes	,	unknown	personnel>organisation
related_dataset_relation_type	yes			related_dataset>relation_type
related_dataset_id	yes			related_dataset>id
iso_topic_category	yes	,		iso_topic_category
keywords_resource	yes	,		keywords>resource

Glossary of Terms and Names

Term	Description
Application service	TBC
CDM dataset	A dataset that “may be a NetCDF, HDF5, GRIB, etc. file, an OPeNDAP dataset, a collection of files, or anything else which can be accessed through the NetCDF API.” Unidata Common Data Model [https://www.unidata.ucar.edu/software/netcdf-java/v4.6/CDM/index.html]
Configuration metadata	See Configuration metadata definition in Table 2
Controlled vocabulary	A carefully selected list of terms (words and phrases) controlled by some authority. They are used to tag information elements (such as datasets) so that they are easier to search for. (see Wikipedia article [https://en.wikipedia.org/wiki/Controlled_vocabulary]) A basic element in the implementation of the Semantic web .
Data life cycle management	“Data life cycle management (DLM) is a policy-based approach to managing the flow of an information system’s data throughout its life cycle: from creation and initial storage to the time when it becomes obsolete and is deleted.” Excerpt from TechTarget [https://searchstorage.techtarget.com/definition/data-life-cycle-management] article. Alias: life cycle management
Data Management Plan	“A data management plan (DMP) is a written document that describes the data you expect to acquire or generate during the course of a research project, how you will manage, describe, analyse, and store those data, and what mechanisms you will use at the end of your project to share and preserve your data.” Stanford Libraries [https://library.stanford.edu/research/data-management-services/data-management-plans]
Data centre	A combination of a (distributed) data repository and the data availability services and information about them (e.g., a metadata catalog). A data centre may include contributions from several other data centres.
Data management	How data sets are handled by the organisation through the entire value chain - include receiving, storing, metadata management and data retrieval.

Term	Description
Data provenance	“The term ‘data provenance’ refers to a record trail that accounts for the origin of a piece of data (in a database, document or repository) together with an explanation of how and why it got to the present place.” (Gupta, 2009). See also Boohers (2015) [https://www.theboohers.org/2015/03/03/provenance/]
Data repository	A set of distributed components that will hold the data and ensure they can be queried and accessed according to agreed protocols. This component is also known as a Data Node.
Dataset	<p>A dataset is a pre-defined grouping or collection of related data for an intended use. Datasets may be categorised by:</p> <p><i>Source</i>, such as observations (in situ, remotely sensed) and numerical model projections and analyses;</p> <p><i>Processing level</i>, such as “raw data” (values measured by an instrument), calibrated data, quality-controlled data, derived parameters (preferably with error estimates), temporally and/or spatially aggregated variables;</p> <p><i>Data type</i>, including point data, sections and profiles, lines and polylines, polygons, gridded data, volume data, and time series (of points, grids, etc.).</p> <p>Data having all of the same characteristics in each category, but different independent variable ranges and/or responding to a specific need, are normally considered part of a single dataset. In the context of data preservation a dataset consists of the data records and their associated knowledge (information, tools). In practice, our datasets should conform to the Unidata CDM dataset definition, as much as possible.</p>
Discovery metadata	See Discovery metadata definition in Table 2
Dynamic geodata	Data describing geophysical processes which are continuously evolving over time. Typically these data are used for monitoring and prediction of the weather, sea, climate and environment.

Term	Description
FAIR principles	The four foundational principles of good data management and stewardship: *F*indability, *A*ccessibility, *I*nteroperability and *R*eusability. Nature article [RD3 [https://www.nature.com/articles/sdata201618]], FAIR Data Principles [https://www.go-fair.org/fair-principles/], FAIR metrics proposal [https://github.com/FAIRMetrics/Metrics], EU H2020 Guidelines [https://ec.europa.eu/research/participants/data/ref/h2020/grants_manual/hi/oa_pilot/h2020-hi-oa-data-mgt_en.pdf]
Feature type	A categorisation of data according to how they are stored, for example, grid, time series, profile, etc. It has been formalised in the NetCDF/CF feature type table [https://www.nodc.noaa.gov/data/formats/netcdf/v2.0/#templatesexamples], which currently defines eight feature types.
Geodataloven	“Norwegian regulation toward good and efficient access to public geographic information for public and private purposes.” See https://www.regjeringen.no/no/tema/plan-bygg-og-eiendom/plan&#8212;&#8203;og-bygningsloven/kart/geodataloven/id749728/ ” De ling av geodata – Geodataloven [font size="0.85em"> https://www.regjeringen.no/no/tema/plan-bygg-og-eiendom/plan&#8212;&#8203;og-bygningsloven/kart/geodataloven/id749728/].
Geonorge	“Geonorge is the national website for map data and other location information in Norway. Users of map data can search for any such information available and access it here.” See Geonorge [https://www.geonorge.no/en/].
Geographic Information System	A geographic information system (GIS) is a system designed to capture, store, manipulate, analyze, manage and present spatial or geographic data. (Clarke, K. C., 1986) GIS systems have lately evolved in distributed Spatial Data Infrastructures (SDI)
Glossary	Terms and their definitions, possibly with synonyms.
Interoperability	The ability of data or tools from non-cooperating resources to integrate or work together with minimal effort.

Term	Description
[linked-data]]Linked data	A method of publishing structured data so that they can be interlinked and become more useful through semantic queries [https://en.wikipedia.org/wiki/Semantic_query], i.e., through machine-machine interactions. (see Wikipedia article [https://en.wikipedia.org/wiki/Linked_data])
Ontology	A set of concepts with attributes and relationships that define a domain of knowledge.
OpenSearch	A collection of simple formats for the sharing of search results (OpenSearch [https://github.com/dewitt/opensearch])
Product	“Product” is not a uniquely defined term among the various providers of dynamical geodata, either nationally or internationally. It is often used synonymously with “dataset.” For the sake of clarity, “product” is not used in this handbook. The term “dataset” is adequate for our purpose.
Semantic web	“The Semantic Web provides a common framework that allows data to be shared and reused across application, enterprise, and community boundaries”. W3C [https://www.w3.org/2001/sw/] (see Wikipedia article [https://en.wikipedia.org/wiki/Semantic_Web])
Site metadata	See Site metadata definition in Table 2
Spatial Data Infrastructure	“Spatial Data Infrastructure (SDI) is defined as a framework of policies, institutional arrangements. technologies, data, and people that enables the sharing and effective usage of geographic information by standardising formats and protocols for access and interoperability.” (Tonchovska et al, 2012) SDI has evolved from GIS . Among the largest implementations are: NSDI in the USA, INSPIRE in Europe and UNSDI as an effort by the United Nations. For areas in the Arctic, there is arctic-sdi.org [https://arctic-sdi.org/].
Unified data management	A common approach to data management in a grouping of separate data management enterprises.
Use metadata	See Use metadata definition in Table 2
Web portal	A central website where all users can search, browse, access, transform, display and download datasets irrespective of the data repository in which the data are held.

Term	Description
Web service	Web services are used to communicate metadata, data and to offer processing services. Much effort has been put on standardisation of web services to ensure they are reusable in different contexts. In contrast to web applications, web services communicate with other programs, instead of interactively with users. (See TechTerms article [https://techterms.com/definition/web_service])
Workflow management	Workflow management is the process of tracking data, software and other actions on data into a new form of the data. It is related to data provenance, but is usually used in the context of workflow management systems .
(Scientific) Workflow management systems	A scientific workflow system is a specialised form of a workflow management system designed specifically to compose and execute a series of computational or data manipulation steps, or workflow, in a scientific application. (Wikipedia [https://en.wikipedia.org/wiki/Scientific_workflow_system]) As of today, many different frameworks exist with their own proprietary languages, these might eventually get connected by using a common workflow definition language [https://www.commonwl.org/].