

# Lecture 22 – Data Models for Discovery (Meta-data models)

*Open Data Management & the Cloud*

(Data Science & Scientific Computing / UniTS – DMG)

# Lecture summary



- Data discovery
- Why data discovery
- How do we use data: some hints
- How do we collect data: some examples and a demo
- Where we collect data (repos, catalogues, archives recap)
- Metadata representations to make possible and ease the data collection

# What is Data Discovery?



- **Data discovery** involves the **collection** and **evaluation** of data from various sources and is often (in business) used to understand trends and patterns in the data.

The data discovery process includes

- **connecting multiple data sources**,
- **cleansing** and preparing the data,
- **sharing** the data throughout the organization
- **performing analysis** to gain insights into business processes.

# What is Data Discovery?



- Huge amounts of data are collected in business on customers, markets, suppliers, production processes, and more.
- Data are collected differently depending from the environment:
  - from online and traditional transactions systems, sensors, social media, mobile devices, and other sources.
  - In Astrophysics from ground- and space-based instrumentsr also from simulations

# Why Data Discovery?



Insights are hidden within the data.

- Huge amounts of data are collected in business on customers, markets, suppliers, production processes, and more.
- Need to collect data from different resources
- Decision makers has to extract insight from data
  - Business Intelligence
  - Visual Analitics
- Researchers also extract knowledge from data

# Knowledge, Wisdom, Insight



**Knowledge** is the accumulation of facts and data that you have learned about or experienced. It's being aware of something, and having information. Knowledge is really about facts and ideas that we acquire through study, research, investigation, observation, or experience.

**Wisdom** is the ability to discern and judge which aspects of that knowledge are true, right, lasting, and applicable to your life. It's the ability to apply that knowledge to the greater scheme of life. It's also deeper; knowing the meaning or reason; about knowing why something is, and what it means to your life.

**Insight** is the deepest level of knowing and the most meaningful to your life. Insight is a deeper and clearer perception of life, of knowledge, of wisdom. It's grasping the underlying nature of knowledge, and the essence of wisdom. Insight is a truer understanding of your life and the bigger picture of how things intertwine.

<https://www.lifehack.org/articles/communication/what-are-the-differences-between-knowledge-wisdom-and-insight.html>

Christopher Reiss does a great job of summing up the differences on Quora

<https://www.quora.com/Wisdom/What-are-the-valuable-differences-between-knowledge-wisdom-and-insight-Beyond-their-basic-definitions-what-benefits-do-they-hold/answer/Christopher-Reiss?share=1&srid=ot>

- **Knowledge** is measuring that a desert path is 12.4 miles long.
- **Wisdom** is packing enough water for the hike.
- **Insight** is building a lemonade stand at mile 6.

# What is Business Intelligence?



- **Business intelligence (BI)** includes the applications, infrastructure, tools, and best practices that enable
  - access to and
  - analysis ofinformation to **improve and optimize decisions and performance.**
- Data are collected from online and traditional transactions systems, sensors, social media, mobile devices, and other sources.

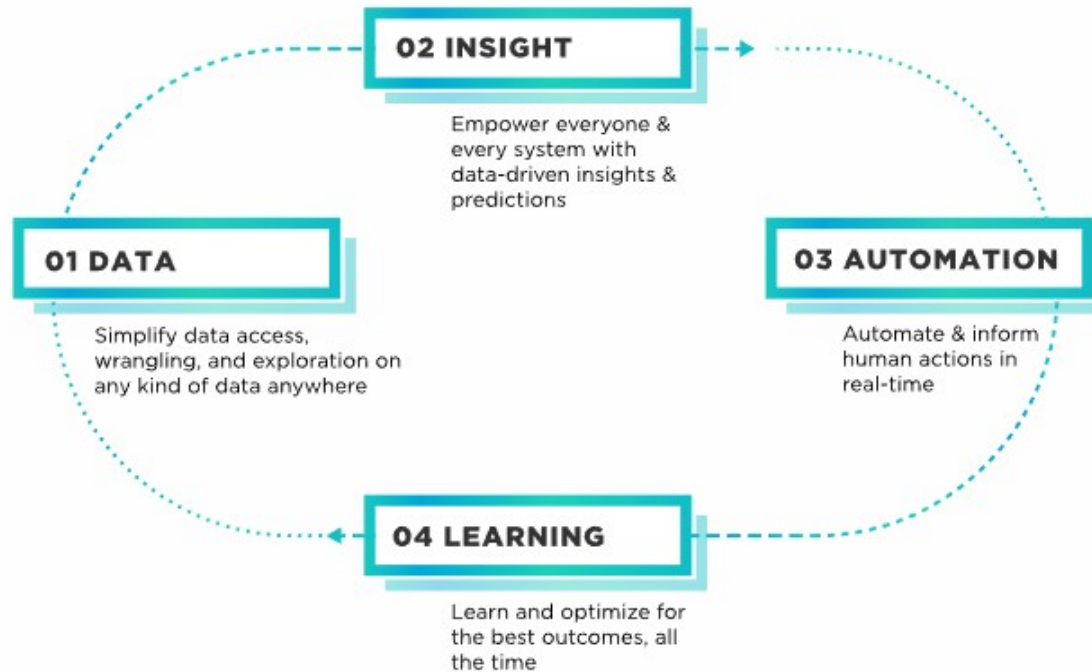
**Insights are hidden within that data.**

Decision makers need to extract insight from data

- Business Intelligence
- Visual Analytics



# How Business Intelligence works



<https://www.tibco.com/reference-center/what-is-business-intelligence>

Business intelligence enables data visualizations

**Data visualization** is the graphical representation of data to help people understand context and significance.

Interactive (clickable) data visualization enables companies to

- drill down to explore details,
- identify patterns and outliers
- change which data is processed and/or excluded.

When data is visualized, it's easier to identify emerging trends, the very first step in deriving insight.

# Visual Analytics



Visual analytics is a form of reasoning that uses interactive, visual interfaces.

Visual analytics uses data analytics and interactive visual representations of the data and dashboarding to enable users to interpret large volumes of data.

Visual analytics combine visualization, human factors, and data analysis to gain knowledge from data.

When presented visually, analytics are easier and faster for users to interpret.

Visual analytics make complex issues much easier to understand

The interactive and visual elements are often helpful in communicating what one sees in the data to others and in making better-informed business decisions.

Data scientists can find visual analytics useful to show and clarify business trends and other concepts to those users that are not data scientists or not know complex statistical algorithms.

# Business Intelligence vs Visual Analytics



- **BI** concentrates on **current and past events** recorded in the data
- **BA** focuses mostly on what is more likely to happen in the **future**.

Data visualizations answer the “what” questions, but visual analytics help you get to the “why.”

# Useful links



<https://www.tibco.com/reference-center/what-is-data-discovery>

<https://www.tibco.com/reference-center/what-is-business-intelligence>

<https://www.tibco.com/reference-center/guide-to-data-visualization>

<https://www.tibco.com/reference-center/what-is-visual-analytics>

<https://www.simplilearn.com/business-intelligence-vs-business-analytics-article>

# Registry vs Repository



- A **registry** is a list of items with pointers for where to find the items, like the index on a database table or the card catalog for a library.
- A **repository** stores the actual items, like a database table itself or a library's shelves of books.

If you lose a registry, the items still exist; you just may need to *reindex* them.

If you lose a repository, the items are gone.

# Catalog vs Registry



- A **catalog** is a list of something, or a book containing a list. Examples of a catalog:
  - a library's list of all of the books it has available.
  - a booklet showing everything a store has for sale.
- A **registry** is a list of items with pointers for where to find the items. Examples of registry:
  - the index on a database table
  - the card catalog for a library. It is a database containing the network locations of service instances.

The service/data registry is a key part of service/data discovery.

A service registry needs to be highly available and up to date.

Clients can cache network locations obtained from the service/data registry.

# How to collect data



Chapter 5: Collecting data | Cochrane Training - Chromium

training.cochrane.org/handbook/current/chapter-05

Do you have feedback about the new online Handbook? [Open feedback form](#) [Close](#)

**Cochrane Training**  
Trusted evidence.  
Informed decisions.  
Better health.

Search...

Online learning   Learning events   Guides and handbooks   Trainers' Hub   [Log in](#)

## Chapter 5: Collecting data

Tianjing Li, Julian PT Higgins, Jonathan J Deeks

**Key Points:**

- Systematic reviews have studies, rather than reports, as the unit of interest, and so multiple reports of the same study need to be identified and linked together before or after data extraction.
- Because of the increasing availability of data sources (e.g. trials registers, regulatory documents, clinical study reports), review authors should decide on which sources may contain the most useful information for the review, and have a plan to resolve discrepancies if information is inconsistent across sources.
- Review authors are encouraged to develop outlines of tables and figures that will appear in the review to facilitate the design of data collection forms. The key to successful data collection is to construct easy-to-use forms and collect sufficient and unambiguous data that faithfully represent the source in a structured and organized manner.
- Effort should be made to identify data needed for meta-analyses, which often need to be calculated or converted from data reported in diverse formats.
- Data should be collected and archived in a form that allows future access and data sharing.

Cite this chapter as: Li T, Higgins JPT, Deeks JJ (editors). Chapter 5: Collecting data. In: Higgins JPT, Thomas J, Chandler J, Cumpston M, Li T, Page MJ, Welch VA (editors). *Cochrane Handbook for Systematic Reviews of*

Search Handbook

Supplementary material: Appendix of resources

- ♦ Chapter 5: Collecting data
  - ♦ 5.1 Introduction
  - ♦ 5.2 Sources of data
  - ♦ 5.3 What data to collect
  - ♦ 5.4 Data collection tools
  - ♦ 5.5 Extracting data from reports
  - ♦ 5.6 Extracting study results and converting to the desired format
  - ♦ 5.7 Managing and sharing data

<https://training.cochrane.org/handbook/current/chapter-05>



# How to collect data – Web scraping



**Web scraping**, web harvesting, or web data extraction is a way to extract data from websites.

Web scraping software may directly access the World Wide Web using the HyperText Transfer Protocol or a web browser. Can be done

- manually by a software user,
- automatically by a processes implemented using a bot or web crawler.

Specific **collected** data is gathered and **copied** from the web, into a central local database or spreadsheet, for later retrieval or **analysis**.

*[https://en.wikipedia.org/wiki/Web\\_scraping](https://en.wikipedia.org/wiki/Web_scraping)*

# IVOA-SAMP implementation demo



IVOA: International Virtual Observatory Alliance (see Lecture 4 on FAIR Principles for an introduction)

SAMP: Simple Application Messaging Protocol

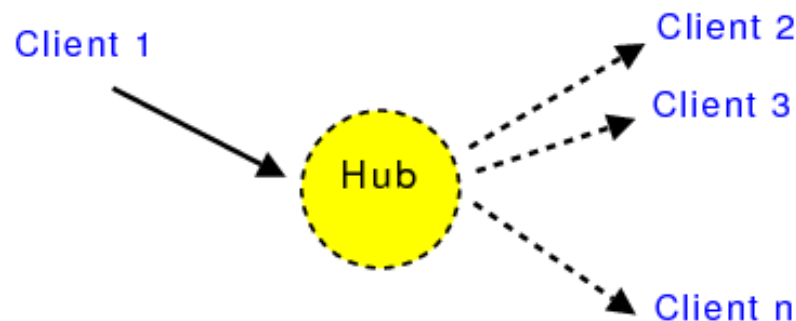
is a messaging protocol that enables astronomy software tools to interoperate and communicate. It supports communication between applications on the desktop and in web browsers.

<https://www.ivoa.net/documents/SAMP/20120411/REC-SAMP-1.3-20120411.html>

# IVOA-SAMP for data collection



SAMP has a hub-based architecture.



The hub is a single service used to route all messages between clients.

This makes application discovery more straightforward in that each client only needs to locate the hub, and the services provided by the hub are intended to simplify the actions of the client.

A disadvantage of this architecture is that the hub may be a message bottleneck and potential single point of failure.

SAMP may not be suitable for extremely high throughput requirements; may be mitigated by an appropriate strategy for hub restart if failure is likely.

<https://www.ivoa.net/documents/SAMP/20120411/REC-SAMP-1.3-20120411.html>

# IVOA-SAMP implementation example



ESCAPE EU Project <https://projectescape.eu/>

European Science Cluster of Astronomy & Particle physics ESFRI research infrastructures

brings together the astronomy, astroparticle, and particle physics communities to address fundamental challenges in data-driven research, inspired by the goals and needs of major European research infrastructures, or ESFRIs

ESAP <https://sdc-dev.astron.nl/esap-gui/>

<https://projectescape.eu/services/esfris-science-analysis-platform-esap>

ESFRI Science Analysis Platform

is a platform-service gateway with the capability to **access and combine data from multiple collections** and stage for subsequent **processing and analysis**. It allows **data discovery and handling of large and distributed data collections**. It is a flexible science platform for the analysis of open access data available through EOSC

# IVOA-SAMP implementation demo



Open Topcat

```
java -jar /home/bertocco/work/VO_School_feb2021/topcat-full.jar
```

Open esap-samp

<https://sdc-dev.astron.nl/esap-gui/samp>

Hit register

Follow the first example of the tutorial:

```
~/work/VO_School_feb2021/slides/topcat$ evince topcat_tutorial.pdf
```

# Metadata and Repositories



- Data Modelling and metadata modelling are first step in archiving, creating a repository of products
  - Custom ones for specific purposes
  - Common/shared ones to
    - Reach larger communities
    - Interoperate within or outside a research domain
- Models can be standardized exactly as can protocols or other technical specification
  - If not even more
    - Identifiers, vocabularies, formats, ...
- Better if standardization is open
  - Communities and organizations exist which have this goal
- (examples and details follow)

# Metadata Standards



- Standards for metadata change by domain and granularity
- Keeping track of them is hard work
  - An example

- RDA: Research Data Sharing without barriers
- <http://rd-alliance.github.io/metadata-directory/standards/>

## Physical Sciences & Mathematics

**AgMES (Agricultural Metadata Element Set)** [Edit](#)  
A semantic standard developed by the Food and Agriculture Organization (FAO) of the United Nations, AgMES enables description, resource discovery, interoperability and data exchange. Sponsored by the UN AIMS - Agricultural Information Management Standards, the current standard was issued in November 2010.

**AVM (Astronomy Visualization Metadata)** [Edit](#)  
The AVM scheme supports the cross-searching of collections of print-ready and screen-ready astronomical imagery rendered from telescopic observations (also known as 'pretty pictures'). Such images can combine data acquired at different wavelengths and from different observatories. While the primary intent is to cover data-derived astronomical images, there are broad

## General Research Data

**CERIF (Common European Research Information Format)** [Edit](#)  
The Common European Research Information Format is the standard that the EU recommends to its member states for recording information.

**Data Package** [Edit](#)  
The Data Package specification is a generic wrapper format for exchanging data. Although it supports arbitrary metadata, the format defines a separate but linked specification provides a way to describe the columns of a data table; descriptions of this form can be included directly in the data package.

**DataCite Metadata Schema** [Edit](#)  
A set of mandatory metadata that must be registered with the DataCite Metadata Store when minting a DOI persistent identifier for a dataset. Sponsored by the DataCite consortium, version 3.0 was recently released in 2013.

**DCAT (Data Catalog Vocabulary)** [Edit](#)  
By using DCAT to describe datasets in data catalogs, publishers increase discoverability and enable applications easily to consume metadata.

**Dublin Core** [Edit](#)  
A basic, domain-agnostic standard which can be easily understood and implemented, and as such is one of the best known and most widely used. Sponsored by the Dublin Core Metadata Initiative, Dublin Core was published as ISO Standard 15836 in February 2009.

**OAI-ORE (Open Archives Initiative Object Reuse and Exchange)** [Edit](#)  
The goal of these standards is to expose the rich content in aggregations of Web resources to applications that support authoring, depositing and retrieving content from popular social networks of "Web 2.0".

**Observations and Measurements** [Edit](#)  
This encoding is an essential dependency for the OGC Sensor Observation Service (SOS) Interface Standard. More specifically, this standard defines the identification, the extent, the quality, the spatial and temporal characteristics of observations. Since 1915-1-2014 contains the fundamental structure of NeXus files is extremely flexible, allowing the storage of heterogeneous data. Specifically, this standard defines XML schemas for observations, proteins, nucleic acids, and complex assemblies, managed by the International Union of Pure and Applied Chemistry (IUPAC). PDBx has been extended by the wwPDB requirements.

**PREMIS** [Edit](#)  
The PREMIS (Preservation Metadata: Implementation Strategies) Data Dictionary defines a set of metadata that most repositories of digital content should implement. Influence the creation of local application profiles, an XML Schema is provided to allow the metadata to be serialized independently. PREMIS was initially developed by the Preservation Metadata: Implementation Strategies Working Group, convened by OCLC and RLG.

**PROV** [Edit](#)  
Provenance is information about entities, activities, and people involved in producing a piece of data or thing, which can be used to form a record of the history of the data.

**RDF Data Cube Vocabulary** [Edit](#)  
The standard provides a means to publish multi-dimensional data, such as statistics, on the web in such a way that it can be linked to related metadata among organizations.

**Repository-Developed Metadata Schemas** [Edit](#)  
Some repositories have decided that current standards do not fit their metadata needs, and so have created their own requirements.

**UKEOF** [Edit](#)  
A metadata standard for describing environmental monitoring activities, programmes, networks and facilities published by the UK Environmental Observation Framework (UKEOF).

## Arts and Humanities

**Encoded Archival Description (EAD)** [Edit](#)  
A standard for encoding archival descriptions.

**DDI (Data Documentation Initiative)** [Edit](#)  
A widely used, international standard for describing data from the social, behavioral, and economic sciences.

**MIDAS-Heritage** [Edit](#)  
A British cultural heritage standard for recording information on buildings, archaeological sites, and other cultural heritage.

**ISA-TI** [Edit](#)  
A widely used, international standard for describing data from the social, behavioral, and economic sciences.

**OAI-ORE (Open Archives Initiative Object Reuse and Exchange)** [Edit](#)  
The goal of these standards is to expose the rich content in aggregations of Web resources to applications that support authoring, depositing and retrieving content from popular social networks of "Web 2.0".

**MIBBI** [Edit](#)  
Both versions are XML-based and defined using XML Schemas. They were developed by the International Metadata Interchange Working Group.

**MIDAS-Heritage** [Edit](#)  
A British cultural heritage standard for recording information on buildings, archaeological sites, and other cultural heritage.

**NeXus** [Edit](#)  
A metadata standard for describing environmental monitoring activities, programmes, networks and facilities published by the UK Environmental Observation Framework (UKEOF).

**OAI-ORE (Open Archives Initiative Object Reuse and Exchange)** [Edit](#)  
The goal of these standards is to expose the rich content in aggregations of Web resources to applications that support authoring, depositing and retrieving content from popular social networks of "Web 2.0".

**QuDEX (Qualitative Data Exchange Format)** [Edit](#)  
The QuDEX standard/schema is a software-neutral format for qualitative data.

**SDMX (Statistical Data and Metadata Exchange)** [Edit](#)  
A set of common technical and statistical standards and guidelines to be used by sponsoring institutions.

**Repository-Developed Metadata Schemas** [Edit](#)  
Some repositories have decided that current standards do not fit their metadata needs, and so have created their own requirements.

**UKEOF** [Edit](#)  
A metadata standard for describing environmental monitoring activities, programmes, networks and facilities published by the UK Environmental Observation Framework (UKEOF).

## Life Sciences

**ABCD (Access to Biological Collection Data)** [Edit](#)  
The Access to Biological Collections Data (ABCD) Schema is a free-text standard that can be accommodated.

**Darwin Core** [Edit](#)  
A body of standards, including a glossary of terms (in both English and Spanish), and a set of data standards.

**EML (Ecological Metadata Language)** [Edit](#)  
Ecological Metadata Language (EML) is a metadata schema for describing ecological data.

**CSMD (Core Scientific Metadata)** [Edit](#)  
A standard for encoding archival descriptions.

**DDI (Data Documentation Initiative)** [Edit](#)  
A widely used, international standard for describing data from the social, behavioral, and economic sciences.

**MIDAS-Heritage** [Edit](#)  
A British cultural heritage standard for recording information on buildings, archaeological sites, and other cultural heritage.

**NeXus** [Edit](#)  
A metadata standard for describing environmental monitoring activities, programmes, networks and facilities published by the UK Environmental Observation Framework (UKEOF).

**OAI-ORE (Open Archives Initiative Object Reuse and Exchange)** [Edit](#)  
The goal of these standards is to expose the rich content in aggregations of Web resources to applications that support authoring, depositing and retrieving content from popular social networks of "Web 2.0".

**QuDEX (Qualitative Data Exchange Format)** [Edit](#)  
The QuDEX standard/schema is a software-neutral format for qualitative data.

**SDMX (Statistical Data and Metadata Exchange)** [Edit](#)  
A set of common technical and statistical standards and guidelines to be used by sponsoring institutions.

**Repository-Developed Metadata Schemas** [Edit](#)  
Some repositories have decided that current standards do not fit their metadata needs, and so have created their own requirements.

**UKEOF** [Edit](#)  
A metadata standard for describing environmental monitoring activities, programmes, networks and facilities published by the UK Environmental Observation Framework (UKEOF).

# Metadata Standards



## General Research Data

[CERIF \(Common European Research Information Format\)](#)

The Common European Research Information Format is the standard that the EU recommends to its member states for recording information.

[Data Package](#)

The Data Package specification is a generic wrapper format for exchanging data. Although it supports arbitrary metadata, the format defines a set of mandatory metadata that must be registered with the DataCite Metadata Store when minting a DOI persistent identifier for a data package.

A separate but linked specification provides a way to describe the columns of a data table; descriptions of this form can be included directly in the data package.

[DataCite Metadata Schema](#)

A set of mandatory metadata that must be registered with the DataCite Metadata Store when minting a DOI persistent identifier for a data package.

Sponsored by the DataCite consortium, version 3.0 was recently released in 2013.

[DCAT \(Data Catalog Vocabulary\)](#)

By using DCAT to describe datasets in data catalogs, publishers increase discoverability and enable applications easily to consume metadata.

[Dublin Core](#)

A basic, domain-agnostic standard which can be easily understood and implemented, and as such is one of the best known and most widely used.

Sponsored by the Dublin Core Metadata Initiative, Dublin Core was published as ISO Standard 15836 in February 2009.

[OAI-ORE \(Open Archives Initiative Object Reuse and Exchange\)](#)

The goal of these standards is to expose the rich content in aggregations of Web resources to applications that support authoring, depositing, and retrieving content from popular social networks of "Web 2.0".

[Observations and Measurements](#)

This encoding is an essential dependency for the OGC Sensor Observation Service (SOS) Interface Standard. More specifically, this standard defines the metadata for the observation and measurement data.

[PREMIS](#)

The PREMIS (Preservation Metadata: Implementation Strategies) Data Dictionary defines a set of metadata that most repositories of digital objects use. The influence of the creation of local application profiles, an XML Schema is provided to allow the metadata to be serialized independently.

PREMIS was initially developed by the Preservation Metadata: Implementation Strategies Working Group, convened by OCLC and RLG.

[PROV](#)

Provenance is information about entities, activities, and people involved in producing a piece of data or thing, which can be used to form a chain of custody.

[RDF Data Cube Vocabulary](#)

The standard provides a means to publish multi-dimensional data, such as statistics, on the web in such a way that it can be linked to related data and metadata among organizations.

[Repository-Developed Metadata Schemas](#)

Some repositories have decided that current standards do not fit their metadata needs, and so have created their own requirements.

## Arts and Humanities

[Encoded Archival Description \(EAD\)](#)

A standard for e

[DDI \(Data Documentation Initiative\)](#)

A widely used, it

- DDI Code
- DDI Lifecycle

Both versions are

[MIDAS-Heritage](#)

A British cultural

Sponsored by th

[OAI-ORE \(Open Archives Initiative\)](#)

The goal of these

popular social ne

## Engineering

[CIF \(Crystallographic Information File\)](#)

A well-established stand

Sponsored by the Intern

[CSMD \(Core Scientific Metadata\)](#)

A s

## Social ar

[ISA-TI \(International Standard Archival\)](#)

A widely used, it

- DDI Code
- DDI Lifecycle

[MIBBI \(Metadata Interchange\)](#)

Both versions are

[MIDAS-Heritage](#)

A British cultural

Sponsored by th

[OAI-ORE \(Open Archives Initiative\)](#)

The goal of these

popular social ne

[QuDEX \(Qualitative Data Exchange\)](#)

The QuDEX stan

[SDMX \(Statistical Data and Metadata eXchange\)](#)

A set of common

Sponsoring insti



# Metadata Standards



## ● Model extensions/integration

### Dublin Core

A basic, domain-agnostic standard which can be easily understood and implemented, and as such is c

Sponsored by the Dublin Core Metadata Initiative, Dublin Core was published as ISO Standard 15836

#### Summary Edit

##### Standard Website

<http://dublincore.org>

##### Specification

<http://dublincore.org/specifications/>

##### Related Vocabularies

[DCMI Vocabulary Managment Community](#)

##### Mappings

[UK AGMAP \(Academic Geospatial Metadata Application Profile\)](#)

[DataCite Metadata Schema](#)

[PROV](#)

[DDI \(Data Documentation Initiative\)](#)

[MARC \(Machine-Readable Cataloging\)](#)

##### Subjects

[General Research Data](#)

##### Disciplines

[Multi-disciplinary](#)

### Extensions Add

#### [AGLS Metadata Profile](#) Edit

An application of [Dublin Core](#) designed to improve visibility and availability of online resources, orig

#### [AGRIS Application Profile](#) Edit

A metadata standard drawing on [Dublin Core](#) and [AgMES](#) created specifically to enhance the desc

#### [ANZLIC Metadata Profile](#) Edit

A profile of [ISO 19115](#), also mapping to the AGLS profile of [Dublin Core](#), designed to facilitate effici

#### [Dryad Metadata Application Profile](#) Edit

An application profile based on the [Dublin Core](#) Metadata Initiative Abstract Model, used to describ

#### [eBank UK Metadata Application Profile](#) Edit

A [Dublin Core](#) Metadata Application Profile created for the eBank UK project, which provides acces

#### [OpenAIRE Guidelines for publication repositories, data archives and CRIS systems](#) Edit

The OpenAIRE Guidelines are a suite of application profiles designed to allow research institutions the OAI-PMH metadata harvesting protocol:

- The OpenAIRE Guidelines for Literature Repositories are based on [Dublin Core](#);
- The OpenAIRE Guidelines for Data Archives are based on the [DataCite Metadata Schema](#);
- The OpenAIRE Guidelines for CRIS Managers is based on [CERIF](#).

While the focus of each profile is different, they allow for interlinking and the contextualization of res

#### [Resource Metadata for the Virtual Observatory](#) Edit

Defines metadata terms and concepts necessary for discovery and use of astronomical data collec

The extension is based on Dublin Core, but with astronomy-specific extensions. Resource Metadat and maintained by IVOA Resource Registry Working Group and NVO Metadata Working Group

# Dublin Core - History



- Formally standardized in 1995 in an invitational workshop in Dublin.

"Dublin" refers to Dublin, Ohio, USA where the schema originated during the 1995 invitational OCLC/NCSA Metadata Workshop, hosted by the OCLC (known at that time as Online Computer Library Center), a library consortium based in Dublin, and the National Center for Supercomputing Applications (NCSA).

"Core" refers to the metadata terms as "broad and generic being usable for describing a wide range of resources". T

- The semantics of Dublin Core were established and are maintained by an international, cross-disciplinary group of professionals from librarianship, computer science, text encoding, museums, and other related fields of scholarship and practice.
- In 1999, the first Dublin Core encoding standard was in HTML.
- Starting in 2000, the Dublin Core community focused on "application profiles" – the idea that metadata records would use Dublin Core together with other specialized vocabularies to meet particular implementation requirements. During that time, the World Wide Web Consortium's work on a generic data model for metadata, the Resource Description Framework (RDF), was maturing. As part of an extended set of DCMI metadata terms, Dublin Core became one of the most popular vocabularies for use with RDF, more recently in the context of the linked data movement.
- The [Dublin Core Metadata Initiative \(DCMI\)](#) provides an open forum for the development of interoperable online metadata standards for a broad range of purposes and of business models.
- In 2008, DCMI separated from OCLC and incorporated as an independent entity.
- Currently, any and all changes that are made to the Dublin Core standard, are reviewed by a DCMI Usage Board within the context of a DCMI Namespace Policy (DCMI-NAMESPACE). This policy describes how terms are assigned and also sets limits on the amount of editorial changes allowed to the labels, definitions, and usage comments.

[https://en.wikipedia.org/wiki/Dublin\\_Core](https://en.wikipedia.org/wiki/Dublin_Core)

# Dublin Core - Definition



- The Dublin Core, also known as the Dublin Core Metadata Element Set, is a set of fifteen "core" elements (properties) for describing resources.
  - "core" because its elements are broad and generic, usable for describing a wide range of resources
  - not anymore only electronic
  - one of the top metadata vocabularies on the web
- Dublin Core Metadata Element Set, Version 1.1  
<http://dublincore.org/documents/dces/>

# Dublin Core – Original Elements



- 15 generic elements for describing resources
  - Contributor – "An entity responsible for making contributions to the resource".
  - Coverage – "The spatial or temporal topic of the resource, the spatial applicability of the resource, or the jurisdiction under which the resource is relevant".
  - Creator – "An entity primarily responsible for making the resource".
  - Date – "A point or period of time associated with an event in the lifecycle of the resource".
  - Description – "An account of the resource".
  - Format – "The file format, physical medium, or dimensions of the resource".
  - Identifier – "An unambiguous reference to the resource within a given context".
  - Language – "A language of the resource".
  - Publisher – "An entity responsible for making the resource available".
  - Relation – "A related resource".
  - Rights – "Information about rights held in and over the resource".
  - Source – "A related resource from which the described resource is derived".
  - Subject – "The topic of the resource".
  - Title – "A name given to the resource". <https://www.dublincore.org/specifications/dublin-core/dces/>
  - Type – "The nature or genre of the resource".

# Dublin Core – Evolution



- The core properties are part of a larger set of DCMI Metadata Terms.
- Later formally standardized. “Later” because no semantic web (and, e.g., RDF) was available at the time
- Current version

Refers to a set of metadata vocabularies and technical specifications  
Maintained by the Dublin Core Metadata Initiative (DCMI)

The full set of vocabularies includes sets of resource classes,  
vocabulary encoding schemes, and syntax encoding schemes

The terms in DCMI vocabularies are intended to be used in  
combination with terms from other, compatible vocabularies in the  
context of application profiles and on the basis of the DCMI  
Abstract Model [DCAM].



- Semantic web evolution
  - DCMI includes formal domains and ranges in the definitions of its properties
  - not to affect the conformance of existing implementations of "simple Dublin Core"
    - domains and ranges have not been specified for the "initial" fifteen properties
      - namespace dc:
        - <http://purl.org/dc/elements/1.1/>
    - fifteen new properties with "names" identical to those of the Dublin Core Metadata Element Set Version 1.1 have been created
      - new namespace dcterms:
        - <http://purl.org/dc/terms/>
    - These fifteen new properties have been defined as subproperties of the corresponding properties of DCES Version 1.1 and assigned domains and ranges

# Dublin Core – Evolution



- Dublin Core Metadata Element Set, Version 1.1

<https://www.dublincore.org/specifications/dublin-core/dces/>



<https://www.dublincore.org/specifications/dublin-core/dcmi-terms/>

# Dublin Core – semantic evolution



## ● Property

<b>Term Name:</b> contributor	
<b>URI:</b>	<a href="http://purl.org/dc/elements/1.1/contributor">http://purl.org/dc/elements/1.1/contributor</a>
<b>Label:</b>	Contributor
<b>Definition:</b>	An entity responsible for making contributions to the resource.
<b>Comment:</b>	Examples of a Contributor include a person, an organization, or a service. Typically, the name of a Contributor should be used to indicate the entity.

## ● Term

<b>Term Name:</b> contributor	
<b>URI:</b>	<a href="http://purl.org/dc/terms/contributor">http://purl.org/dc/terms/contributor</a>
<b>Label:</b>	Contributor
<b>Definition:</b>	An entity responsible for making contributions to the resource.
<b>Comment:</b>	Examples of a Contributor include a person, an organization, or a service.
<b>Type of Term:</b>	<a href="#">Property</a>
<b>Refines:</b>	<a href="http://purl.org/dc/elements/1.1/contributor">http://purl.org/dc/elements/1.1/contributor</a>
<b>Has Range:</b>	<a href="http://purl.org/dc/terms/Agent">http://purl.org/dc/terms/Agent</a>
<b>Version:</b>	<a href="http://dublincore.org/usage/terms/history/#contributorT-001">http://dublincore.org/usage/terms/history/#contributorT-001</a>



# Dublin Core – OAI-PMH usage



## ● OAI-PMH

- Open Archives Initiative Protocol for Metadata Harvesting
- application-independent interoperability framework based on metadata harvesting
- two classes of participants
  - Data Providers support OAI-PMH as a means of exposing metadata
  - Service Providers harvest metadata via the OAI-PMH
    - for building value-added services
    - Harvest: issue OAI-PMH requests

## ● OAI-PMH supports the dissemination of records in multiple metadata formats from a repository

- metadataPrefix specifies the format to be used to reply to a request

it is used together with the requests methods available in the protocol: ListRecords, ListIdentifiers, and GetRecord to retrieve the records, or the headers of the records that include metadata in the format specified by the metadataPrefix

- For purposes of interoperability, repositories must disseminate Dublin Core, without any qualification
  - metadataPrefix “oai\_dc” reserved
  - XML namespace URI → [http://www.openarchives.org/OAI/2.0/oai\\_dc/](http://www.openarchives.org/OAI/2.0/oai_dc/)
  - URL → [http://www.openarchives.org/OAI/2.0/oai\\_dc.xsd](http://www.openarchives.org/OAI/2.0/oai_dc.xsd).

# Dublin Core – OAI-PMH validation schema



- OAI-PMH

- Open Archives Initiative

- application

- two classes

- Data

- Services

- OAI-PMH schema
- from a repository

- metadata requests in different formats

- For purposes of any quality

- metadata

- XML

- URL

A XML schema for validating Unqualified Dublin Core metadata associated with the reserved oai\_dc metadataPrefix

```
<schema targetNamespace="http://www.openarchives.org/OAI/2.0/oai_dc/"
  xmlns:oai_dc="http://www.openarchives.org/OAI/2.0/oai_dc/"
  xmlns:dc="http://purl.org/dc/elements/1.1/"
  xmlns="http://www.w3.org/2001/XMLSchema"
  elementFormDefault="qualified" attributeFormDefault="unqualified">

  <annotation>
    <documentation>
      XML Schema 2002-03-18 by Pete Johnston.
      Adjusted for usage in the OAI-PMH.
      Schema imports the Dublin Core elements from the DCMI schema for unqualified Dublin Core.
      2002-12-19 updated to use simpledc20021212.xsd (instead of simpledc20020312.xsd)
    </documentation>
  </annotation>

  <import namespace="http://purl.org/dc/elements/1.1/"
    schemaLocation="http://dublincore.org/schemas/xmls/simpledc20021212.xsd"/>

  <element name="dc" type="oai_dc:oai_dcType"/>

  <complexType name="oai_dcType">
    <choice minOccurs="0" maxOccurs="unbounded">
      <element ref="dc:title"/>
      <element ref="dc:creator"/>
      <element ref="dc:subject"/>
      <element ref="dc:description"/>
      <element ref="dc:publisher"/>
      <element ref="dc:contributor"/>
      <element ref="dc:date"/>
      <element ref="dc:type"/>
      <element ref="dc:format"/>
      <element ref="dc:identifier"/>
      <element ref="dc:source"/>
      <element ref="dc:language"/>
      <element ref="dc:relation"/>
      <element ref="dc:coverage"/>
      <element ref="dc:rights"/>
    </choice>
  </complexType>

</schema>
```

This Schema is available at [http://www.openarchives.org/OAI/2.0/oai\\_dc.xsd](http://www.openarchives.org/OAI/2.0/oai_dc.xsd)

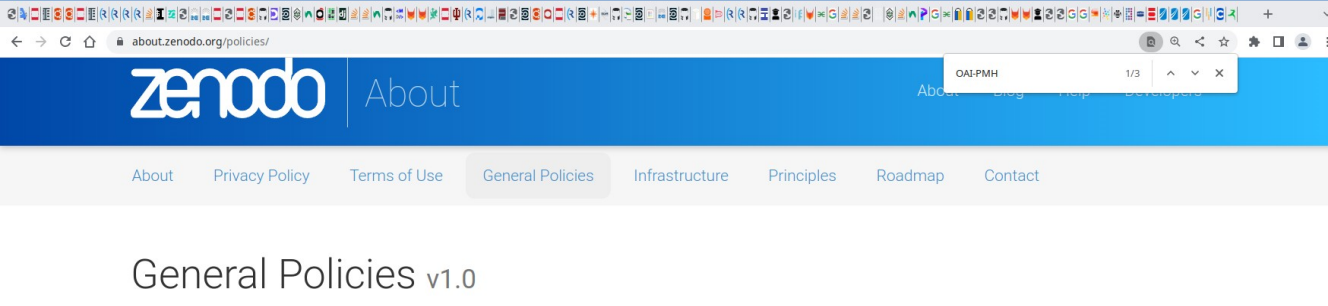
a harvesting

ata formats

d GetRecord  
etadata in the

Core, without

# OAI-PMH & Zenodo



## General Policies v1.0

### Content

- **Scope:** All file agreements f
- **Status of res**
- **Eligible depo**
- **Ownership:** B the property c
- **Data file form** preservation
- **Volume and :**
- **Data quality:** harmless in c
- **Metadata typ:** several stand
- **Language:** Fc
- **Licenses:** Us

### Access and Reuse

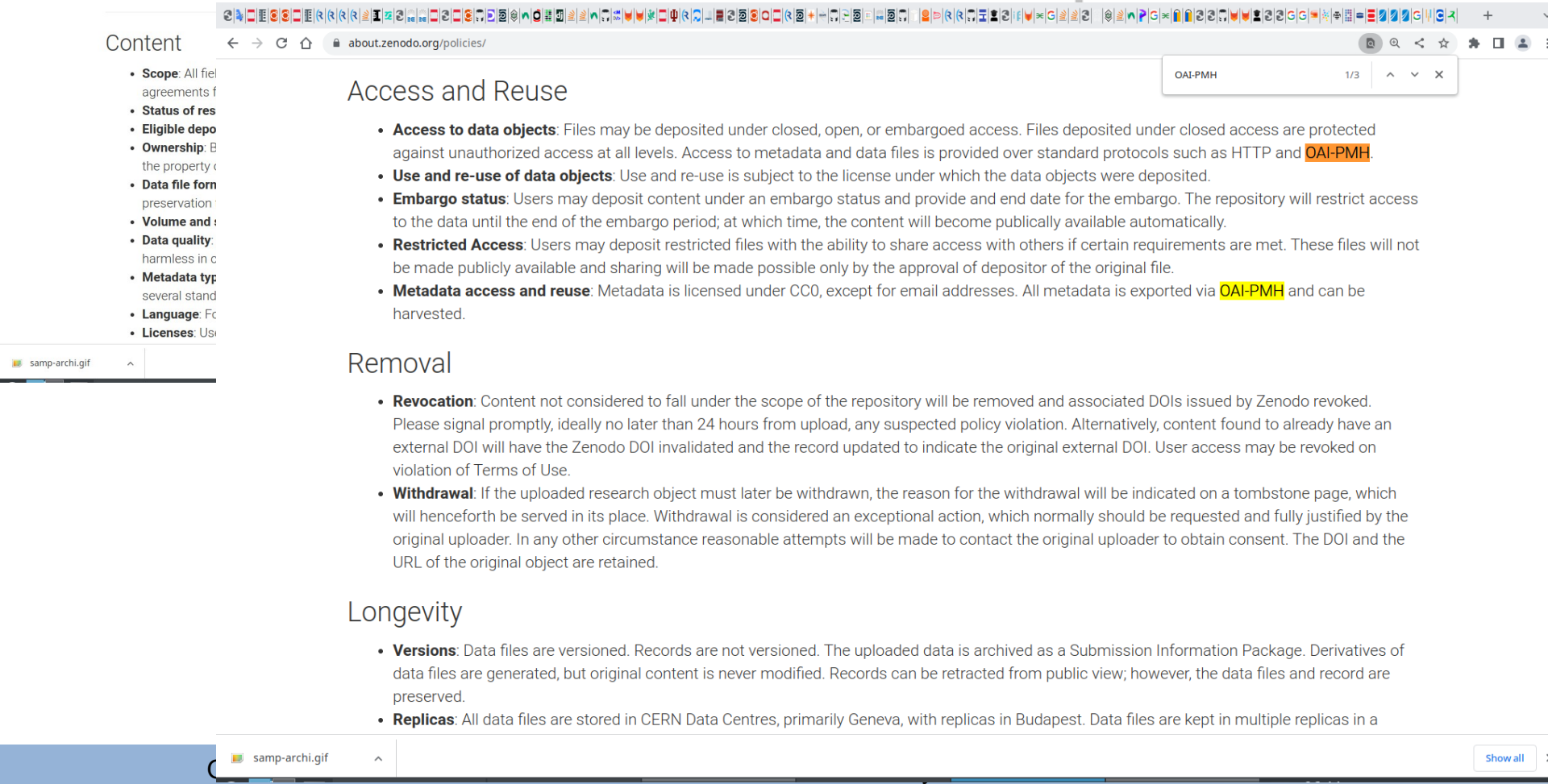
- **Access to data objects:** Files may be deposited under closed, open, or embargoed access. Files deposited under closed access are protected against unauthorized access at all levels. Access to metadata and data files is provided over standard protocols such as HTTP and **OAI-PMH**.
- **Use and re-use of data objects:** Use and re-use is subject to the license under which the data objects were deposited.
- **Embargo status:** Users may deposit content under an embargo status and provide an end date for the embargo. The repository will restrict access to the data until the end of the embargo period; at which time, the content will become publicly available automatically.
- **Restricted Access:** Users may deposit restricted files with the ability to share access with others if certain requirements are met. These files will not be made publicly available and sharing will be made possible only by the approval of depositor of the original file.
- **Metadata access and reuse:** Metadata is licensed under CC0, except for email addresses. All metadata is exported via **OAI-PMH** and can be harvested.

### Removal

- **Revocation:** Content not considered to fall under the scope of the repository will be removed and associated DOIs issued by Zenodo revoked. Please signal promptly, ideally no later than 24 hours from upload, any suspected policy violation. Alternatively, content found to already have an external DOI will have the Zenodo DOI invalidated and the record updated to indicate the original external DOI. User access may be revoked on violation of Terms of Use.
- **Withdrawal:** If the uploaded research object must later be withdrawn, the reason for the withdrawal will be indicated on a tombstone page, which will henceforth be served in its place. Withdrawal is considered an exceptional action, which normally should be requested and fully justified by the original uploader. In any other circumstance reasonable attempts will be made to contact the original uploader to obtain consent. The DOI and the URL of the original object are retained.

### Longevity

- **Versions:** Data files are versioned. Records are not versioned. The uploaded data is archived as a Submission Information Package. Derivatives of data files are generated, but original content is never modified. Records can be retracted from public view; however, the data files and record are preserved.
- **Replicas:** All data files are stored in CERN Data Centres, primarily Geneva, with replicas in Budapest. Data files are kept in multiple replicas in a



# DCAT (W3C technology)



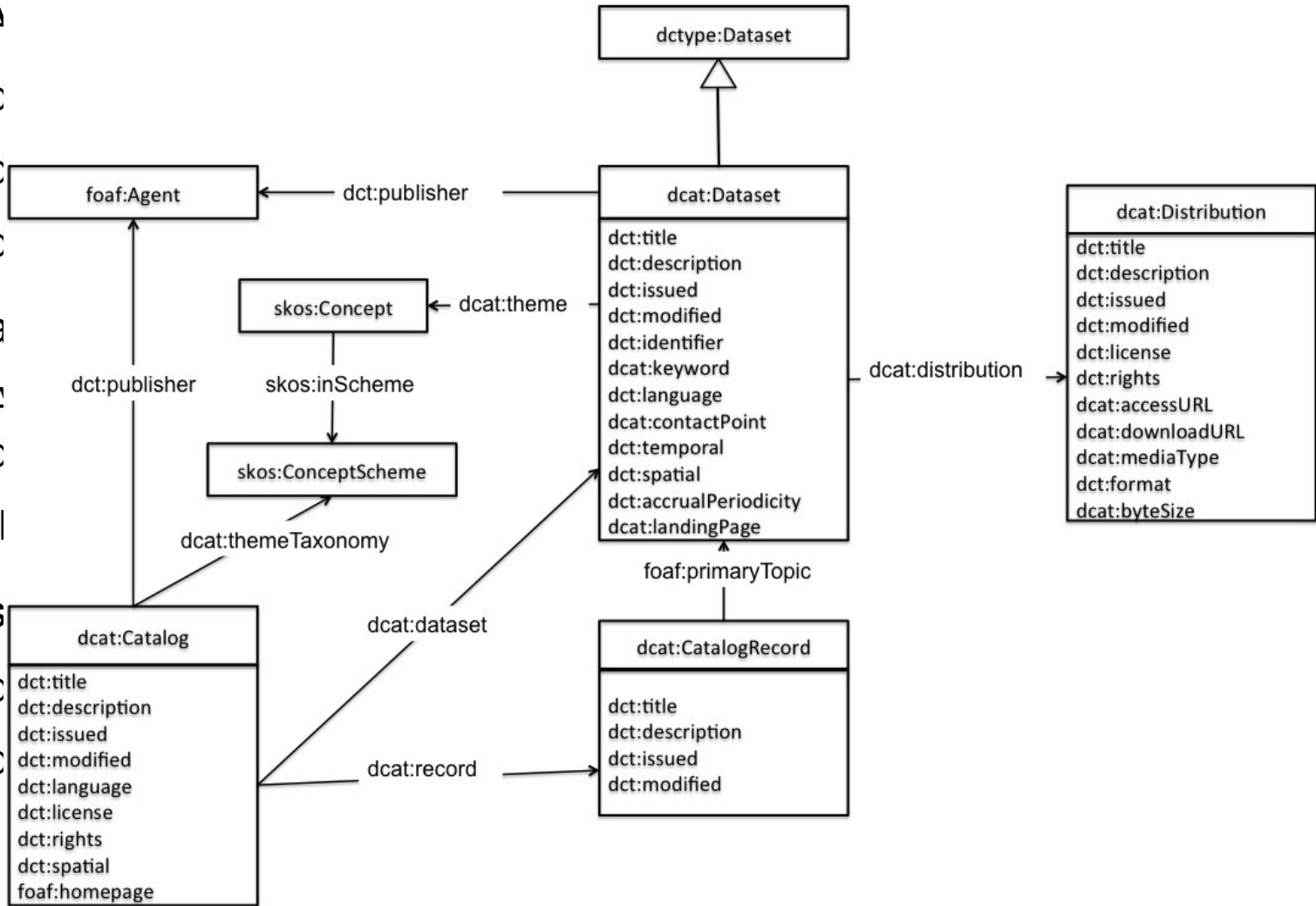
- DCAT is an RDF vocabulary **designed** to facilitate interoperability between **data catalogs** published on the Web
- Using DCAT to describe datasets in data catalogs, publishers increase discoverability and enable applications easily to consume metadata from multiple catalogs
- DCAT does not make any assumptions about the format of the datasets described in a catalog
  - Other, complementary vocabularies may be used together with DCAT to provide more detailed format-specific information
- <https://www.w3.org/TR/vocab-dcat/>
  - <https://www.w3.org/TR/vocab-dcat-2/>

- DCAT defines three main classes:
  - dcat:Catalog represents the catalog
  - dcat:Dataset represents a dataset in a catalog.
  - dcat:Distribution represents an accessible form of a dataset
- A dataset does not have to be available as a downloadable file.
  - A dataset that is available via an API can be defined as an instance of dcat:Dataset
  - the API can be defined as an instance of dcat:Distribution
- Class dcat:CatalogRecord describes a dataset entry in the catalog
  - dcat:Dataset represents the dataset itself
  - dcat:CatalogRecord represents the record that describes a dataset in the catalog
    - is optional
    - is used to capture provenance information

# DCAT (W3C): Metadata Model



- DCA
  - C
  - C
  - C
- A da
  - A
  - C
  - t
- Clas
  - C
  - C



catalog

Prefix dct: namespace to Dublin Core elements 1.1

# DataCite Metadata Model



- The DataCite Metadata Schema
  - list of core metadata properties chosen for an accurate and consistent identification of a resource for citation and retrieval purposes
  - The **resource** that is being identified can be **of any kind**, but it is typically a dataset
  - term ‘dataset’: its broadest sense
- Collaborate with the Dublin Core Metadata Initiative (DCMI) to maintain a Dublin Core Application Profile for the schema
- Presents 3 different levels of obligation for the metadata properties
  - **Mandatory** (M) properties must be provided
  - **Recommended** (R ) properties are optional, but strongly recommended for interoperability
  - **Optional** (O) properties are optional and provide richer description
- <http://schema.datacite.org/meta/kernel-4.3/>
  - v.4.1: explicit changes to improve software citation
  - v.4.2: alternate identifiers and machine readable licences
  - v.4.3: “funder” identifiers...

# DataCite Metadata – Properties



Table 1: DataCite Mandatory Properties

<i>ID</i>	<i>Property</i>	<i>Obligation</i>
1	Identifier (with mandatory type sub-property)	M
2	Creator (with optional given name, family name, name identifier and affiliation sub-properties)	M
3	Title (with optional type sub-properties)	M
4	Publisher	M
5	PublicationYear	M
10	ResourceType (with mandatory general type description sub-property)	M

Table 2: DataCite Recommended and Optional Properties

<i>ID</i>	<i>Property</i>	<i>Obligation</i>
6	Subject (with scheme sub-property)	R
7	Contributor (with optional given name, family name, name identifier and affiliation sub-properties)	R
8	Date (with type sub-property)	R
9	Language	O
11	AlternateIdentifier (with type sub-property)	O
12	RelatedIdentifier (with type and relation type sub-properties)	R
13	Size	O
14	Format	O
15	Version	O
16	Rights	O
17	Description (with type sub-property)	R
18	GeoLocation (with point, box and polygon sub-properties)	R
19	FundingReference (with name, identifier, and award related sub-properties)	O

- Among Recommended

- Description (17) is considered the most important
- Especially in connected usage with the Recommended sub-property
  - descriptionType = "Abstract"



# DataCite Metadata – Mandatory+SubProperties



ID	DataCite-Property	Occ	Definition	Allowed values, examples, other constraints
1	Identifier	1	The Identifier is a unique string that identifies a resource. For software, determine whether the identifier is for a specific version of a piece of software, (per the Force11 Software Citation Principles <sup>13</sup> ), or for all versions.	DOI (Digital Object Identifier) registered by a DataCite member. Format should be "10.1234/foo"
1.1	identifierType	1	The type of Identifier.	<i>Controlled List Value:</i> DOI
2	Creator	1-n	The main researchers involved in producing the data, or the authors of the publication, in priority order. To supply multiple creators, repeat this property.	May be a corporate/institutional or personal name. Note: DataCite infrastructure supports up to 8000-10000 names. For name lists above that size, consider attribution via linking to the related metadata.
2.1	creatorName	1	The full name of the creator.	Examples: Charpy, Antoine; Foo Data Center  Note: The personal name, format should be: family, given. Non-roman names may be transliterated according to the ALA-LC schemas <sup>14</sup> .
2.1.1	nameType	0-1	The type of name	<i>Controlled List Values:</i> Organizational Personal

ID	DataCite-Property	Occ	Definition	Allowed values, examples, other constraints
2.2	givenName	0-1	The personal or first name of the creator.	Examples based on the 2.1 names: Antoine; Mae
2.3	familyName	0-1	The surname or last name of the creator.	Examples based on the 2.1 names: Charpy; Jemison
2.4	nameIdentifier	0-n	Uniquely identifies an individual or legal entity, according to various schemas.	The format is dependent upon schema.
2.4.1	nameIdentifierScheme	1	The name of the name identifier schema.	If nameIdentifier is used, nameIdentifierScheme is mandatory.  Examples: ORCID <sup>15</sup> , ISNI <sup>16</sup>
2.4.2	schemeURI	0-1	The URI of the name identifier schema.	Examples: <a href="http://www.isni.org">http://www.isni.org</a> <a href="http://orcid.org">http://orcid.org</a>
2.5	affiliation	0-n	The organizational or institutional affiliation of the creator.	Free text.
3	Title	1-n	A name or title by which a resource is known. May be the title of a dataset or the name of a piece of software.	Free text.
3.1	titleType	0-1	The type of Title.	<i>Controlled List Values:</i> AlternativeTitle Subtitle TranslatedTitle  Other

● Properties 4,5 have occurrence 1 (being mandatory) without mandatory sub-properties

● Property 10 has mandatory resourceTypeGeneral sub-property, with values in a controlled list:

- Audiovisual, Collection, DataPaper, Dataset, Event, Image, InteractiveResource, Model, PhysicalObject, Service, Software, Sound, Text, Workflow, Other

- ...some details
- Most Recommended/Optional properties and sub-properties
  - Have values within controlled list vocabularies
    - 7 Contributor [0-n]: Free text and optional
      - 7.1 contributorType [1]: controlled list
        - ContactPerson, DataCollector, DataCurator, DataManager, Distributor, Editor, HostingInstitution, Producer, ProjectLeader, ProjectManager, ProjectMember, RegistrationAgency, RegistrationAuthority, RelatedPerson, Researcher, ResearchGroup, RightsHolder, Sponsor, Supervisor, WorkPackageLeader, Other
  - Specify free text values through (optional) schema & value URI identifiers
    - 6 Subject [0-n]: Free text
      - 6.1 subjectScheme [0-1] The name of the subject scheme: Free text
      - 6.2 schemeURI [0-1] The URI of the subject identifier scheme
      - 6.3 valueURI [0-1] The URI of the subject term
  - Point to external standard formats, models, schemas, ...
    - 9 Language [0-1]: allowed values from IETF BCP 47, ISO 639-1 language codes
      - Examples: en, de, fr

- Metadata expressed through XSD documents (and associated Recommendation documents)
  - “Resource Metadata” describes the basic concepts
  - VOResource brings it to XSD and provides a technical entry point
    - Multiple extensions follow: standards, simple access protocols, collections and services, ...
  - Connected interfaces and identifiers specifications

ReR	VOA Identifiers	2.0		2.0 2.0 2.0 2.0 1.12 1.11 1.10 1.10 1.10 1.00
	VOA Registry Interfaces	1.1		1.1 1.1 1.1 1.1 1.1 1.1 1.0 1.0 1.00 1.02 1.01 1.00
	RM - Resource Metadata for the Virtual Observatory	1.12		1.12 1.12 1.10 1.10 1.01 1.01 1.00 1.00
	StandardsRegExt: a VOResource Schema Extension for Describing IVOA Standards	1.0		1.0 1.0 1.0 1.0 1.0 1.0 1.0
	SimpleDALRegExt - Describing Simple Data Access Services	1.1		1.1 1.1 1.1 1.1 1.0 1.0 1.0 1.0 1.0
	VOResource - an XML Encoding Schema for Resource Metadata	1.1		1.1 1.1 1.1 1.1 1.1 1.03 1.02 1.02 1.01 1.00
	VODataService - A VOResource Schema Extension for Describing Collections and Services	1.1	1.2	1.2 1.1 1.1 1.1 1.1 1.1 1.10
	RegTAP - Registry Relational Schema	1.0	1.1	1.1 1.1 1.0 1.0 1.0 1.0 1.0 1.0 1.0

- <http://ivoa.net/documents/> (ReR section in the table)
  - But TAPRegExt in the DAL part...
  - (future) maybe protocol dedicated extensions will end up in the protocol document itself

# Resource Metadata for the VO



- Starts out from FITS (Flexible Image Transport System) usage scenario
- General concepts are or map directly Dublin Core
  - The harvesting interface to the Registry is OAI-PMH
- Hierarchical system for metadata management
  - Lower levels provide more extensive and complex metadata
    - description of query syntax, access protocols, and usage policies
- Basic concepts
  - Resource is a general term
    - Described in terms of who curates or maintains it
    - Can be given a name and a unique identifier
  - Organisation is specific type of resource that brings people together to pursue participation in VO applications
    - Can be hierarchical and range greatly in size and scope
      - University, observatory, or government agency, ..., scientific project, space mission, or individual researcher
      - A provider is an organisation that makes data and/or services available to users over the network
  - Service is any VO resource that can be invoked by the user to perform some action on their behalf
    - Query service supports a query/response protocol
    - Non-query services: copy or delete files on remote files systems, mail information, kill existing jobs, authorize actions, ...
    - Registry is a query service for which the response is a structured description of resources
- Resource metadata include
  - Identity metadata (name, identifier, ... )
  - Curation metadata (who supports the resource, its availability, ... )
  - Content metadata (types of data, sky coverage, spectral coverage, ... )



Article Talk

Read Edit View history

Search Wikipedia



WIKIPEDIA The Free Encyclopedia

- Main page
- Contents
- Current events
- Random article
- About Wikipedia
- Contact us
- Donate

Contribute

- Help
- Learn to edit
- Community portal
- Recent changes
- Upload file

Tools

- What links here
- Related changes
- Special pages

Get citation

# FITS

From Wikipedia, the free encyclopedia

*For other uses, see [FITS \(disambiguation\)](#).*

**Flexible Image Transport System (FITS)** is an [open standard](#) defining a digital [file format](#) useful for storage, transmission and processing of data: formatted as multi-dimensional arrays (for example a 2D image), or tables.<sup>[3]</sup> FITS is the most commonly used digital [file format](#) in [astronomy](#). The FITS standard was designed specifically for astronomical data, and includes provisions such as describing [photometric](#) and spatial calibration information, together with image origin metadata.

The FITS format was first standardized in 1981;<sup>[4]</sup> it has evolved gradually since then, and the most recent version (4.0) was standardized in 2016. FITS was designed with an eye towards long-term archival storage, and the maxim *once FITS, always FITS* represents the requirement that developments to the format must be [backward compatible](#).

[Image metadata](#) is stored in a human-readable [ASCII](#) header. The information in this header is designed to calculate the byte offset of some information in the subsequent data unit to support direct access to the data cells. Each FITS file consists of one or more headers containing ASCII [card images](#)<sup>[a]</sup> that carry keyword/value pairs, interleaved between data blocks. The keyword/value pairs provide information such as size, origin, coordinates, binary data format, free-form comments, history of the data, and anything else the creator desires: while many keywords are reserved for FITS use, the standard allows arbitrary use of the rest of the name-space.

FITS is also often used to store non-image data, such as [spectra](#), [photon](#) lists, [data cubes](#), or [structured data](#) such as multi-table [databases](#). A FITS file may contain several extensions, and each of these may contain a data object. For example, it is possible to store [x-ray](#) and [infrared](#) exposures in the same file.

## FITS

<b>Filename extension</b>	<code>.fits</code> , <code>.fit</code> , <code>.fts</code>
<b>Internet media type</b>	<code>image/fits</code> <code>application/fits</code> <sup>[1]</sup>
<b>Developed by</b>	IAU FITS Working Group <sup>[2]</sup>
<b>Initial release</b>	1981; 41 years ago
<b>Latest release</b>	4.0 July 2016; 6 years ago
<b>Type of format</b>	<a href="#">image format</a> , <a href="#">structured data</a>
<b>Website</b>	<a href="https://fits.gsfc.nasa.gov">fits.gsfc.nasa.gov</a> <sup>↗</sup>

See: <https://fits.gsfc.nasa.gov/>  
<https://www.vaticanlibrary.va/it/il-patrimonio/il-progetto-di-digitalizzazione.html>  
<https://www.vaticanlibrary.va/it/il-patrimonio/fits-files.html>

# Resource Metadata - Structure



- Identity
  - Title, Shortname, Identifier (IVOID)
- Curation
  - Publisher (with PublisherID), Creator, Contributor
  - Date, Version
  - Contact
- General Content
  - Subject (controlled IAU vocabulary), Description, Source (Bibliographic reference), ReferenceURL, Type (controlled vocabulary), ContentLevel (target user), Relationship
- Collection & Service
  - Facility, Instrument
  - Coverage: spatial, spectral, bounds, resolution
  - UCD (Unified Content Descriptors), format, rights
  - Quality flags, validation, uncertainties
- Interface & Capabilities
  - Interface: BaseURL and other URLs
  - Capability: identified by a StandardID (IVOID)

- Specifies through XSD hierarchical structure of Resource Metadata
  - What's a timestamp?
    - vr:UTCTimestamp

```
<xs:simpleType name="UTCTimestamp" >
  <xs:restriction base="xs:dateTime" >
    <xs:pattern
      value="\d{4}-\d\d-\d\dT\d\d:\d\d:\d\d(\.\d+)?Z?" />
    </xs:restriction>
  </xs:simpleType>
```

- Relation among Interface and Capability elements

```
<capability xsi:type="ex:ExampleCapType"
  standardID="ivo://example.com/std/exampleAccess"
  xmlns:ex="http://ivoa.net/std/example-1.xsd">
  ...
</capability>

<capability>
  <interface xsi:type="vr:WebBrowser">
    <accessURL use="full"
      >http://example.org/browser-service</accessURL>
    </interface>
  </capability>
```

- Provide guidelines to extend the VOResource schema, when needed

# Standards Extensions



- Extend VOResource to add 3 resource types

- vstd:Standard
- vstd:ServiceStandard
- vstd:StandardKeyEnumeration

## vstd:Standard Type Schema Definition

```
<xs:complexType name="Standard" >
  <xs:complexContent >
    <xs:extension base="vr:Resource" >
      <xs:sequence >
        <xs:element name="endorsedVersion" type="vstd:EndorsedVersion"
          maxOccurs="unbounded" />
        <xs:element name="schema" type="vstd:Schema" minOccurs="0"
          maxOccurs="unbounded" >
        <xs:element name="deprecated" type="xs:token" minOccurs="0" />
        <xs:element name="key" type="vstd:StandardKey" minOccurs="0"
          maxOccurs="unbounded" />
      </sequence>
    </extension>
  </complexContent>
</complexType>
```

## An example of a Standard resource that summarizes this specification

```
<?xml version="1.0" encoding="UTF-8"?>
<ri:Resource xsi:type="vstd:Standard" status="active"
  created="2012-02-17T11:15:00" updated="2012-02-17T11:15:00"
  xmlns:ri="http://www.ivoa.net/xml/RegistryInterface/v1.0"
  xmlns:vstd="http://www.ivoa.net/xml/StandardsRegExt/v1.0"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance">
  <title> StandardsRegExt: a VOResource Schema Extension for Describing IVOA Standards </title>
  <shortName> StandardsRegExt </shortName>
  <identifier> ivo://ivoa.net/std/StandardsRegExt </identifier>
  <curation>
    .....
  </curation>
  <content>
    .....
  </content>
  <endorsedVersion status="pr"> 1.0 </endorsedVersion>
  <schema namespace="http://www.ivoa.net/xml/StandardsRegExt/v1.0">
    <location>http://www.ivoa.net/xml/StandardsRegExt/v1.0</location>
    <description>
      the VOResource extension XML Schema for registering standards
    </description>
    <example>http://rofr.ivoa.net/examples/StandardsRegExt.xml</example>
    <example>http://rofr.ivoa.net/examples/SIA.xml</example>
    <example>http://rofr.ivoa.net/examples/VOSpace.xml</example>
  </schema>
</ri:Resource>
```



# Observational Core Metadata Model

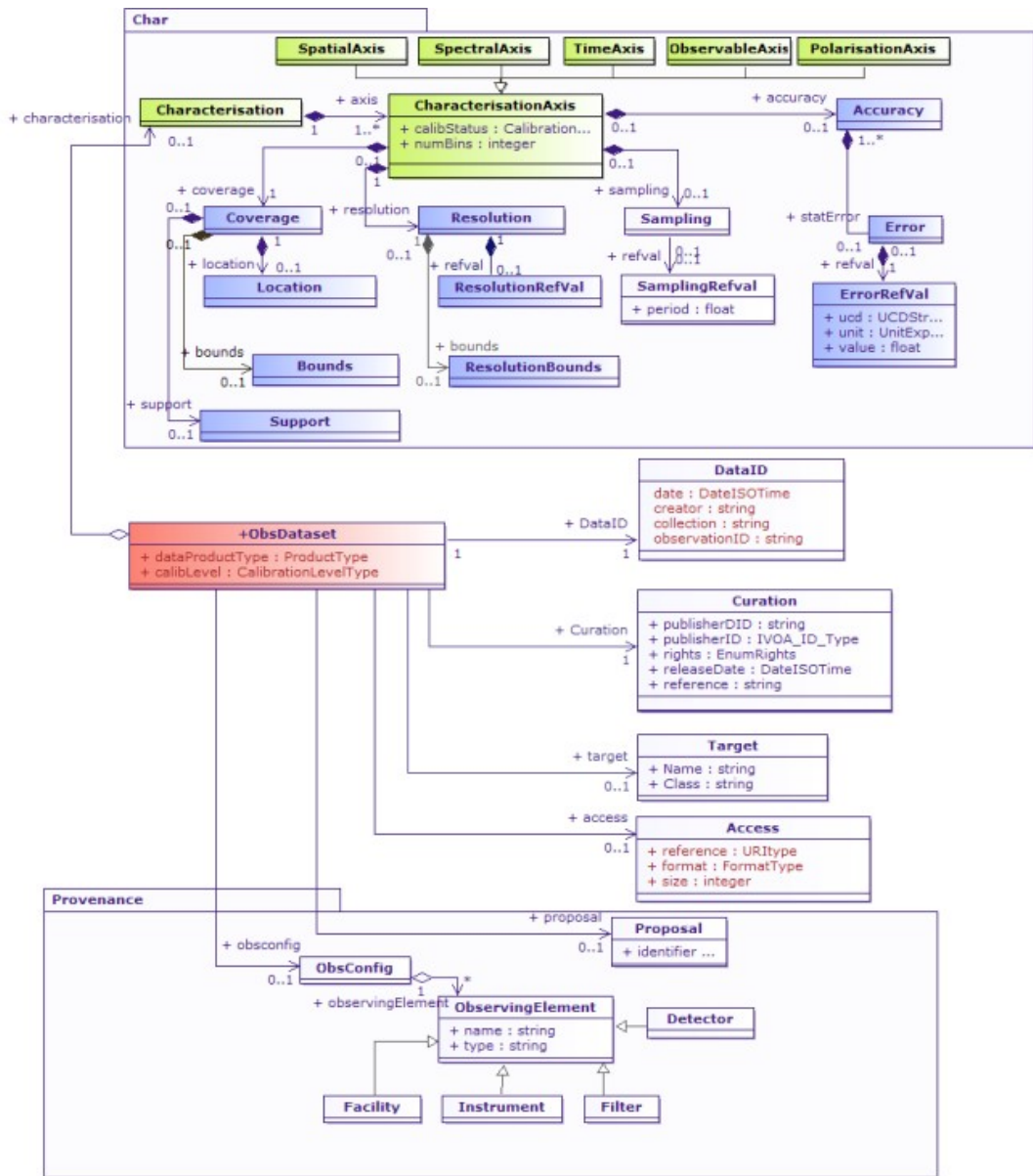


- Core components of the Observation data model necessary to perform **data discovery when querying data centers for astronomical observations** of interest
- Focus is on data discovery
  - A number of use-cases have been defined
  - Aimed at finding observational data products
  - Broadcasting the same query to multiple archives
    - **global data discoverability and accessibility**
- Need to give data providers a set of metadata attributes that they can easily map to their database system in order to support queries
- <http://www.ivoa.net/documents/ObsCore/20170509/REC-ObsCore-v1.1-20170509.pdf>

# Observational Core Metadata: UML



- Core compc discovery w/
- Focus is on
  - A number
  - Aimed at
  - Broadcas
  - global
- Need to give map to their
- <http://www.ivoa.i>



perform data  
operations of interest

they can easily

# ObsCore – Flat View



- Flat table approach
- Mandatory Structure but NULL-able values
  - Exceptions
    - `calib_level`, `obs_collection`, `obs_id`, `obs_publisher_did`
- Mandatory
  - Units
  - Data domain
  - Coordinate frames
- Comprehensive usage of
  - Vocabularies
  - Identifiers
- Limited number of mandatory elements
  - Optional standardized ones
  - Custom additions allowed

<i>Column Name</i>	<i>Unit</i>	<i>Type</i>	<i>Description</i>
<code>dataprodect_type</code>	unitless	String	Logical data product type (image etc.)
<code>calib_level</code>	unitless	enum integer	Calibration level {0, 1, 2, 3, 4}
<code>obs_collection</code>	unitless	String	Name of the data collection
<code>obs_id</code>	unitless	String	Observation ID
<code>obs_publisher_did</code>	unitless	String	Dataset identifier given by the publisher
<code>access_url</code>	unitless	String	URL used to access (download) dataset
<code>access_format</code>	unitless	String	File content format (see in App. BB.5.2)
<code>access_estsize</code>	kbyte	integer	Estimated size of dataset in kilo bytes
<code>target_name</code>	unitless	String	Astronomical object observed, if any
<code>s_ra</code>	deg	double	Central right ascension, ICRS
<code>s_dec</code>	deg	double	Central declination, ICRS
<code>s_fov</code>	deg	double	Diameter (bounds) of the covered region
<code>s_region</code>	unitless	String	Sky region covered by the data product (expressed in ICRS frame)
<code>s_xel1</code>	unitless	integer	Number of elements along the first spatial axis
<code>s_xel2</code>	unitless	integer	Number of elements along the second spatial axis
<code>s_resolution</code>	arcsec	double	Spatial resolution of data as FWHM
<code>t_min</code>	d	double	Start time in MJD
<code>t_max</code>	d	double	Stop time in MJD
<code>t_exptime</code>	s	double	Total exposure time
<code>t_resolution</code>	s	double	Temporal resolution FWHM
<code>t_xel</code>	unitless	integer	Number of elements along the time axis
<code>em_min</code>	m	double	Start in spectral coordinates
<code>em_max</code>	m	double	Stop in spectral coordinates
<code>em_res_power</code>	unitless	double	Spectral resolving power
<code>em_xel</code>	unitless	integer	Number of elements along the spectral axis
<code>o_ucd</code>	unitless	String	UCD of observable (e.g. phot.flux.density, phot.count, etc.)
<code>pol_states</code>	unitless	String	List of polarization states or NULL if not applicable
<code>pol_xel</code>	unitless	integer	Number of polarization samples
<code>facility_name</code>	unitless	String	Name of the facility used for this observation
<code>instrument_name</code>	unitless	String	Name of the instrument used for this observation

- Common Archive Observation Model, enables
  - **Storage of observational metadata** from the complete set of telescopic data
  - **Searching through that metadata** using a **single interface**
- The generalized capability of CAOM comes at the expense of some model complexity and the requirement of adopting a language that is unfamiliar to users
- To decrease the learning curve for users
  - expose CAOM via a simplified search web page interface
  - expose via a Table Access Protocol (TAP) web service
    - for users requiring access to more details of the observations and greater flexibility in query construction

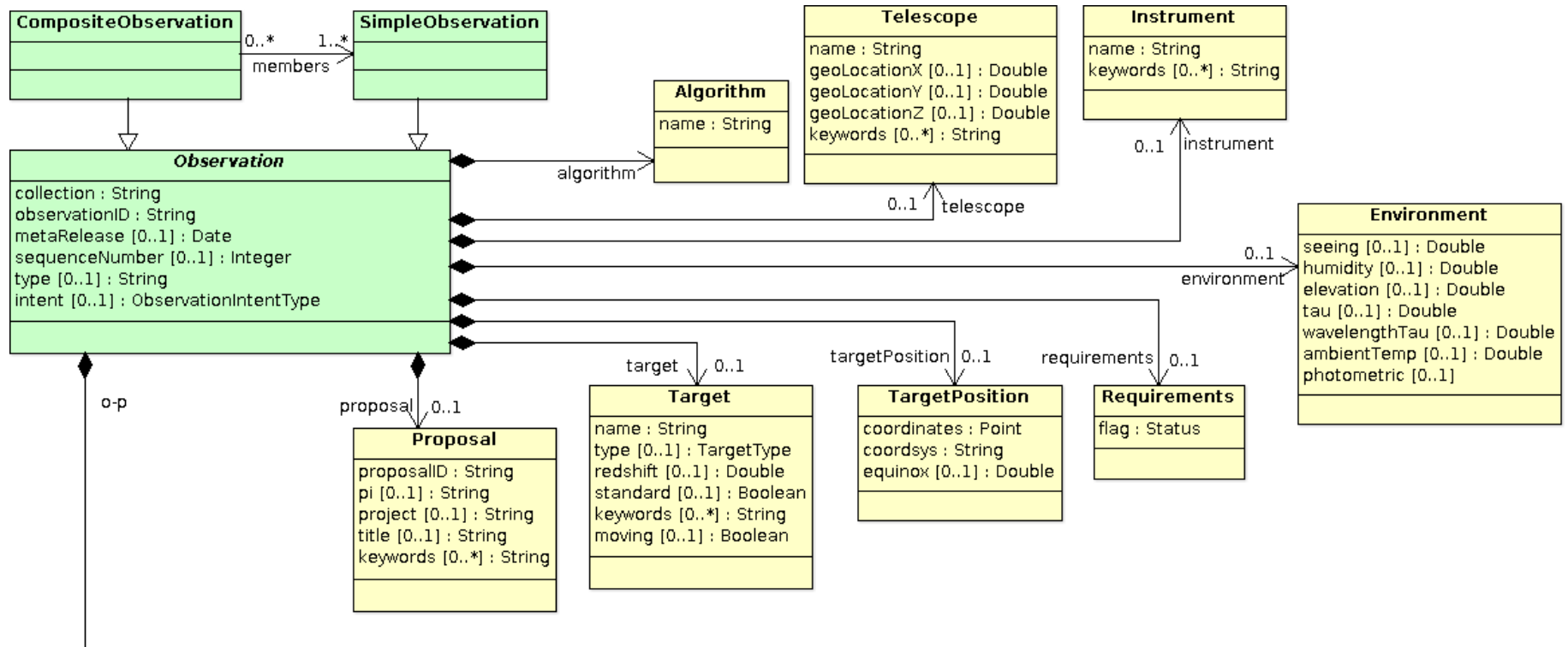
## ● Model structure

- Observation: overall container for all associated datasets (top level of the model)
- Plane: to store each dataset associated to an Observation
- Artifact: the actual data files containing the observational data (e.g. FITS)
- Part: each describable part within an Artifact that has a complex data structure (e.g. FITS header data unit)
  - Description and discovery of the Part(s) rely on the Artifact's internal metadata content
- Chunk: further fine-grain level if also Part is a complex data structure (rare)
  - Usually not clearly separated in term of Artifact metadata

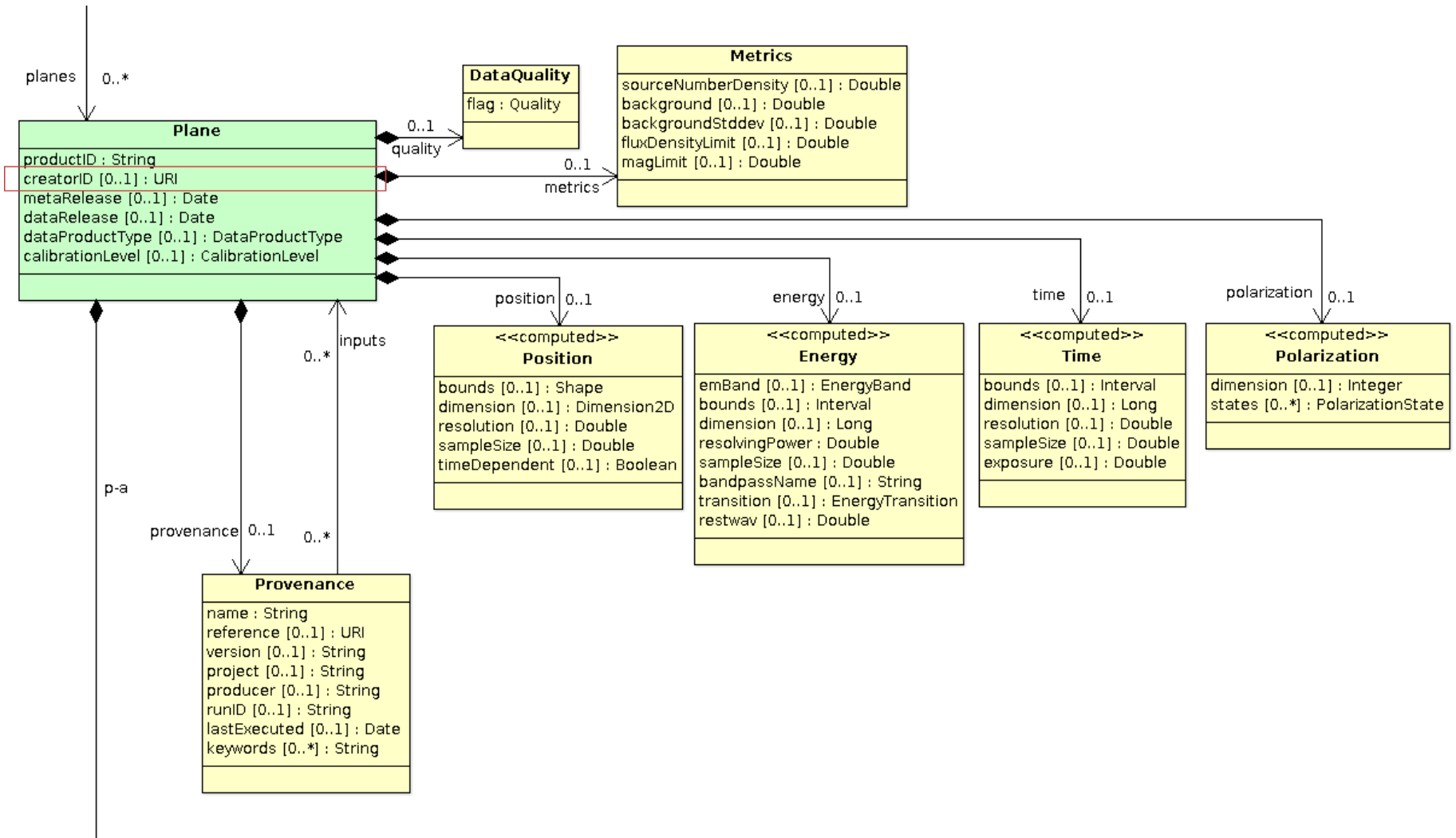
```
Observation
-> Plane
    -> Artifact
        -> Part
            -> Chunk
        -> Part
            -> Chunk
        -> Part
            ...
    -> Plane
        -> Artifact
            ...
-> Plane
    ...
```

- <http://www.opencadc.org/caom2/>

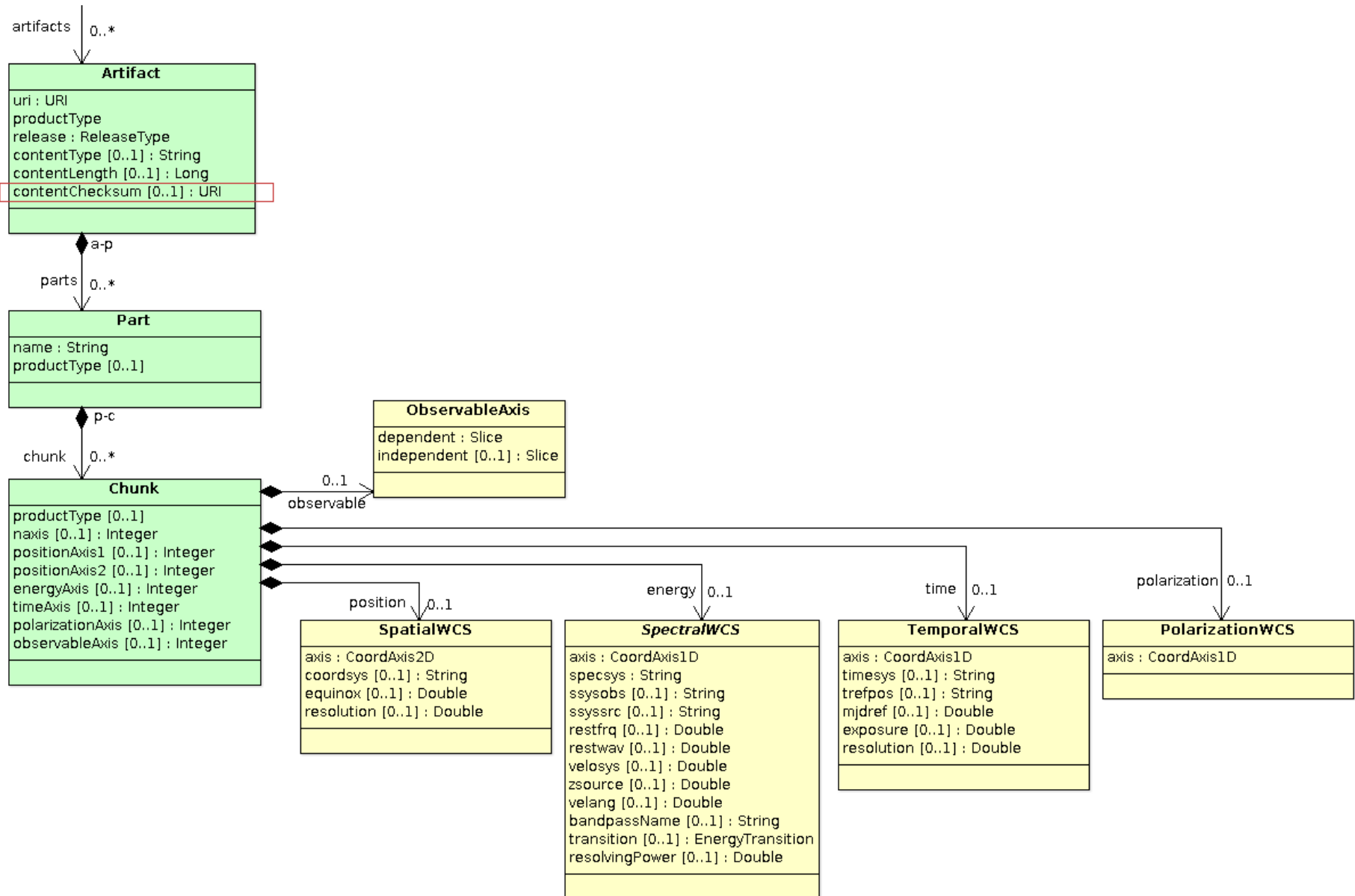
# CAOM – Observation



# CAOM – Plane



# CAOM – Artifact, Part, Chunk



# CAOM – Access to Instances

