
datS-doc Documentation

Release 0.1

Philippe Rocca-Serra and Alejandra Gonzalez-Beltran

Mar 27, 2017

Contents

1	Introduction:	1
1.1	First Steps with DATS	1
1.2	DATS Model	2
1.3	Dataset Distribution	5
1.4	Dataset Creator(s)	7
1.5	Dataset About	7
1.6	Dataset Provenance	9
1.7	Frequently Asked Questions	9
2	License:	11
3	Contributing:	13
4	Indices and tables:	15

CHAPTER 1

Introduction:

DATS, which stands for DAta Tag Suite, is a data description model designed and developed to describe datasets being ingested in [DataMed](#), a prototype for data discovery developed as part of the [NIH Big Data 2 Knowledge bioCADDIE project](#). For more information about the objectives of the bioCADDIE project, please have a look at the [bioCADDIE White Paper](#).

This documentation describes the DATS model and how to use it. More details about how DATS was designed and how it relates to other models can be found in the [documents accompanying each of the releases](#).

Table of Contents:

First Steps with DATS

This document gives an overview of the DATS components from a practical perspective, i.e. considering how DATS is used to describe a specific dataset considering a set of questions that determine the dataset provenance.

Who produced the dataset:

DATS records the *Person* (s) and *Organization* (s) associated with the dataset. In addition, it supports documenting their roles (e.g. creator, curator, developer, funder, principal investigator)

When was the data produced:

DATS records key *Date**(s) about the **Dataset*.

Each *Date* can specify its type, related to the event related to the key date (e.g. creation, update, validation, verification, deprecation of the dataset).

This mechanism of providing a generic *Date* indicating its type allows for extensions to new types of dates, which may be required in specific scenarios.

What is the dataset about:

DATS records the nature of information available in a dataset through the data type object.

Why was the data produced:

DATS supports to document the purpose, objective or hypothesis that gave origin to the dataset.

Where and How the dataset can be accessed:

DATS Model

Entity	Property	Definition
dataset	identifier	Primary identifier
	relatedIdentifiers	Related identifiers
	alternateIdentifiers	Alternate identifiers
	title	The name of the dataset
	types	A term, idea or concept
	creators	The person(s) who created the dataset
	dates	Relevant dates
	distributions	The distribution of the dataset
	dimensions	The dimensions of the dataset
	isCitedBy	The relevant dataset cited by the dataset
	producedBy	A study produced by the dataset
	isAbout	Different entities about which the dataset is about
	hasPart	A Dataset that is part of the dataset
	keywords	Tags associated with the dataset
	acknowledges	The grant(s) that funded the dataset
	extraProperties	Extra properties of the dataset
DatasetDistribution		“A specific distribution of the dataset”
	identifiers	Primary identifier
	alternateIdentifiers	Alternate identifiers
	relatedIdentifiers	Related identifiers
	title	“The name of the dataset”
	description	A textual name of the dataset
	dates	“Relevant dates”
	“storedIn “	The data repository
	version	A release point
	accessModalities	The information on how to access the dataset
	licenses	The terms of use
	curationStatus	The level of curation
	conformsTo	A data standard
	format	The technical format
	qualifiers	“One or more”
	“size “	The size of the dataset
	unit	“The unit of measurement”

Entity	Property	Definition
DataStandard	extraProperties	Extra properties
	identifiers	Primary identifiers
	alternateIdentifiers	Alternate identifiers
	relatedIdentifiers	Related identifiers
	name	“The name of the standard”
	type	“The nature of the standard”
	description	A textual description
	licenses	The terms of use
	version	A release point
	extraProperties	Extra properties
DataRepository	identifiers	Primary identifiers
	alternateIdentifiers	Alternate identifiers
	relatedIdentifiers	Related identifiers
	name	The name of the repository
	description	A textual description
	dates	Relevant dates
	scopes	“Information about the scope of the repository”
	types	“A description of the types of data in the repository”
	licenses	The terms of use
	version	“A release point”
	publishers	The person(s) who publish the data
	aggregatorOf	The DataRepository that aggregates this repository
	accessModalities	The information about how to access the data
	extraProperties	Extra properties
	extraProperties	Extra properties
Software	identifiers	Primary identifiers
	alternateIdentifiers	Alternate identifiers
	relatedIdentifiers	Related identifiers
	name	The name of the software
	licenses	The terms of use
	isUsedBy	The data accession numbers
	manufacturer	The person or organization that manufactured the software
	version	A release point
	extraProperties	Extra properties
Publication	identifiers	Primary identifiers
	alternateIdentifiers	Alternate identifiers
	relatedIdentifiers	Related identifiers
	title	“The name of the publication”
	“dates “	“Relevant dates”
	type	“Publication type”
	publicationVenue	The name of the publication venue
	authorsList	The list of authors
	authors	The person(s) who authored the publication
	acknowledges	The grant(s) that funded the publication
	licenses	The terms of use
	extraProperties	Extra properties

Entity	Property	Definition
IdentifiersInformation		Information
	identifier	A code unique
	identifierSource	The identifier
AlternateIdentifiersInformation		Information
	alternateIdentifier	An identifier
	alternateIdentifierSource	The identifier
RelatedIdentifiersInformation		Information
	relatedIdentifier	An identifier
	relatedIdentifierSource	The identifier
	relationType	The type of
Annotation		“A pair of va
	“value “	A label or va
	ontologyTermIRI /suggested renaming = ValueIRI	The IRI of a
Date		“Information
	date	A date follow
Access		Information
	identifiers	Primary iden
	alternateIdentifiers	Alternate id
	relatedIdentifiers	Related iden
	landingPage	A web page
	accessURL	“A URL from
	types	“Method to
	authorizations	Types of ver
	authentications	“Types of ve
	licenses	Terms of usa
	extraProperties	Extra proper
Grant		An allocated
	identifiers	Primary iden
	alternateIdentifiers	Alternate id
	relatedIdentifiers	Related iden
	name	The name of
	funds	The study or
	funders	The person(s
	awardees	The person(s
	extraProperties	Extra proper
License		“A legal doc
	identifiers	Primary iden
	alternateIdentifiers	Alternate id
	relatedIdentifiers	Related iden
	name	The name of
	version	The version
	creators	The person(s
	extraProperties	Extra proper
Dimension		“A feature o
	identifiers	Primary iden
	alternateIdentifiers	Alternate id
	relatedIdentifiers	Related iden
	name	“The name o
	types	“A term, ide

Entity	Property	Definition
	partOf	The dataset(
	description	A textual na
	values	The actual c
	unit	“A reference
	“isAbout “	“A material
	extraProperties	Extra proper
	information	The measure
	method	The procedu
	platform	“The set of i
	instrument	The specific
Material	extraProperties	Extra proper
		“A physical
	identifiers	Primary iden
	alternateIdentifiers	Alternate id
	relatedIdentifiers	Related iden
	name	The name of
	derivesFrom	A material f
	bearerOfDisease	The patholo
	taxonomicInformation	The taxonom
	involvedInBiologicalEntity	A biological
Person	characteristics	The characte
	roles	The roles pl
	extraProperties	Extra proper
		A human be
	identifiers	Primary iden
	alternateIdentifiers	Alternate id
	relatedIdentifiers	Related iden
	fullName	“The first na
	firstName	The given na
	middleInitial	The first lett
	lastName	The person’s
	email	An electroni
	affiliations	The organiz
	roles	“The roles a
	extraProperties	Extra proper
	identifiers	Primary iden
	alternateIdentifiers	Alternate id
	relatedIdentifiers	Related iden
	name	The name of
	abbreviation	“The shortna
	postalAddress	“The postal,
	roles	“The roles o
	extraProperties	Extra proper

Dataset Distribution

Where and How (can the dataset be accessed):

- Document DataSet Distribution options. This encompasses specifying:

- data availability (boolean choice: available, unavailable)
- data formats or mime-types ([terminology needs to be specified] ‘resource: https://github.com/lukaszsliva/friendly_mime/blob/master/mimes.csv>‘_)
- data access conditions
- data compression (boolean choice: compressed, uncompressed)
- data encryption (boolean choice: encrypted, non-encrypted)
- data privacy protection (fully identifiable, pseudo-anonymized, full anonymized...[terminology needs to be specified])

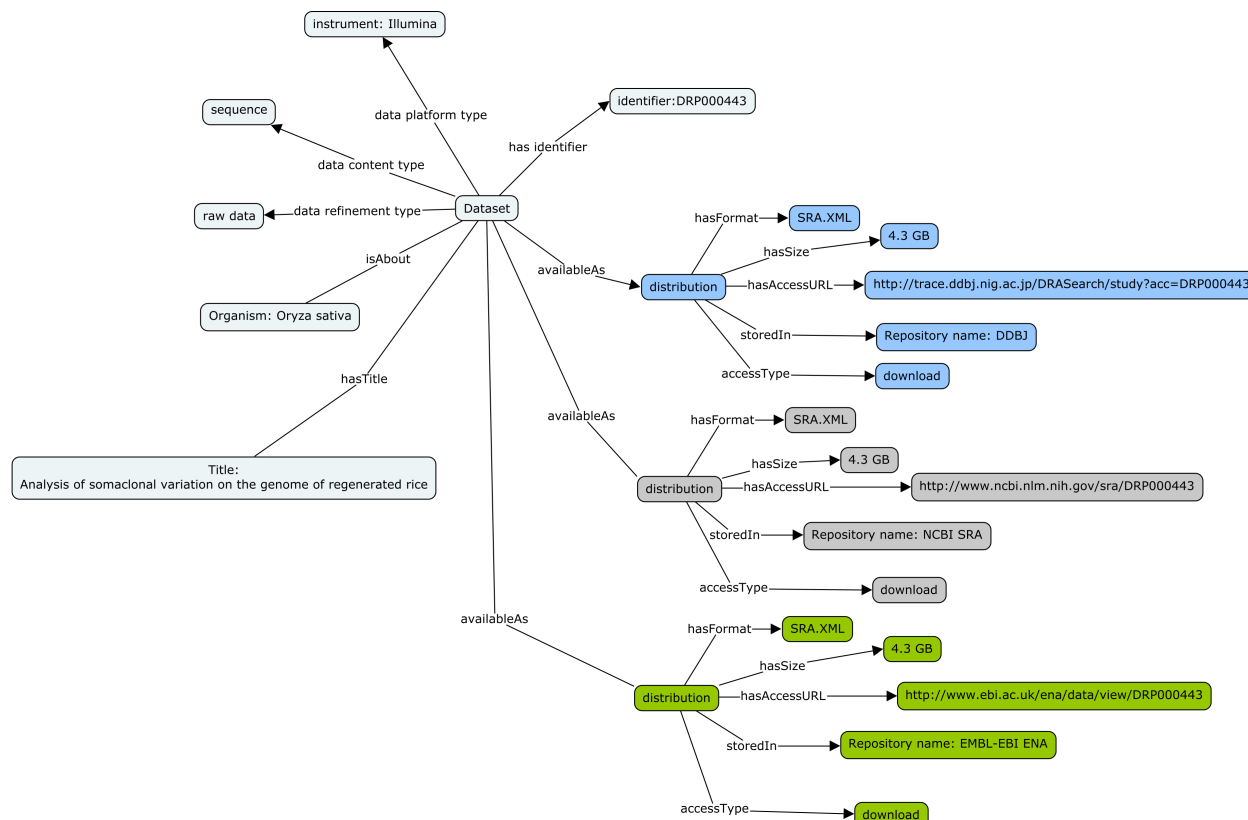
The image below provides an graphical overview of how to use Biocaddie DATS objects to encode information about dataset availability in a similar file format but from 3 distinct data repositories, each with it own access modalities.

The three INSDC sequence databases (DDBJ, SRA and ENA) exchange their data and provide the same datasets it in the three sites. Let’s consider an example dataset.

The same Dataset identified by accession number DRP000443 can be accessed through the following 3 access URI pages:

- [DDBJ](#):
- [SRA](#):
- [ENA](#):

While the distributions use the same Format, the accessURL are different as are the Repository but these distributions are all about the same dataset



The block below shows a snippet of a bioCADDIE DATS JSON document holding key information about dataset distribution. Note the link to *access information* and *data file format* information.

Dataset Creator(s)

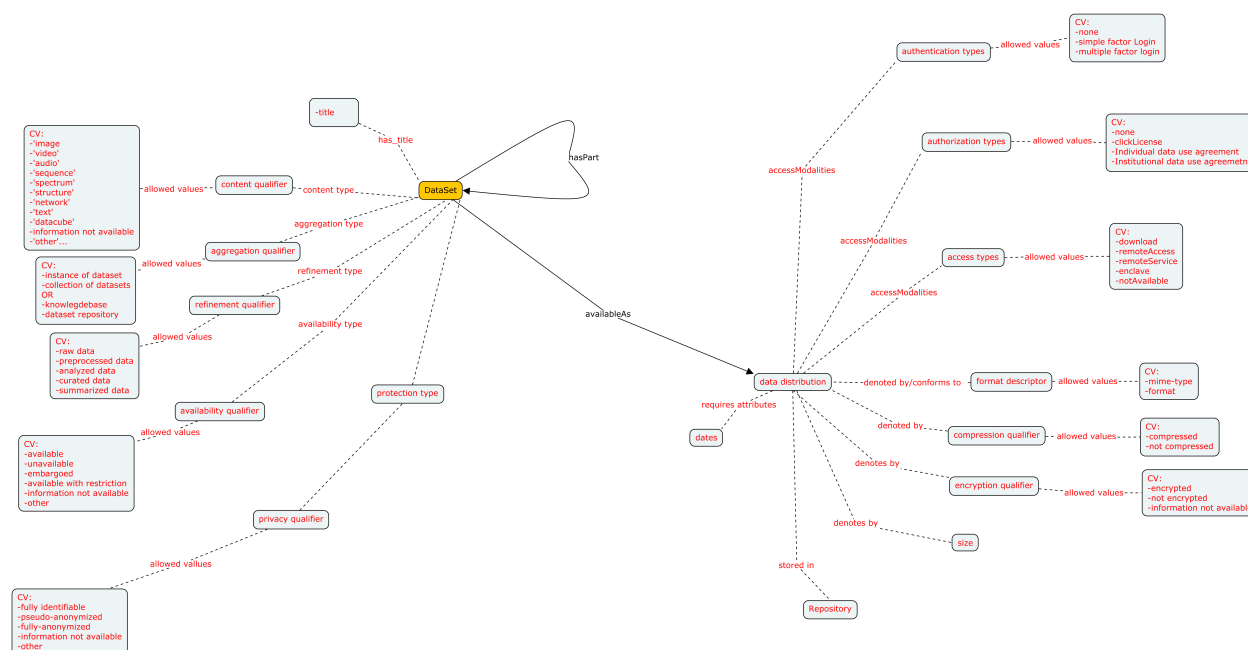
Who (produced the dataset):

- Document the Person(s) or Organization(s) which contributed to the creation of the Dataset.
- Document their roles (creator,curator,developer,funder,principal investigator. . . [terminology needs to be specified])

Dataset About

Describing what the dataset is about (i.e what was the scope, objective, materials) and providing information about the type of data associated with the given dataset:

- Document the nature of information available in a dataset through the Biocaddie **‘data type’** object.



In this context, the **‘data type’** required to annotate a DataSet should be viewed as a *content type* [terminology needs to be specified]). This encompasses the nature of the signal recorded in a dataset or information content of interest. For instance: gene expression data or phenotypic data, electronic health records But mime-type may be used. * chemical * sequence * spectrum * audio * image * video * ...

but other descriptors may be used such as Biosharing, Scicrunch or re3data category/data domain descriptors.

- Data aggregation type:

In the context of DataMed indexing, the information obtained from repositories may correspond to datasets served individually or may correspond to collections or records. As these 2 situations represent a very different metadata context, the Biocaddie DATS model allows to distinguish between the two cases.

- collection (as in ‘collection of instances’)
- singleton (as in ‘individual instance’)
- Data refinement type:

To describe the level of data processing associated with the data available from the dataset and its distributions....[terminology needs to be specified])

- raw data
- preprocessed data
- analyzed data
- summarized data
- curated data
- reannotated data
- ...
- data privacy protection type: (applicable only to human/clinical data)
 - fully identifiable none
 - pseudo-anonymized data
 - fully anonymized data
 - not information available
 - ...
- Document the Material, object, scope and Biological Entities the dataset is about and their characteristics or properties.
- Document the nature of intervention and Treatment applied to the Material, if any or if applicable.
- Data Types and specific Platform

Currently, in DataMed, datasets can be search according to Data Type (.e.g Proteomics data) and/or by Platform (e.g. Illumina) DATS provides a mechanism via DataType object to qualify the nature of the data collected in a Dataset. The 4 facets/attributes allow to incrementally specify the type of information contained by the data and how it has been produced

- **data acquisition / method type:** This attribute allows to indicate the technique or technology , also known sometimes as data modality used to acquire the signal. For instance:
 - ‘crystallography’,
 - ‘mass spectrometry’
 - ‘nucleic acid sequencing’,
 - ‘computational simulation’
 - ‘questionnaire based survey’
 - ‘nuclear magnetic resonance spectroscopy’
 - ‘nuclear magnetic resonance imaging’
 - ‘questionnaire’
 - ...
- **platform/instrument type**
 - Agilent, Bruker, Affymetrix, Illumina, SeaHorse
 - HumanHap550v3.0
 - HumanExome-12 v1.1 BeadChip

- Sentrix Human-6 Expression BeadChip
- SureSelect Human All Exon v2 - 44Mb
- HiSeq 2000
- ...

Dataset Provenance

In order to proceed with indexing a data source under bioCADDIE DataMed, it is essential to provide information about the actual source of information. This means unambiguously identifying the repository, the actual material from that resource used as input to the transformation allowing processing by DataMed software agents.

This falls under the provenance information section of the DATS for DataMed.

- identify the repository
- document the **url** or **filename** and address of the source information
- document the **date of last access** to the resource as input to the data transformation
- document the data transformation pipeline in the datamed infrastructure, ideally by pointed to the biocaddie [github repository](#) .

Frequently Asked Questions

Why are some properties (e.g. “title” and “description”) included in both Dataset and DataDistribution?

When designing DATS we chose to be flexible and consider some redundancy by including properties in both Dataset as well as DatasetDistribution, even though in some cases it might be expected that a Dataset property should be inherited by their DatasetDistributions. We followed this approach to cover cases where repositories may have different information. For example, it would be possible that each DatasetDistribution has more information in its “description” on how the distribution was produced, adding more details to the general information in the corresponding Dataset.

CHAPTER 2

License:

BioCADDIE DATS is licensed under [Creative Commons Attribution Share-Alike 4.0](#).

CHAPTER 3

Contributing:

If you wish to contribute to DATS and/or this documentation, please report issues in our [tracker](#) or contact us directly ([agbeltran](#) and [proccaserra](#)).

The different releases of DATS are available in the [bioCADDIE Working Group 3 Github Repository](#), including documents and appendixes, JSON schemas, JSON-LD context files and JSON-LD instance files.

CHAPTER 4

Indices and tables:

- `genindex`
- `modindex`
- `search`