

# **Metadata Specification to support Innovation Gateway Minimum Viable Product (MVP)**

**November 2019**

## Purpose

The following document has been developed to define the metadata needed for the Innovation Gateway MVP stage of the onboarding process.

The specification is initially focused on Summary and Business level metadata (see proposed metadata levels in principles below)

The metadata specification has been created and generated from:

<https://hdronboard.metadata.works/#/122894/dataClass/123142/main?path=122894-all-123142>

Document Control Information							
Version	Status	Author(s)	Date Created	Reviewer(s)	Date Reviewed	Approved	Notes
1.1.7	Trial	A.Milward, A.Tripathi, M.Jones,	14/11/2019	A.Milward, J.Davies, J.Welch, S.Varma, D.Seymor, N.Sebire, G.Reilly	15/11/2019	15/11/2019	
<b>Minimum Retention Period</b>	Permanent						
<b>Disposal Required</b>	Archive when new version created						
<b>Classification</b>	PUBLIC						
<b>Next Review due:</b>	22/11/2019						

# Contents

<b><u>METADATA SPECIFICATION TO SUPPORT INNOVATION GATEWAY MINIMUM VIABLE PRODUCT (MVP) .....</u></b>	<b><u>1</u></b>
<b><u>NOVEMBER 2019 .....</u></b>	<b><u>1</u></b>
<b><u>PURPOSE .....</u></b>	<b><u>2</u></b>
<b><u>CONTENTS .....</u></b>	<b><u>3</u></b>
<b><u>INTRODUCTION .....</u></b>	<b><u>5</u></b>
<b><u>SUMMARY METADATA .....</u></b>	<b><u>6</u></b>
IDENTIFIER .....	6
TITLE.....	6
ABSTRACT .....	6
PUBLISHER .....	7
CONTACT POINT .....	7
ACCESS RIGHTS.....	7
GROUP.....	8
<b><u>BUSINESS METADATA .....</u></b>	<b><u>9</u></b>
<b><u>REQUIRED .....</u></b>	<b><u>9</u></b>
DESCRIPTION .....	9
RELEASE DATE .....	9
ACCESS REQUEST COST .....	10
ACCESS REQUEST DURATION.....	10
DATA CONTROLLER .....	10
DATA PROCESSOR .....	10
LICENSE .....	11
DERIVED DATASETS .....	11
LINKED DATASET .....	12
<b><u>RECOMMENDED .....</u></b>	<b><u>13</u></b>
COVERAGE AND DETAIL .....	13
GEOGRAPHIC COVERAGE .....	13
PERIODICITY .....	13
DATASET END DATE .....	14
DATASET START DATE.....	14
JURISDICTION .....	15
POPULATION TYPE.....	15
STATISTICAL POPULATION.....	15
AGE BAND .....	15
PHYSICAL SAMPLE AVAILABILITY .....	16
KEYWORDS .....	16
FORMAT AND STRUCTURE .....	16
CONFORMS TO .....	16
CONTROLLED VOCABULARY.....	16

LANGUAGE.....	17
FORMAT .....	18
FILE SIZE .....	18
ATTRIBUTION .....	18
CREATOR .....	18
CITATIONS .....	19
DOI .....	19
<b>TECHNICAL METADATA .....</b>	<b>20</b>
TABLE NAME .....	20
TABLE DESCRIPTION.....	20
COLUMN NAME.....	20
COLUMN DESCRIPTION .....	20
DATA TYPE.....	20
SENSITIVE .....	20

## Introduction

As part of the development of the Innovation Gateway Minimum Viable Product (MVP), Health Data Research UK (HDR UK) must complete an initial onboarding activity to collect high level information that describes the datasets that are available for research and innovation across members of the UK Health Data Research Alliance and the seven health data research hubs. This information must be human and machine-readable allowing researchers to understand the data that is available and allowing the Gateway to index the information to return relevant search results.

HDR UK's initial requirements have been used as a baseline specification and iterated based on feedback and suggestions from IBM (supporting MVP portal), Oxford University (providing the Metadata catalogue) and HDR UK. This is part of an MVP. We have tried to provide coverage for organisation who are already collecting metadata and would like to provide it to the gateway, but we expect as the project progresses there will be further iterations of the specification.

The following principles have governed the work and decisions made:

- Specification that will be distributed in first iteration will be "good enough MVP" rather than perfect
- Metadata should be based on a user need and a user story
- Metadata should use standard terminologies whenever possible and 2.1. DCAT is the de-Facto standard and links to schema.org should be provided where available
  - If an identified need does not have an equivalent DCAT entry, another standard should be selected, including ONS standards
  - Only if no standard can be found should we create a new metadata entry
- Naming convention: Metadata "titles" will be lowercase, underscore separated (note that in the human readable specification we are using Capitalised for ease of reading)
- Prioritisation of Metadata has been based on Impact vs Effort (effort needs to consider potential for automation in the longer-term vs manual in nearer term) however, this will be easier to quantify as the project progresses

The following different levels of metadata quality have been defined:

- Summary Metadata (Mandatory)
- Business Metadata
  - Required
  - Recommended
- Technical Metadata (variable level metadata i.e. data elements, types etc.)
  - Required
  - Recommended
- Enhanced Metadata (not in scope for initial on boarding)
  - Profiling
  - Quality
  - Other

As the project progresses metadata will be categorised accordingly based on use need.

## Summary metadata

Summary metadata must be completed for [data custodians](#) onboarding metadata into the Innovation Gateway MVP.

<b>Identifier</b>	<p><b>Completion Guidance:</b> Please provide a local identifier used to identify the dataset (if available).  <i>Note: if a DOI is available this can be provided in the Attribution section.</i></p> <p><b>Definition:</b> A unique identifier of the item.  <i>Note: The identifier might be used as part of the URI of the item, but still having it represented explicitly is useful.</i></p> <p>Vocabulary: <a href="#">Data Catalog Vocabulary (DCAT)</a></p> <p>RDF Property: <a href="#">dc:identifier</a></p> <p>Range: <a href="#">rdfs:Literal</a></p> <p>Source: <a href="#">dct:identifier</a><a href="#">dct:resource_identifier</a></p>	<p>xs:string</p> <p>Min Occurs: 0 Max Occurs: 1</p>
	As A [user]	Researcher
	I Want To	know the name of identifier for the dataset
	So That	I can identify, refer to it

<b>Title</b>	<p><b>Completion Guidance:</b> This is the name of the dataset.</p> <p><b>Definition:</b> A name given to the item.</p> <p>Vocabulary: <a href="#">Data Catalog Vocabulary (DCAT)</a></p> <p>RDF Property: <a href="#">dc:title</a></p> <p>Range: <a href="#">rdfs:Literal</a></p> <p>Source: <a href="#">dct:title</a></p> <p>See also: schema.org <a href="#">Dataset/name</a></p>	<p>xs:string</p> <p>Min Occurs: 1 Max Occurs: 1</p>
	As A [user]	Researcher
	I Want To	know the name of the dataset
	So That	I can identify and find it

<b>Abstract</b>	<p><b>Completion Guidance:</b> Provide a short summary of the dataset (limited to 256 characters)</p> <p><b>Definition:</b> A summary of the resource.</p> <p>Vocabulary: <a href="#">DCMI Metadata Terms</a></p> <p>RDF Property: <a href="#">rdf-syntax-ns#Property</a></p> <p>Refines: <a href="#">dc:description</a></p> <p>Source: <a href="#">dc:abstract</a></p> <p>See also: schema.org <a href="#">Dataset/abstract</a></p>	<p>xs:string</p> <p>// min 5 // max 256 x =~ /^.{5,256}\$/</p> <p>Min Occurs: 1 Max Occurs: 1</p>
	As A [user]	Researcher
	I Want To	read a plain English summary of the dataset
	So That	I can decide if I want to find out more

<p><b>Completion Notes:</b> This is the organisation that is the data custodian. They are responsible for the data access request process, as well as publishing and maintaining the metadata</p> <p><b>Definition:</b> The entity responsible for making the item available.</p> <p>Vocabulary: <a href="#">Data Catalog Vocabulary (DCAT)</a></p> <p>RDF Property: <a href="#">dc:publisher</a></p> <p>Source: <a href="#">dct:publisher</a></p> <p>See also: schema.org <a href="#">Dataset/publisher</a></p>		xs:string	Min Occurs: 1 Max Occurs: 1
<b><u>Publisher</u></b>			
As A [user]	Researcher		
I Want To	know the publisher		
So That	I can understand the context of the dataset		

<p><b>Completion Guidance:</b> A URL with the information for and/or the email address of, the person or role within the data custodian organisation, who is responsible for the data access process and metadata maintenance</p> <p><b>Definition:</b> Relevant contact information for the catalogued resource. Use of vCard is recommended [VCARD-RDF].</p> <p>Vocabulary: <a href="#">Data Catalog Vocabulary (DCAT)</a></p> <p>RDF Property: <a href="#">dct:contactPoint</a></p> <p>Range: <a href="#">vcard:Kind</a></p> <p>Source: <a href="#">dct:contact point</a></p>		xs:string	Min Occurs: 1 Max Occurs: 1
<b><u>Contact Point</u></b>			
As A [user]	Researcher		
I Want To	contact person		
So That	I can understand the datasets and request access		

<p><b>Completion Guidance:</b> Text or link to website/documentation where data access request process and/or guidance is provided</p> <p><b>Definition:</b> Information about who can access the resource or an indication of its security status.</p> <p>Vocabulary: <a href="#">Data Catalog Vocabulary (DCAT)</a></p> <p>RDF Property: <a href="#">dct:accessRights</a></p> <p>Range: <a href="#">dct:RightsStatement</a></p> <p>Usage note: Access Rights MAY include information regarding access or restrictions based on privacy, security, or other policies.</p> <p>Source: <a href="#">access rights</a></p> <p>See also: schema.org <a href="#">Dataset/conditionsOfAccess</a></p>		xs:string	Min Occurs: 1 Max Occurs: *
<b><u>Access Rights</u></b>			
As A [user]	Researcher		
I Want To	know what I can do with the data		
So That	I have a legal and ethical basis for analysing it in my research		

**Group**

**Completion Guidance:** Please complete if the dataset is part of a group family or collection i.e. Hospital Episode Statistics has several constituents

**Definition:** A collection is an entity that provides a structure to some constituents that must themselves be entities. These constituents are said to be member of the collections.

More specifically, Influence is the capacity of an entity, activity, or agent to have an effect on the character, development, or behaviour of another by means of usage, start, end, generation, invalidation, communication, derivation, attribution, association, or delegation.

Vocabulary: [PROV-O: The PROV Ontology](#)

RDF Property: [prov:Collection](#)

Sub-property of: [prov:qualifiedInfluence](#)

Domain: [prov:Entity](#)

Range: [prov:Entity](#)

Source: schema.org [Dataset/isPartOf](#)

Min Occurs: 0  
Max Occurs: 1

As A [user]	Researcher
I Want To	know whether the dataset is part of a wider group
So That	I can find other datasets of the group



## Business Metadata

Business metadata provides context for datasets that improve their discoverability through the Innovation Gateway. HDR UK has two categories of business metadata they ask [data custodians](#) to provide.

- **Required:** this information is required as part of the onboarding process and is key to making your data more discoverable for researchers
- **Recommended:** this information is recommended for the onboarding process as it gives researchers key contextual information and will improve the dataset rank

### Required

This information is required as part of the onboarding process and is key to making your data more discoverable for researchers

<p><b>Completion Guidance:</b> A free-text account of the data (limited to 50000 characters) and/or a resolvable URL of a document (webpage, pdf, word) that describing the dataset. This provides detailed context about the dataset.</p>		
<b>Description</b>	<b>Definition:</b> A free-text account of the record.	xs:string
	Vocabulary: <a href="#">Data Catalog Vocabulary (DCAT)</a>	// min 5
	RDF Property: <a href="#">dc:description</a>	// max 50000
	Range: <a href="#">rdfs:Literal</a>	x ==~/^.{5,50000}\$/
	Source: <a href="#">dcat:record_description</a>	Min Occurs: 0
	See schema.org <a href="#">Dataset:description</a>	Max Occurs: 1
As A [user]	Researcher	
I Want To	know the description	
So That	I can understand information about the dataset	

<p><b>Completion Guidance:</b> Date of the latest release of the dataset. If this is a regular release i.e. quarterly, please complete this alongside <a href="#">periodicity</a>. Periodicity will be used to determine when the next update is expected. <b>Definition:</b> Date of formal issuance (e.g., publication) of the distribution.</p>		
<b>Release Date</b>	<b>Definition:</b> Date of formal issuance (e.g., publication) of the distribution, encoded using the relevant ISO 8601 Date and Time compliant string [DATETIME] and typed using the appropriate XML Schema datatype [XMLSCHEMA11-2] (xsd:gYear, xsd:gYearMonth, xsd:date, or xsd:dateTime).	xs:date
	<i>Usage note: This property SHOULD be set using the first known date of issuance.</i>	Min Occurs: 0
	Vocabulary: <a href="#">Data Catalog Vocabulary (DCAT)</a>	Max Occurs: 1
	RDF Property: <a href="#">dc:issued</a>	
	Range: <a href="#">rdfs:Literal</a>	
	Source: <a href="#">dcat:distribution_release_date</a>	
	See schema.org <a href="#">Dataset/datePublished</a>	
As A [user]	Researcher	
I Want To	publication date	

So That	I know how recent and relevant the dataset is
<b><u>Access Request Cost</u></b>	<p><b>Completion Guidance:</b> Please provide a free text description of the cost, or an indication of the range/calculation of costs, for processing a data access request.</p> <p><b>Definition:</b> Indication of cost (in GBP) for processing each data access request by the <a href="#">data custodian</a>.</p> <p>Source: Schema.org <a href="#">Offer:price</a></p>
As A [user]	Researcher
I Want To	know how much it will cost to process my data access request
So That	So that I can budget and decide accordingly

<b><u>Access Request Duration</u></b>	<p><b>Completion Guidance:</b> Please provide an indicative value and/or estimate of the typical processing time.</p> <p><b>Definition:</b> Indication of the typical duration of a data access request.</p> <p>Source: Schema.org <a href="#">Offer:deliveryLeadTime</a></p>	<p>xs:string</p> <p>Min Occurs: 0 Max Occurs: 1</p>
As A [user]	Researcher	
I Want To	know how much time it will take to process my data access request	
So That	I can plan and decide accordingly	

<b><u>Data Controller</u></b>	<p>Data Controller means a person who (either alone or jointly or in common with other persons) determines the purposes for which and the way any Data Subject data, specifically personal data or are to be processed.</p> <p><b>Please complete if the data custodian is not the Data Controller.</b></p> <p>Source: <a href="#">ICO Key Terms: Controllers and Processors, Article 4.7 of GDPR</a></p> <p>Vocabulary: <a href="#">Data Privacy Vocabulary</a></p> <p>Source: <a href="#">dpv:DataController</a></p>	<p>xs:string</p> <p>Min Occurs: 0 Max Occurs: 1</p>
As A [user]	Processor	
I Want To	know who has ultimately responsibility for the dataset (if not the publisher)	
So That	I can continue to process their data	

<b><u>Data Processor</u></b>	<p>A Data Processor, in relation to any Data Subject data, specifically personal data, means any person (other than an employee of the data controller) who processes the data on behalf of the data controller.</p> <p><b>Please complete if the data custodian is the Data Processor.</b></p> <p>Source: <a href="#">ICO Key Terms: Controllers and Processors, Article 4.7 of GDPR</a></p> <p>Vocabulary: <a href="#">Data Privacy Vocabulary</a></p>	<p>xs:string</p> <p>Min Occurs: 0 Max Occurs: *</p>
------------------------------	--	---

	Source: <a href="#">dpv:DataProcessor</a>
As A [user]	Administrator
I Want To	know who the data processor is
So That	I that can understand the legal basis for processing

<b><u>License</u></b>	<p>Definition: A legal document under which the distribution is made available.</p> <p>Usage note: Information about licenses and rights SHOULD be provided on the level of Distribution. Information about licenses and rights MAY be provided for a Dataset in addition to but not instead of the information provided for the Distributions of that Dataset. Providing license or rights information for a Dataset that is different from information provided for a Distribution of that Dataset SHOULD be avoided as this can create legal conflicts.</p> <p>Vocabulary: <a href="#">Data Catalog Vocabulary (DCAT)</a></p> <p>RDF Property: <a href="#">dct:license</a></p> <p>Range: <a href="#">dct:LicenseDocument</a></p> <p>Source: <a href="#">dcat:distribution_license</a></p> <p>Source: <a href="#">Dataset/license</a></p>	xs:string	Min Occurs: 0 Max Occurs: 1
	As A [user]	Researcher	
	I Want To	know the license that applies to the data	
	So That	I know the legal information that describes how the distribution is made available	

<b><u>Derived Datasets</u></b>	<p><b>Completion Guidance:</b> indicate if derived datasets are available and type of derivation available</p> <p><b>Definition:</b> A derivation is a transformation of an entity into another, an update of an entity resulting in a new one, or the construction of a new entity based on a pre-existing entity. If this is Entity :e, then it can qualify how it was derived using prov:qualifiedDerivation [a prov:Derivation; prov:entity :e; :foo :bar ].</p> <p>Vocabulary: <a href="#">PROV-O: The PROV Ontology</a></p> <p>RDF Property: <a href="#">prov:qualifiedDerivation</a></p> <p>Sub-property of: <a href="#">prov:qualifiedInfluence</a></p> <p>Domain: <a href="#">prov:Entity</a></p> <p>Range: <a href="#">prov:Derivation</a></p> <p>PROV-DM term: <a href="#">prov:wasDerivedFrom</a></p> <p>See also: Schema.org <a href="#">Dataset/isBasedOn</a></p>	xs:string	Min Occurs: 0 Max Occurs: *
	As A [user]	Researcher	
	I Want To	Know if there are any derivations	
	So That	I can reuse pre-processed information	

**Linked Dataset**

**Completion Guidance:** if applicable, describe other datasets that have previously been linked to this dataset and their availability. If no datasets have currently been linked, please indicate if there is the possibility to link the datasets as part of the data access request process.

**Definition:** Link to a description of a relationship with another resource.

*Usage note: Used to link to another resource where the nature of the relationship is known but does not match one of the standard [DCTERMS] properties (dct:hasPart, dct:isPartOf, dct:conformsTo, dct:isFormatOf, dct:hasFormat, dct:isVersionOf, dct:hasVersion, dct:replaces, dct:isReplacedBy, dct:references, dct:isReferencedBy, dct:requires, dct:isRequiredBy) or [PROV-O] properties (prov:wasDerivedFrom, prov:wasInfluencedBy, prov:wasQuotedFrom, prov:wasRevisionOf, prov:hadPrimarySource, prov:alternateOf, prov:specializationOf).*

xs:string

Min Occurs: 0  
Max Occurs: \*

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

RDF Property: [dcat:qualifiedRelation](#)

Sub-property of: [prov:qualifiedInfluence](#)

Domain: [dcat:Resource](#)

Range: [dcat:Class:Relationship](#)

Source: [dcat:resource\\_qualified\\_relation](#)

As A [user]	Researcher
I Want To	know if there are any datasets that have already been linked
So That	I can use more pre-processed data in my research

## Recommended

This information is recommended for the onboarding process as it gives researchers key contextual information and will improve the dataset rank:

### Coverage and detail

<b>Geographic coverage</b>	<p><b>Definition:</b> The geographical area covered by the dataset.</p> <p><i>Usage note:</i> The spatial coverage of a dataset may be encoded as an instance of <a href="#">dct:Location</a> or may be indicated using a URI reference (link) to a resource describing a location. It is recommended that links are to entries in a well-maintained gazetteer such as Geonames.</p>	xs:string	Min Occurs: 0 Max Occurs: 1
	Vocabulary: <a href="#">Data Catalog Vocabulary (DCAT)</a>		
	RDF Property: <a href="#">dct:spatial</a>		
	Range: <a href="#">dct:Location</a> (A spatial region or named place)		
	Source: <a href="#">dcat:dataset_spatial</a>		
See also: schema.org <a href="#">Dataset/spatialCoverage</a>			
As A [user]	Researcher		
I Want To	know the geographic are covered by the dataset		
So That	I can contextualise my research		

<b>Periodicity</b>	<p><b>Completion guidance:</b> please indicate the frequency of publishing and/or the rate that the whole dataset is updated.</p> <p><b>Definition:</b> The frequency at which dataset is published.</p> <p><i>Usage note:</i> The value of <a href="#">dct:accrualPeriodicity</a> gives the rate at which the dataset-as-a-whole is updated. This may be complemented by <a href="#">dct:temporalResolution</a> to give the time between collected data points in a time series.</p>	xs:string	Min Occurs: 0 Max Occurs: 1
	Vocabulary: <a href="#">Data Catalog Vocabulary (DCAT)</a>		
	RDF Property: <a href="#">dct:accrualPeriodicity</a>		
	Range: <a href="#">dct:Frequency</a> (A rate at which something recurs)		
	Source: <a href="#">dcat:dataset_frequency</a>		
See also: Schema.org <a href="#">Dataset/temporalCoverage</a>			
As A [user]	Researcher		
I Want To	know the frequency of the dataset		
So That	I know when the next release is		

<b><u>Dataset end date</u></b>	<b>Completion Guidance:</b> The end of the time period that the dataset provides coverage for.		
	<b>Definition:</b> The end of the period.		
	Vocabulary: <a href="#">Data Catalog Vocabulary (DCAT)</a>		
	RDF Property: <a href="#">dcat:endDate</a>		
	Domain: <a href="#">dct:PeriodOfTime</a>	xs:date	Min Occurs: 0
	Source: <a href="#">dcat:period_end_date</a>		Max Occurs: 1
	Range: <a href="#">rdfs:Literal</a> encoded using the relevant ISO 8601 Date and Time compliant string [DATETIME] and typed using the appropriate XML Schema datatype [XMLSCHEMA11-2] ( <a href="#">xsd:gYear</a> , <a href="#">xsd:gYearMonth</a> , <a href="#">xsd:date</a> , or <a href="#">xsd:dateTime</a> ).		
	Source : Schema.org <a href="#">Dataset/temporalCoverage</a>		
As A [user]	researcher		
I Want To	Know the time period for which data is available		
So That	I can use data that is within the time period I am interested in /know there is enough history.		

<b><u>Dataset start date</u></b>	<b>Completion Guidance:</b> The start of the time period that the dataset provides coverage for.		
	<b>Definition:</b> The start of the period.		
	Vocabulary: <a href="#">Data Catalog Vocabulary (DCAT)</a>		
	RDF Property: <a href="#">dcat:startDate</a>		
	Domain: <a href="#">dct:PeriodOfTime</a>	xs:date	Min Occurs: 0
	Source: <a href="#">dcat:period_start_date</a>		Max Occurs: 1
	Range: <a href="#">rdfs:Literal</a> encoded using the relevant ISO 8601 Date and Time compliant string [DATETIME] and typed using the appropriate XML Schema datatype [XMLSCHEMA11-2]		
	Source: <a href="#">dcat:period_start_date</a>		
	See also: Schema.org <a href="#">Dataset/temporalCoverage</a>		
As A [user]	Researcher		
I Want To	know about potential constraints on the data due to jurisdiction		
So That	I can tailor my request and research accordingly (or realise the data isn't suitable)		

	<p><b>Completion Guidance:</b> Please use ISO country code for country under whose laws the data subjects' data is collected, processed and stored</p> <p><b>Definition:</b> A named and identified geospatial area with defined borders which is used for exercising the action of the Rule. An IRI MUST be used to represent this value.</p> <p><b>Jurisdiction</b> Note: A code value for the area and source of the code must be presented in the Right Operand. Example: the [iso3166] Country Codes or the Getty Thesaurus of Geographic Names. Narrower terms: spatialCoordinates.</p> <p>Vocabulary: <a href="#">ODRL Vocabulary &amp; Expression</a></p> <p>Source: <a href="#">odrl:spatial</a></p> <p>See also: Schema.org <a href="#">Dataset/locationCreated</a></p>	xs:string	Min Occurs: 0 Max Occurs: 1
As A [user]	Researcher		
I Want To	know about potential constraints on the data due to jurisdiction		
So That	I can tailor my request and research accordingly (or realise the data isn't suitable)		

	<p><b>Completion Guidance:</b> Please provide a description(s) of the primary population type within the dataset i.e. participants in a study, or images with certain characteristics.</p> <p><i>Usage Note: Used with Statistical Population, which is the number of instances of a certain population type that satisfy some set of constraints.</i></p> <p>Source: Schema.org <a href="#">Population Type</a></p>	xs:string	Min Occurs: 0 Max Occurs: 1
As A [user]	Researcher		
I Want To	Know the population type within the data		
So That	I can understand whether data is suitable for my need		

	<p><b>Completion Guidance:</b> Please provide the primary population size within the dataset i.e. x number of participants in a study, or y number of images with certain characteristics.</p> <p><i>Usage Note: Used with Population Type, which specifies the type of the population in the dataset. Statistical Population is the number of instances of a certain population type that satisfy some set of constraints.</i></p> <p>Source: Schema.org <a href="#">Statistical Population</a></p>	xs: decimal	Min Occurs: 0 Max Occurs: 1
As A [user]	Researcher		
I Want To	Know the population size and type within data		
So That	I can understand whether data is suitable for my need		

	<p>Harmonised Principle Age Band Groups 1 &amp; 2</p> <p>Source: <a href="#">gss.civilservice.gov.uk</a></p>	xs:string	Min Occurs: 0 Max Occurs: 1
As A [user]	Researcher		

I Want To	Know the age range of individuals included within data
So That	I can understand whether data is suitable for my need

<b>Physical Sample Availability</b>	<p>Availability of physical samples associated with the dataset as Yes/No. If Yes provide descriptions of the samples available and process for access.</p> <p><b>Please provide this only if the dataset has any associated samples, otherwise omit this property.</b></p>	xs:string	Min Occurs: 0 Max Occurs: 1
As A [user]	Researcher		
I Want To	Know whether physical samples are available		
So That	I can understand if they can be requested for further research		

<b>Keywords</b>	<p>Definition: A keyword or tag describing the resource.</p> <p>Vocabulary: <a href="#">Data Catalog Vocabulary (DCAT)</a></p> <p>RDF Property: <a href="#">dcat:keyword</a></p> <p>Range: <a href="#">rdfs:Literal</a></p> <p>Source: <a href="#">dcat:resource_keyword</a></p>	xs:string	Min Occurs: 0 Max Occurs: 1
As A [user]	Researcher		
I Want To	know the keywords		
So That	I can navigate to it easily		

## Format and Structure

<b>Conforms To</b>	<p><b>Completion Guidance:</b> data standard such as OMOP or FHIR</p> <p><b>Definition:</b> An established standard to which the described resource conforms.  <i>"A basis for comparison; a reference point against which other things can be evaluated." [DCTERMS])</i></p> <p><i>Usage note: This property SHOULD be used to indicate the model, schema, ontology, view or profile that the catalogued resource content conforms to."</i></p> <p>Vocabulary: <a href="#">Data Catalog Vocabulary (DCAT)</a></p> <p>RDF Property: <a href="#">dct:conformsTo</a></p> <p>Range: <a href="#">dct:Standard</a></p> <p>Source: <a href="#">dcat:resource_conforms_to</a></p>	xs:string	Min Occurs: 0 Max Occurs: 1
As A [user]	Researcher		
I Want To	Know which standards the data set conforms to		
So That	I can use it to link to other datasets that use the same standards		

<b>Controlled Vocabulary</b>	<p><b>Completion Guidance:</b> Standard terminology / Ontology / Controlled Vocabulary such as ICD 10 Codes, NHS Data Dictionary or SNOMED CT International</p>	xs:string	Min Occurs: 0 Max Occurs: 1
------------------------------	---	-----------	--------------------------------



**Definition:** The nature or genre of the resource.

*Usage note: The value SHOULD be taken from a well governed and broadly recognised controlled vocabulary, such as: DCMI Type vocabulary [DCTERMS] [ISO-19115-1] scope codes DataCite resource types [DataCite] PARSE. Insight content-types used by re3data.org [RE3DATA-SCHEMA] (see item 15 contentType) MARC intellectual resource types Some members of these controlled vocabularies are not strictly suitable for datasets or data services (e.g. DCMI Type Event, PhysicalObject; [ISO-19115-1] CollectionHardware, CollectionSession, Initiative, Sample, Repository), but might be used in the context of other kinds of catalogues defined in DCAT profiles or applications. Usage note: To describe the file format, physical medium, or dimensions of the resource, use the dct:format element.*

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

RDF Property: [dct:type](#)

Sub-property of: [dct:type](#)

Range: [rdfs:Class](#)

Source: [dcat:resource\\_type](#)

As A [user]	Researcher
I Want To	Know which vocabulary does the dataset use or provide
So That	I can use it to link to other datasets that use the same vocabularies

**Definition:** A language of the item. This refers to the natural language used for textual metadata (i.e. titles, descriptions, etc) of a catalogued resource (i.e. dataset or service) or the textual values of a dataset distribution.

*Resources defined by the Library of Congress (ISO 639-1, ISO 639-2) SHOULD be used. If an ISO 639-1 (two-letter) code is defined for language, then its corresponding IRI SHOULD be used; if no ISO 639-1 code is defined, then IRI corresponding to the ISO 639-2 (three-letter) code SHOULD be used. Usage note: Repeat this property if the resource is available in multiple languages. Usage note: The value(s) provided for members of a catalogue (i.e. dataset or service) override the value(s) provided for the catalogue if they conflict. Usage note: If representations of a dataset are available for each language separately, define an instance of dct:Distribution for each language and describe the specific language of each distribution using dct:language (i.e. the dataset will have multiple dct:language values and each distribution will have just one as the value of its dct:language property).*

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

RDF Property: [dct:language](#)

Range: [dct:LinguisticSystem](#)

Source: [dcat:language](#)

xs:language      Min Occurs: 0  
Max Occurs: 1

## Language

As A [user]	Researcher
I Want To	Know the language of the dataset
So That	I understand the contents using that language

<b>Format</b>	<p><b>Completion Guidance:</b> If multiple formats are available please specify</p> <p><b>Definition:</b> The file format of the distribution.</p> <p><i>Usage note:</i> <i>dct:mediaType SHOULD be used if the type of the distribution is defined by IANA [IANA-MEDIA-TYPES].</i></p> <p>Vocabulary: <a href="#">Data Catalog Vocabulary (DCAT)</a></p> <p>RDF Property: <a href="#">dct:format</a></p> <p>Range: <a href="#">dct:MediaTypeOrExtent</a></p> <p>Source: <a href="#">dcat:format</a></p>	xs:string	Min Occurs: 0 Max Occurs: 1
	As A [user]	Researcher	
	I Want To	know the format	
	So That	I know how to analyse it	

<b>File size</b>	<p><b>Completion Guidance:</b> Typical extract file size range i.e. typically 1TB or 10GB or 100MB etc.</p> <p><i>Usage note:</i> <i>The size in bytes can be approximated (as a decimal) when the precise size is not known.</i></p> <p><b>Related Definition:</b> The size of a distribution in bytes.</p> <p>Vocabulary: <a href="#">Data Catalog Vocabulary (DCAT)</a></p> <p>RDF Property: <a href="#">dct:byteSize</a></p> <p>Domain: <a href="#">dct:Distribution</a></p> <p>Range: <a href="#">rdfs:Literal</a> typed as <a href="#">xsd:decimal</a></p> <p>Source: <a href="#">byteSize</a></p> <p>See also: <a href="#">Dataset/contentSize</a></p>	xs:string	Min Occurs: 0 Max Occurs: 1
	As A [user]	Researcher	
	I Want To	know the file size	
	So That	I can make provisions for storage, compute, transfer etc.	

## Attribution

<b>Creator</b>	<p><b>Completion Guidance:</b> Please provide the text that you would like included as part of any citation that credits this dataset.</p> <p><b>Definition:</b> The entity responsible for producing the resource.</p> <p><i>Usage note:</i> <i>Resources of type foaf:Agent are recommended as values for this property. See also: § 6.11 Class: Organization/Person</i></p> <p>Vocabulary: <a href="#">Data Catalog Vocabulary (DCAT)</a></p> <p>RDF Property: <a href="#">dct:creator</a></p> <p>Range: <a href="#">foaf:Agent</a></p> <p>Source: <a href="#">dcat:resource_creator</a></p> <p>See also: Schema.org <a href="#">Dataset/resource creator</a></p>	xs:string	Min Occurs: 0 Max Occurs: 1
	As A [user]	Researcher	
	I Want To	know the creator	
	So That	I can understand the context and for attribution	

<b>Citations</b>	<p><b>Completion Guidance:</b> Please provide a list of known citations, if available</p> <p><b>Definition:</b> A related resource, such as a publication, that references, cites, or otherwise points to the catalogued resource.</p> <p><i>Usage note: In relation to the use case of data citation, when the catalogued resource is a dataset, the <code>dct:isReferencedBy</code> property allows to relate the dataset to the resources (such as scholarly publications) that cite or point to the dataset. Multiple <code>dct:isReferencedBy</code> properties can be used to indicate the dataset has been referenced by multiple publications, or other resources.</i></p> <p><i>Usage note: This property is used to associate a resource with the resource (of type <code>dct:Resource</code>) in question. For other relations to resources not covered with this property, the more generic property <a href="#">dct:qualifiedRelation</a> can be used.</i></p> <p>Vocabulary: <a href="#">Data Catalog Vocabulary (DCAT)</a></p> <p>RDF Property: <a href="#">dct:isReferencedBy</a></p> <p>Source: <a href="#">dcat:resource_is_referenced_by</a></p>	xs:string	Min Occurs: 0 Max Occurs: *
	As A [user]	Researcher	
	I Want To	know how to cite the asset	
	So That	I can appropriately cite the resource in my research publications	

<b>DOI</b>	<p><b>Completion Guidance</b> Please provide a Digital Object Identifier if available</p> <p><b>Definition</b> Digital Object Identifier. Provides an actionable, interoperable, persistent link Actionable – through use of identifier syntax and network resolution mechanism (Handle System®) Persistent – through combination of supporting improved handle infrastructure (registry database, proxy support, etc) and social infrastructure (obligations by Registration Agencies) Interoperable – through use of a data model providing semantic interoperability and grouping mechanisms.</p> <p>Vocabulary: <a href="#">DOI Data Dictionary</a></p> <p>Source: <a href="#">DOI</a></p>	xs:string	Min Occurs: 0 Max Occurs: 1
	As A [user]	Researcher	
	I Want To	know the persistent link	
	So That	I can cite and discover the asset	

## Technical Metadata

In the phase 1 we will ask sites to provide a data dictionary of technical metadata that includes the following information:

<b><u>Table Name</u></b>	Name of the table in the dataset. Use a fully qualified name if appropriate.	xs:string	Min Occurs: 1 Max Occurs: 1
<b><u>Table Description</u></b>	Description of the table in the dataset.	xs:string	Min Occurs: 1 Max Occurs: 1
<b><u>Column Name</u></b>	Name of the column in the table dataset.	xs:string	Min Occurs: 1 Max Occurs: 1
<b><u>Column Description</u></b>	Description of the column in the table dataset.	xs:string	Min Occurs: 1 Max Occurs: 1
<b><u>Data Type</u></b>	Type of data contained in the column.	xs:string	Min Occurs: 1 Max Occurs: 1
<b>Sensitive</b>	<p><b>Completion guidance:</b> Please indicate (True / False) whether the information must be treated as sensitive and may need additional constraints / removal / anonymisation / masking through the data access request process.</p> <p><b>Definition:</b> An ODRL conformant policy expressing the rights associated with the resource.</p> <p><i>Usage note: Information about rights expressed as an ODRL policy [ODRL-MODEL] using the ODRL vocabulary [ODRL-VOCAB] MAY be provided for the resource.</i></p> <p>Vocabulary: <a href="#">ODRL Vocabulary &amp; Expression</a></p> <p>RDF Property: <a href="#">odrl:hasPolicy</a></p> <p>Range: <a href="#">odrl:Policy</a></p> <p>See also guidance at <a href="#">§ 8 License and rights statements, § 6.4.19 Property: license</a>, <a href="#">§ 6.4.1 Property: access rights, § 6.4.20 Property: rights</a></p>	xs:boolean	Max Occurs: 1 Min Occurs: 1