

# Building & Mining Knowledge Graphs

(KEN4256)

## Lecture 6: Publishing Knowledge Graphs



Maastricht University

Institute of Data Science

© 2024 by Michel Dumontier and the Institute of Data Science at Maastricht University is licensed under Attribution 4.0 International  
To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>

This license requires that reusers give credit to the creator. It allows reusers to distribute, remix, adapt, and build upon the material in any medium or format, even for commercial purposes.

id: KEN4256\_L6

version: 1.2024.0

created: January 13, 2019

last modified: March 26, 2024

published on: March 26, 2024

*Significant* effort is needed to find the right information, make sense of it, understand your rights and responsibilities, and ultimately reuse them for a new purpose



<https://youtu.be/N2zK3sAtr-4>

slido

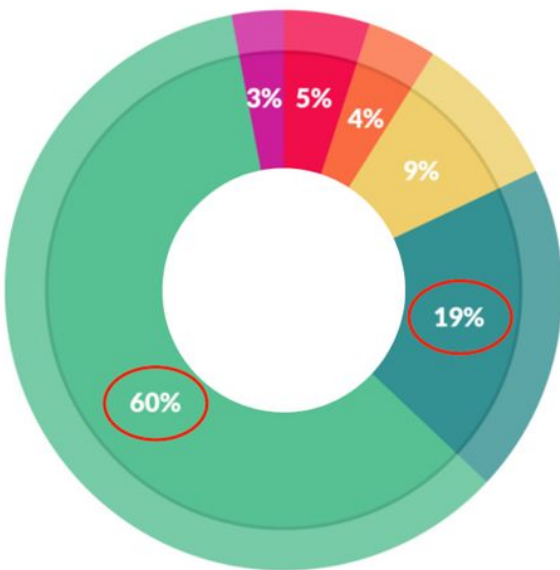


## What problems were encountered

① Click **Present with Slido** or install our [Chrome extension](#) to activate this poll while presenting.

# What problems were there?

- Limited data availability
- Inefficient data access
  - Unreliable (non-replicated) storage
- Proprietary data format
- Non-standardized metadata
  - No description of the data elements
  - Lack of contact information
  - No traceability of data from publication to source

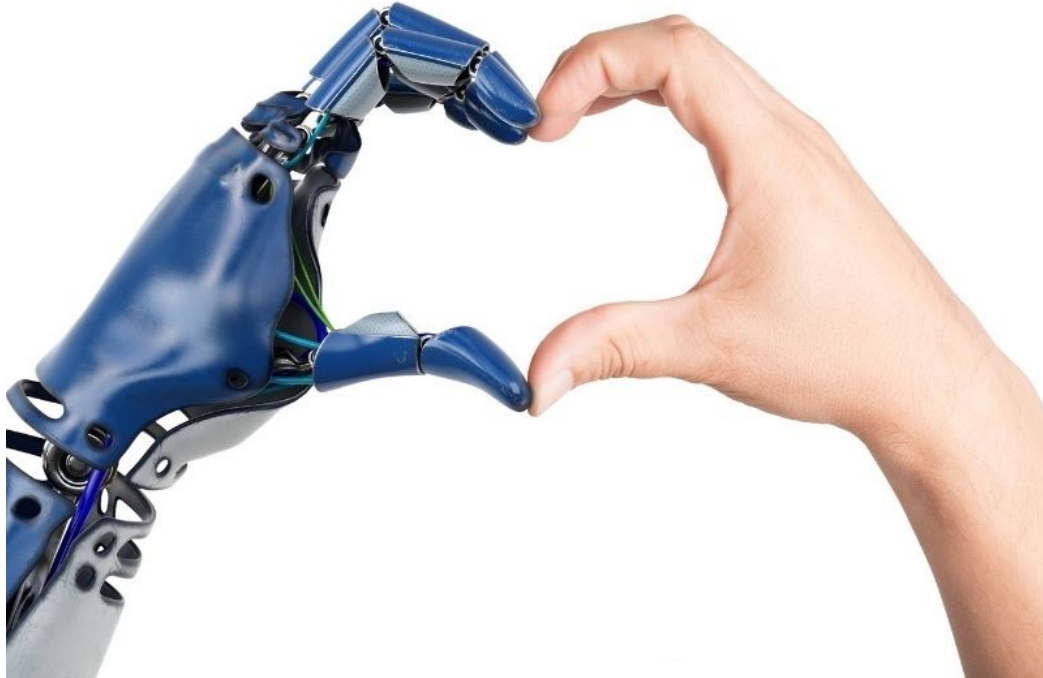


### What data scientists spend the most time doing

- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets; 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
- Other: 5%

[http://visit.crowdfunder.com/rs/416-ZBE-142/images/CrowdFunder\\_DataScienceReport\\_2016.pdf](http://visit.crowdfunder.com/rs/416-ZBE-142/images/CrowdFunder_DataScienceReport_2016.pdf)

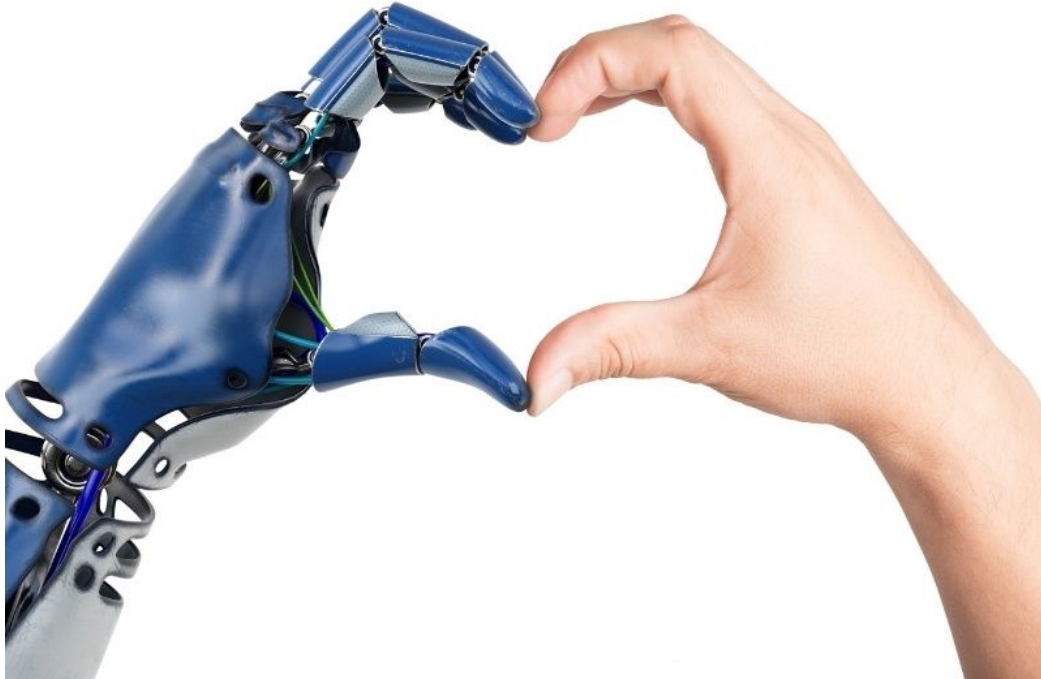
# Human Machine collaboration is crucial to our future work

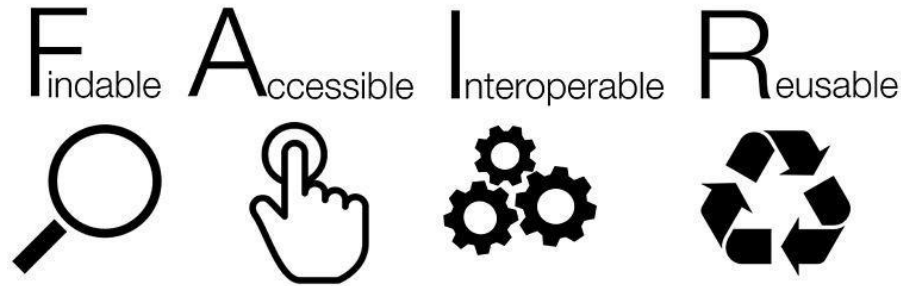




# Machines

need to be able to discover and reuse data  
(and arguably any digital resource)





**An international, bottom-up paradigm for  
the discovery and reuse of digital content  
*by and for people and machines***

## The FAIR Guiding Principles for scientific data management and stewardship

Mark D. Wilkinson, Michel Dumontier [...] Barend Mons

[Affiliations](#) | [Contributions](#) | [Corresponding author](#)

*Scientific Data* 3, Article number: 160018 (2016) | doi:10.1038/sdata.2016.18

Received 10 December 2015 | Accepted 12 February 2016 | Published online 15 March 2016

This article is in the 99<sup>th</sup> percentile (ranked 59<sup>th</sup>) of the 266,804 tracked articles of a similar age in all journals and the 1<sup>st</sup> percentile (ranked 1<sup>st</sup>) of the 1 tracked articles of a similar age in *Scientific Data*

Endorsed by publishers, industry partners, funding agencies

Principles to enhance the value of all digital resources: data, software, repositories, **knowledge graphs**

<http://www.nature.com/articles/sdata201618>



European Commission  
Press Release Database

European Commission > Press releases database > Press Release details

**European Commission - Statement**

**G20 Leaders' Communique Hangzhou Summit**

Hangzhou, 5 September 2016

1. We, the Leaders of the G20, met in Hangzhou, China on 4-5 September 2016.

 **Annex 4:**  
**G7 Expert Group on Open Science**

Turin, Italy, September 28, 2017



## Box 2 | The FAIR Guiding Principles

### To be Findable:

- F1. (meta)data are assigned a globally unique and persistent identifier
- F2. data are described with rich metadata (defined by R1 below)
- F3. metadata clearly and explicitly include the identifier of the data it describes
- F4. (meta)data are registered or indexed in a searchable resource

### To be Accessible:

- A1. (meta)data are retrievable by their identifier using a standardized communications protocol
  - A1.1 the protocol is open, free, and universally implementable
  - A1.2 the protocol allows for an authentication and authorization procedure, where necessary
- A2. metadata are accessible, even when the data are no longer available

### To be Interoperable:

- I1. (meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation.
- I2. (meta)data use vocabularies that follow FAIR principles
- I3. (meta)data include qualified references to other (meta)data

### To be Reusable:

- R1. meta(data) are richly described with a plurality of accurate and relevant attributes
  - R1.1. (meta)data are released with a clear and accessible data usage license
  - R1.2. (meta)data are associated with detailed provenance
  - R1.3. (meta)data meet domain-relevant community standards

<http://www.nature.com/articles/sdata201618>

4 facets (F,A,I,R) to the 15 principles.

# Making FAIR Data

## Collect

Data

## Describe

Standardized  
Metadata

use standard  
metadata format

use ontologies +  
vocabularies

add provenance  
and license

## Transform

Standardized  
Data

Use standard  
data format

use ontologies +  
vocabularies

## Publish

FAIR Data

Data Repository

Persistent Metadata  
Identifier

Standardized  
Metadata

Persistent Data Identifier

Standardized  
Data

# applies to data and their metadata

## Box 2 | The FAIR Guiding Principles

### To be Findable:

- F1. (meta)data are assigned a globally unique and persistent identifier
- F2. data are described with rich metadata (defined by R1 below)
- F3. metadata clearly and explicitly include the identifier of the data it describes
- F4. (meta)data are registered or indexed in a searchable resource

### To be Accessible:

- A1. (meta)data are retrievable by their identifier using a standardized communications protocol
  - A1.1 the protocol is open, free, and universally implementable
  - A1.2 the protocol allows for an authentication and authorization procedure, where necessary
- A2. metadata are accessible, even when the data are no longer available

### To be Interoperable:

- I1. (meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation.
- I2. (meta)data use vocabularies that follow FAIR principles
- I3. (meta)data include qualified references to other (meta)data

### To be Reusable:

- R1. meta(data) are richly described with a plurality of accurate and relevant attributes
  - R1.1. (meta)data are released with a clear and accessible data usage license
  - R1.2. (meta)data are associated with detailed provenance
  - R1.3. (meta)data meet domain-relevant community standards

# Data vs Metadata

**Data** are information made available for machine processing. Data can be quantitative (e.g. height, weight), qualitative (how well do you feel on a scale of 1-5), and multi-modal (text, relational, image, sound, video).

**Metadata** are information about data. They may contain a description, context, provenance, and meaning of the data.

# Metadata

- What is the name or title of the digital resource?
- What is the digital resource about?
- Who contributed to creating or maintaining the digital resource?
- When was it created, modified, released?
- What methodology or tool was used?
- Which language is used?
- Which formats is it available as?
- What license is it released under?



# Metadata



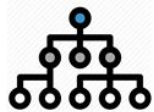
## Administrative metadata

- Provides information about the digital artifact
  - date created, size of file, type of file
  - license, permissions, etc



## Descriptive metadata

- Contains information relevant to find and interpret the data
  - title, description, keywords
  - datatype-specific metadata e.g. protocol, instrument



## Relational metadata

- Captures the relationship between the data item and the entity it is about or is in a contextual relation with
  - The patient for which the MRI scan was taken

# Data and their metadata



Data: jpg image file



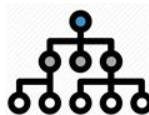
**Informative metadata:**

- Size: 155kb
- Date created: 2015-05-25
- Filetype: jpg



**Descriptive metadata**

- Title: MRI of the head
- Device: Ingenia 3.0T



**Relational metadata**

- About: EHR092376573
- Clinical Study: CT7812356

# (meta)data for machines

The long term objective of FAIR is to make content accessible to machines, to support the everyday work we do

**Data and their metadata** ought to be:

- machine readable - the syntax of the data are formally specified to enable reliable reading/writing of the data.
- machine interpretable - the semantics of the data elements are well defined and can be reasoned about for information retrieval and query answering

# (meta)data for machines

- find relevant digital objects that are published on the web
- facilitate **query** and **filter** content based on specific features of interest.
- easier to **understand** and **compare** data and their provenance
- easier to **replicate** experiments and **reproduce** research results
- easier to **integrate** independently produced data
- **exchange** content between different tools and environments
- **explore** and **visualize** knowledge connections



Technical infrastructure (generic operations)

Social agreements/contracts (domain-specific content)

## Box 2 | The FAIR Guiding Principles

### To be Findable:

- F1. (meta)data are assigned a globally unique and persistent identifier
- F2. data are described with rich metadata (defined by R1 below)
- F3. metadata clearly and explicitly include the identifier of the data it describes
- F4. (meta)data are registered or indexed in a searchable resource

### To be Accessible:

- A1. (meta)data are retrievable by their identifier using a standardized communications protocol
  - A1.1 the protocol is open, free, and universally implementable
  - A1.2 the protocol allows for an authentication and authorization procedure, where necessary
- A2. metadata are accessible, even when the data are no longer available

### To be Interoperable:

- I1. (meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation.
- I2. (meta)data use vocabularies that follow FAIR principles
- I3. (meta)data include qualified references to other (meta)data

### To be Reusable:

- R1. meta(data) are richly described with a plurality of accurate and relevant attributes
  - R1.1. (meta)data are released with a clear and accessible data usage license
  - R1.2. (meta)data are associated with detailed provenance
  - R1.3. (meta)data meet domain-relevant community standards

# **F1: (meta) data are assigned globally unique and persistent identifiers**

**A unique identifier unambiguously refers a single resource, cannot be confused with another resource, and is never re-assigned to another entity.**

An identifier used to denote multiple resources will confound efforts to make describes any one resource, and will lead to confusion when trying to retrieve a corresponding representation.

**A persistent identifier is one where there is a technological capability and a corresponding social commitment to ensure continued access to the resource.**

for instance, the responsibility to provide the resource can be transferred to another party without having to change the identifier, or there is redundancy in network to provide copies from other providers.

# Globally unique and persistent?

26978244

pubmed:26978244

<https://pubmed.ncbi.nlm.nih.gov/26978244/>

doi:10.1038/sdata.2016.18

<https://doi.org/10.1038/sdata.2016.18>

# Globally unique and persistent identifiers

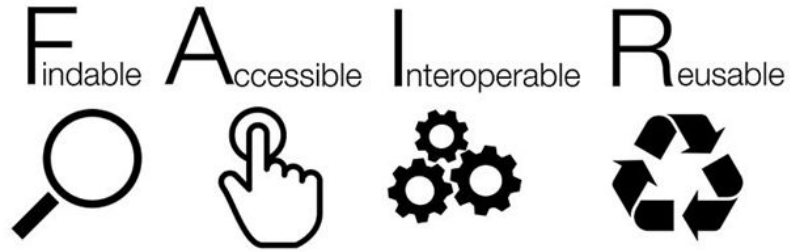
Obtain globally unique and persistent identifiers from a software or service that uses algorithms that can detect changes in the content.

- Persistent URLs: <http://www.purlz.org>
- Identifiers.org: <http://identifiers.org>
- Digital Object Identifier: <http://www.doi.org>
- Archival Resource Key (ARK): [http://n2t.net/e/ark\\_ids.html](http://n2t.net/e/ark_ids.html)
- Global research identifiers: <https://www.grid.ac>

Globally unique and persistent identifiers owing to content (via digital fingerprint)

- Data GUIDs <https://dataguids.org>
- Trusty URIs: <http://trustyuri.net/>

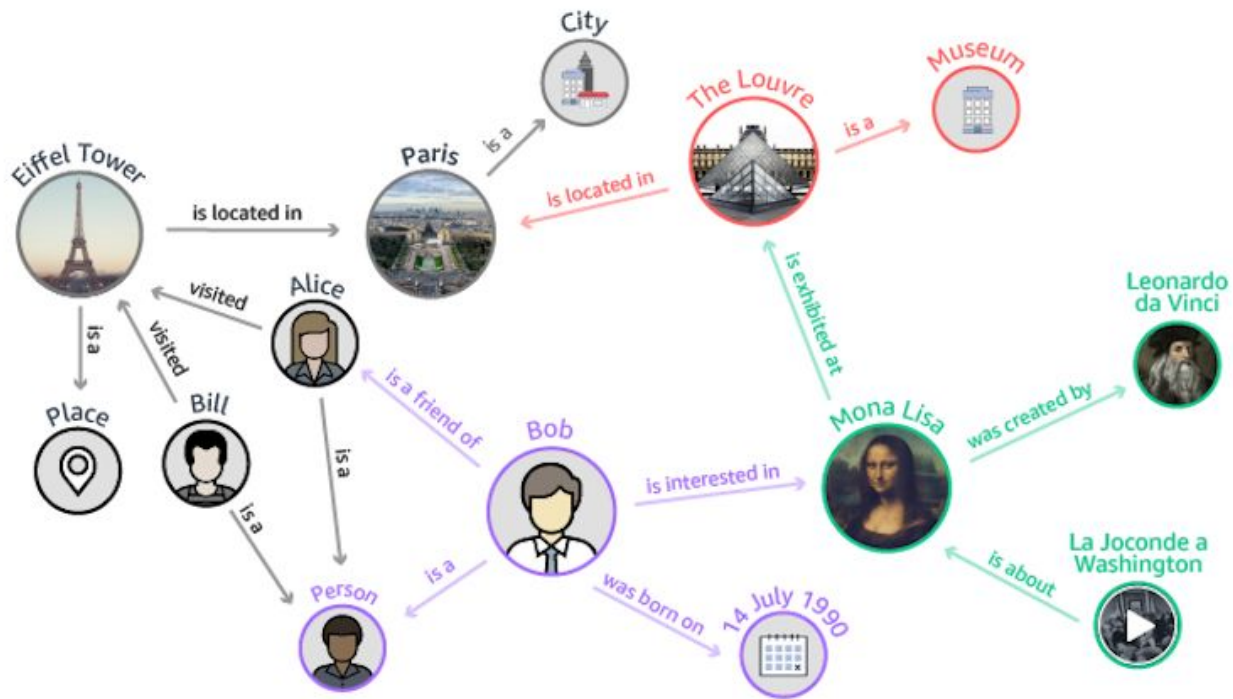




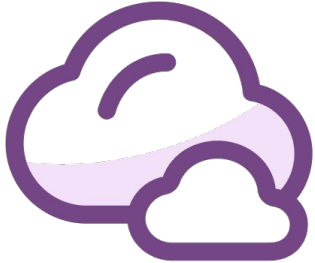
## Principles to enhance the value of *all* digital resources

*data, software, web services, repositories,...* **knowledge graphs!**

Developed and endorsed by researchers, publishers, funding agencies, industry partners.



slido



**Which metadata could be included to describe a KG?**

① Start presenting to display the poll results on this slide.

# Metadata

- What is the name or title of the digital resource?
- What is the digital resource about?
- Who contributed to creating or maintaining the digital resource?
- When was it created, modified, released?
- What methodology or tool was used?
- Which language is used?
- Which formats is it available as?
- What license is it released under?
- Which descriptive or quality metrics are available?
- Who is using it?

## Dataset Descriptions: HCLS Community Profile

### Editors working draft.

#### Editors:

Alasdair J.G. Gray, Heriott-Watt University, UK <[A.J.G.Gray@hw.ac.uk](mailto:A.J.G.Gray@hw.ac.uk)>  
 Joachim Baran, Stanford University, USA <[joachim.baran@stanford.edu](mailto:joachim.baran@stanford.edu)>  
 M. Scott Marshall, MAASTRO Clinic, The Netherlands <[m.scott.marshall@maastro.nl](mailto:m.scott.marshall@maastro.nl)>  
 Michel Dumontier, Stanford University, USA <[michel.dumontier@stanford.edu](mailto:michel.dumontier@stanford.edu)>

#### Contributors:

Vladimir Alexiev, Ontotext Corp, Bulgaria <[vladimir.alexiev@ontotext.com](mailto:vladimir.alexiev@ontotext.com)>  
 Peter Ansell, CSIRO, Australia <[peter.ansell@csiro.au](mailto:peter.ansell@csiro.au)>  
 Gary D. Bader, The Donnelly Centre, University of Toronto, Canada <[gary.bader@utoronto.ca](mailto:gary.bader@utoronto.ca)>  
 Asuka Bando, NBDC, Japan <[bando@biosciencedbc.jp](mailto:bando@biosciencedbc.jp)>  
 Jerven Bolleman, SIB Swiss Institute of Bioinformatics, Switzerland <[jerven.bolleman@isb-sib.ch](mailto:jerven.bolleman@isb-sib.ch)>  
 Alison Callahan, Carleton University, Canada <[alison.callahan@carleton.ca](mailto:alison.callahan@carleton.ca)>  
 José Cruz-Toledo, Carleton University, Canada <[josecruztoledo@gmail.com](mailto:josecruztoledo@gmail.com)>  
 Pascale Gaudet, SIB Swiss Institute of Bioinformatics, Switzerland <[pascale.gaudet@isb-sib.ch](mailto:pascale.gaudet@isb-sib.ch)>  
 Erich Gombocz, IO Informatics, USA <[egombocz@io-informatics.com](mailto:egombocz@io-informatics.com)>  
 Alejandra Gonzalez-Beltran, University of Oxford, UK <[alejandra.gonzalez.beltran@gmail.com](mailto:alejandra.gonzalez.beltran@gmail.com)>  
 Paul Groth, VU University Amsterdam, The Netherlands <[p.t.groth@vu.nl](mailto:p.t.groth@vu.nl)>  
 Melissa Haendel, Oregon Health and Science University, USA <[haendel@ohsu.edu](mailto:haendel@ohsu.edu)>  
 Maori Ito, NIBIO, Japan <[maori@nibio.go.jp](mailto:maori@nibio.go.jp)>  
 Simon Jupp, EMBL-EBI, UK <[jupp@ebi.ac.uk](mailto:jupp@ebi.ac.uk)>  
 Nick Juty, EMBL-EBI, UK <[juty@ebi.ac.uk](mailto:juty@ebi.ac.uk)>  
 Toshiaki Katayama, Database Center for Life Sciences, Japan <[ktym@dbcls.jp](mailto:ktym@dbcls.jp)>  
 Norio Kobayashi, RIKEN, Japan <[norio.kobayashi@riken.jp](mailto:norio.kobayashi@riken.jp)>  
 Kalpana Krishnaswami, Metaome, USA <[kalpana@metaome.com](mailto:kalpana@metaome.com)>  
 Camille Laibe, EMBL-EBI, UK <[laibe@ebi.ac.uk](mailto:laibe@ebi.ac.uk)>  
 Nicolas Le Novère, Babraham Institute, UK <[n.lenovere@gmail.com](mailto:n.lenovere@gmail.com)>  
 Simon Lin, Marshfield Clinic Research Foundation, USA <[lin.simon@mcrf.mfldclin.edu](mailto:lin.simon@mcrf.mfldclin.edu)>  
 James Malone, EMBL-EBI, UK <[malone@ebi.ac.uk](mailto:malone@ebi.ac.uk)>  
 Michael Miller, Institute for Systems Biology, USA <[mmiller@systemsbiology.org](mailto:mmiller@systemsbiology.org)>  
 Chris Mungall, Lawrence Berkeley National Laboratory, USA <[cjm@berkeleylab.org](mailto:cjm@berkeleylab.org)>  
 Laurens Rietveld, VU University Amsterdam, The Netherlands <[laurens.rietveld@vu.nl](mailto:laurens.rietveld@vu.nl)>  
 Sarala M. Wimalaratne, EMBL-EBI, UK <[sarala@ebi.ac.uk](mailto:sarala@ebi.ac.uk)>  
 Atsuko Yamaguchi, Database Center for Life Sciences, Japan <[atsuko@dbcls.jp](mailto:atsuko@dbcls.jp)>

<http://www.w3.org/TR/hcls-dataset/>

## A guide to describing data with RDF vocabularies

- Identifiers
- Descriptors
- Versioning
- Attribution
- Provenance
- Content summarization

Mandatory, recommended, optional descriptors  
 Reference editor and validation

# Metadata element, description, and example of use

## 6.2.2 Title

At least one human-readable title should be provided for a dataset using `dct:title`. Alternative or older titles may be specified using `dct:alternative`.

For example, to provide a title and alternative title for the ChEMBL dataset:

```
:chembl
  dct:title "ChEMBL"@en ;
  dct:alternative "ChEMBLdb"@en ;
```

# Basic Description

- Identifiers
- Title
- Description
- Homepage
- License
- Language
- Keywords
- Concepts and vocabularies used
- Standards
- Reference Documentation/Publication
- Format
- Download URL
- Landing page
- SPARQL endpoint

Prefix	URI	Description
cito:	<a href="http://purl.org/spar/cito/">http://purl.org/spar/cito/</a>	<a href="#">Citation Typing Ontology</a>
dcat:	<a href="http://www.w3.org/ns/dcat#">http://www.w3.org/ns/dcat#</a>	<a href="#">Data Catalog</a>
dctypes:	<a href="http://purl.org/dc/dcmitype/">http://purl.org/dc/dcmitype/</a>	<a href="#">Dublin Core Metadata Types</a>
dct:	<a href="http://purl.org/dc/terms/">http://purl.org/dc/terms/</a>	<a href="#">Dublin Core Metadata Terms</a>
foaf:	<a href="http://xmlns.com/foaf/0.1/">http://xmlns.com/foaf/0.1/</a>	<a href="#">Friend-of-a-Friend</a>
freq:	<a href="http://purl.org/cld/freq/">http://purl.org/cld/freq/</a>	<a href="#">Collection Description Frequency Vocabulary</a>
idot:	<a href="http://identifiers.org/idot/">http://identifiers.org/idot/</a>	<a href="#">Identifiers.org vocabulary</a>
lexvo:	<a href="http://lexvo.org/ontology#">http://lexvo.org/ontology#</a>	<a href="#">Lexical Vocabulary</a>
pav:	<a href="http://purl.org/pav/">http://purl.org/pav/</a>	<a href="#">Provenance Authoring and Versioning ontology</a>
prov:	<a href="http://www.w3.org/ns/prov#">http://www.w3.org/ns/prov#</a>	<a href="#">PROV Ontology</a>
rdf:	<a href="http://www.w3.org/1999/02/22-rdf-syntax-ns#">http://www.w3.org/1999/02/22-rdf-syntax-ns#</a>	<a href="#">RDF Syntax</a>
rdfs:	<a href="http://www.w3.org/2000/01/rdf-schema#">http://www.w3.org/2000/01/rdf-schema#</a>	<a href="#">RDF Schema</a>
schemaorg:	<a href="http://schema.org/">http://schema.org/</a>	<a href="#">schema.org vocabulary</a>
sd:	<a href="http://www.w3.org/ns/sparql-service-description#">http://www.w3.org/ns/sparql-service-description#</a>	<a href="#">SPARQL 1.1 Service Description</a>
sio:	<a href="http://semanticscience.org/resource/">http://semanticscience.org/resource/</a>	<a href="#">Semanticscience Integrated Ontology (SIO)</a>
xsd:	<a href="http://www.w3.org/2001/XMLSchema#">http://www.w3.org/2001/XMLSchema#</a>	<a href="#">XML Schema</a>
void:	<a href="http://www.w3.org/TR/void/">http://www.w3.org/TR/void/</a>	<a href="#">Describing Linked Datasets with the VoID Vocabulary</a>
void-ext:	<a href="http://ldf.fi/void-ext">http://ldf.fi/void-ext</a>	<a href="#">Extensions to the Vocabulary of Interlinked Datasets (VoID)</a>



# Statistics

Statistics	
# of triples	void:triples
# of typed entities	void:entities
# of subjects	void:distinctSubjects
# of properties	void:properties
# of objects	void:distinctObjects
# of classes	void:classPartition
# of literals	void:classPartition
# of RDF graphs	void:classPartition
class frequency	void:classPartition
property frequency	void:propertyPartition
property and subject types	void:propertyPartition
property and object types	void:propertyPartition
property and literals	void:propertyPartition
property subject and object types	void:propertyPartition

6.6.1.1 To specify the **number of triples** in the dataset:

Pattern:

```
:rdfdataset  
  void:triples "###"^^xsd:integer .
```

Example:

```
:chembl17rdf  
  void:triples "409942525"^^xsd:integer .
```

SPARQL query:

```
SELECT (COUNT(*) AS ?triples)  
{ ?s ?p ?o }
```

# Metagraph

Pattern:

```

:rdfdataset
  void:propertyPartition [
    void:property <property-uri> ;
    void:classPartition [
      void:class <subject-class-uri> ;
      void:distinctSubjects "###"^^xsd:integer ;
    ];
    void-ext:objectClassPartition [
      void:class <object-class-uri> ;
      void:distinctObjects "###"^^xsd:integer ;
    ];
  ] .

```

Example:

```

:chembl17rdf
  void:propertyPartition [
    void:property <http://rdf.ebi.ac.uk/terms/chembl#hasAssay> ;
    void:triples "12419715"^^xsd:integer ;
    void:classPartition [
      void:class <http://rdf.ebi.ac.uk/terms/chembl#Activity> ;
      void:distinctSubjects "12419715"^^xsd:integer ;
    ];
  ];
  void-ext:objectClassPartition [
    void:class <http://rdf.ebi.ac.uk/terms/chembl#Assay> ;
    void:distinctObjects "1042288"^^xsd:integer ;
  ] .

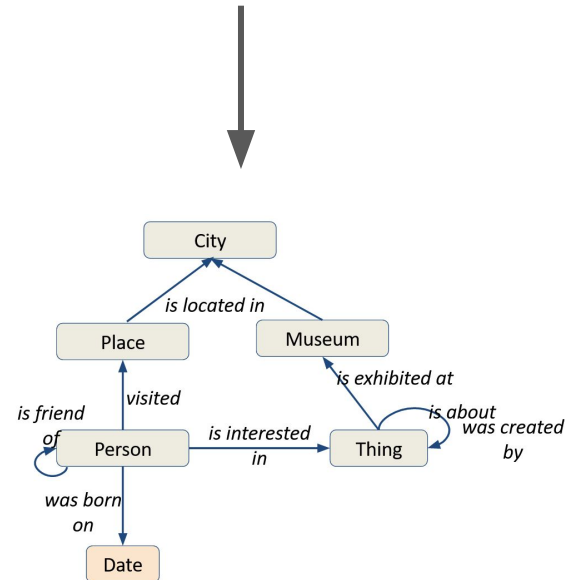
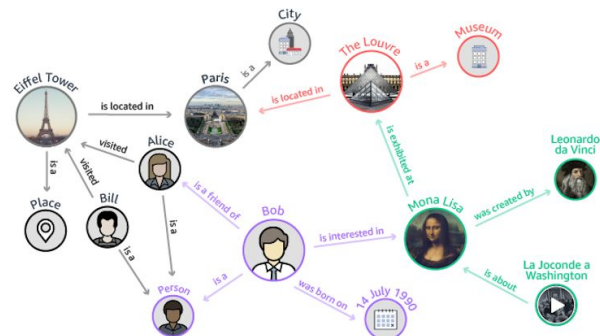
```

SPARQL query:

```

SELECT (COUNT(DISTINCT ?s) AS ?scount) ?stype ?p ?otype (COUNT(DISTINCT ?o) AS ?ocount)
{
  ?s ?p ?o .
  ?s a ?stype .
  ?o a ?otype .
} GROUP BY ?p ?stype ?otype

```





## Describing Linked Datasets with the VoID Vocabulary

W3C Interest Group Note 03 March 2011

**This version:**

<http://www.w3.org/TR/2011/NOTE-void-20110303/>

**Latest version:**

<http://www.w3.org/TR/void/>

```
:DBpedia a void:Dataset;
  dcterms:title "DBPedia";
  dcterms:description "RDF data extracted from Wikipedia";
  dcterms:contributor :FU_Berlin;
  dcterms:contributor :University_Leipzig;
  dcterms:contributor :OpenLink_Software;
  dcterms:contributor :DBpedia_community;
  dcterms:source <http://dbpedia.org/resource/Wikipedia>;
  dcterms:modified "2008-11-17"^^xsd:date;
.
:FU_Berlin a foaf:Organization;
  rdfs:label "Freie Universität Berlin";
  foaf:homepage <http://www.fu-berlin.de/>;
.
# Similar descriptions of the other contributors go here
```

## Table of Contents

### [1. Introduction](#)

#### [1.1 Scope](#)

#### [1.2 Document conventions](#)

#### [1.3 Definition: Dataset](#)

#### [1.4 Definition: Linkset](#)

### [2. General dataset metadata](#)

#### [2.1 Web page links](#)

#### [2.2 Basic Dublin Core metadata](#)

#### [2.3 Contact information](#)

#### [2.4 Announcing the license of a dataset](#)

#### [2.5 Categorizing datasets by subject](#)

#### [2.6 Technical features](#)

### [3. Access metadata](#)

#### [3.1 Resolvable HTTP URIs](#)

#### [3.2 SPARQL endpoints](#)

#### [3.3 RDF data dumps](#)

#### [3.4 Root resources](#)

#### [3.5 URI lookup endpoints](#)

#### [3.6 OpenSearch description documents](#)

### [4. Structural metadata](#)

#### [4.1 Example resources](#)

#### [4.2 Patterns for resource URIs](#)

#### [4.3 Vocabularies used in a dataset](#)

#### [4.4 Describing partitioned datasets](#)

#### [4.5 Partitioning a dataset based on classes and properties](#)

#### [4.6 Providing statistics about datasets](#)

### [5. Describing linksets](#)

#### [5.1 Naming a linkset's two target datasets](#)

#### [5.2 Linksets as part of larger datasets](#)

#### [5.3 Stating the link predicate of a linkset](#)

### [6. Deploying VoID descriptions](#)

#### [6.1 Choosing URIs for datasets](#)

#### [6.2 Publishing a VoID file alongside a dataset](#)

#### [6.3 Multi-document datasets and backlinks](#)

#### [6.4 Describing RDF dumps](#)

#### [6.5 Using VoID with the SPARQL Service Description Vocabulary](#)

### [7. Discovering VoID descriptions](#)

#### [7.1 Discovery via links in the dataset's documents](#)

#### [7.2 Discovery with well-known URI](#)

### [8. Index of VoID classes and properties](#)

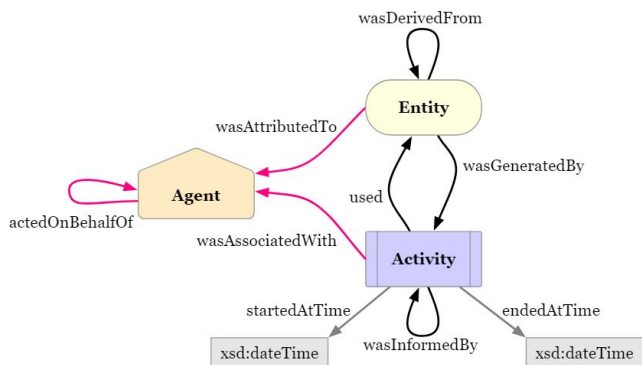
### [9. Acknowledgements](#)

### [References](#)

W3C Recommendation 30 April 2013

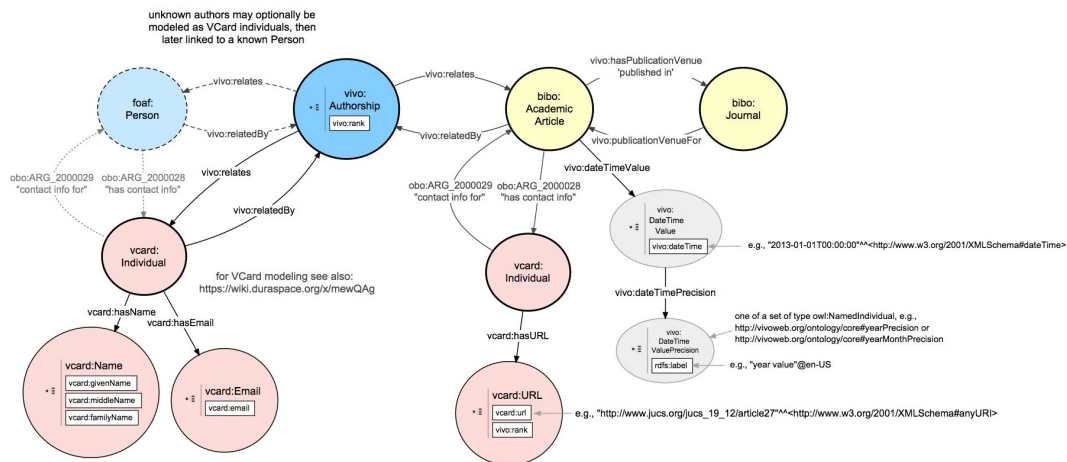
<http://www.w3.org/TR/2013/REC-prov-o-20130430/>

<http://www.w3.org/TR/prov-o/>



8/26/15

## VIVO Authorship: Connecting an Author with a Publication



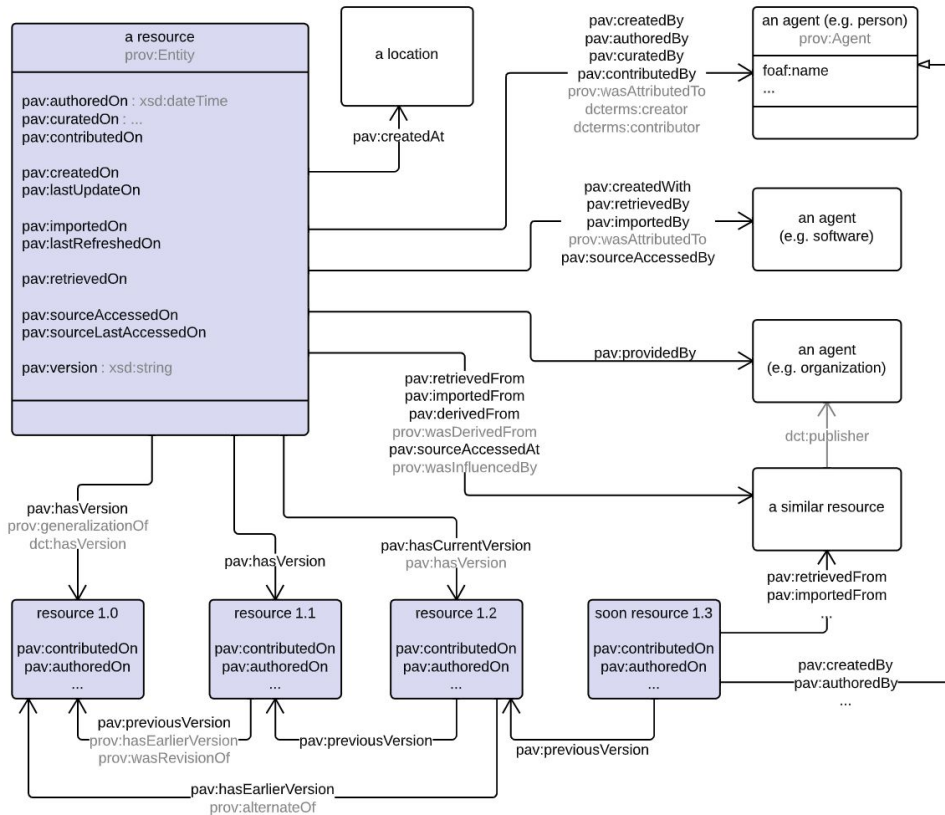
## PAV - Provenance, Authoring and Versioning

**IRI:**

<http://purl.org/pav/>

Version IRI:

<http://purl.org/pav/2.3>



<http://bioschemas.org>

Key to specification table

Schema.org properties where the Expected Types have been changed, or new (i.e., Bioschemas created) properties/types are **green**.

Schema.org properties/types are *red*.

Pending Schema.org properties/types are *blue*.

External (i.e., from 3rd party ontology) properties/types are *black*.

CD = Cardinality

Property	Expected Type	Description	CD	Controlled Vocabulary	Example
<b>Marginality: Minimum.</b>					
<a href="#"><u>description</u></a>	<a href="#"><u>Text</u></a>	<p><b>Schema:</b> A description of the item.</p> <p><b>Bioschemas:</b> A short summary describing a dataset.</p>	ONE		
<a href="#"><u>identifier</u></a>	<a href="#"><u>Property/Value</u></a> <a href="#"><u>Text</u></a> <a href="#"><u>URL</u></a>	<p><b>Schema:</b> The identifier property represents any kind of identifier for any kind of Thing, such as ISBNs, GTIN codes, UUIDs etc. Schema.org provides dedicated properties for representing many of these, either as textual strings or as URL (URI) links. See <a href="#">background notes</a> for more details.</p>	MANY		
<a href="#"><u>keywords</u></a>	<a href="#"><u>Text</u></a>	<p><b>Schema:</b> Keywords or tags used to describe this content. Multiple entries in a keywords list are typically delimited by commas.</p> <p><b>Bioschemas:</b> These keywords provide a summary of the dataset.</p>	ONE		
<a href="#"><u>name</u></a>	<a href="#"><u>Text</u></a>	<p><b>Schema:</b> The name of the item.</p> <p><b>Bioschemas:</b> A descriptive name of the dataset.</p>	ONE		
<a href="#"><u>rdf:type</u></a>	<a href="#"><u>URL</u></a>	<p><b>Bioschemas:</b> This is used by validation tools to identify the profile used. You must use the value specified in the Controlled Vocabulary column.</p>	ONE		
<a href="#"><u>url</u></a>	<a href="#"><u>URL</u></a>	<p><b>Schema:</b> URL of the item.</p> <p><b>Bioschemas:</b> The location of a page describing the dataset.</p>	ONE		
<b>Marginality: Recommended.</b>					
<a href="#"><u>citation</u></a>	<a href="#"><u>Creative/Work</u></a> <a href="#"><u>Text</u></a>	<p><b>Schema:</b> A citation or reference to another creative work, such as another publication, web page, scholarly article, etc.</p>	MANY		



## Dataset

A *Schema.org* Type

Thing > CreativeWork > Dataset

[more...]

A body of structured information describing some topic(s) of interest.

Property	Expected Type	Description
<b>Properties from Dataset</b>		
<b>distribution</b>	DataDownload	A downloadable form of this dataset, at a specific location, in a specific format. This property can be repeated if different variations are available. There is no expectation that different downloadable distributions must contain exactly equivalent information (see also DCAT on this point). Different distributions might include or exclude different subsets of the entire dataset, for example.
<b>includedInDataCatalog</b>	DataCatalog	A data catalog which contains this dataset. Supersedes includedDataCatalog, catalog. Inverse property: dataset
<b>issn</b>	Text	The International Standard Serial Number (ISSN) that identifies this serial publication. You can repeat this property to identify different formats of, or the linking ISSN (ISSN-L) for, this serial publication.
<b>measurementMethod</b>	DefinedTerm or MeasurementMethodEnum or Text or URL	A subproperty of measurementTechnique that can be used for specifying specific methods, in particular via MeasurementMethodEnum.
<b>measurementTechnique</b>	DefinedTerm or MeasurementMethodEnum or Text or URL	<p>A technique, method or technology used in an Observation, StatisticalVariable or Dataset (or DataDownload, DataCatalog), corresponding to the method used for measuring the corresponding variable(s) (for datasets, described using variableMeasured; for Observation, a StatisticalVariable). Often but not necessarily each variableMeasured will have an explicit representation as (or mapping to) an property such as those defined in Schema.org, or other RDF vocabularies and "knowledge graphs". In that case the subproperty of variableMeasured called measuredProperty is applicable.</p> <p>The measurementTechnique property helps when extra clarification is needed about how a measuredProperty was measured. This is oriented towards scientific and scholarly dataset publication but may have broader applicability, it is not intended as a full representation of measurement, but can often serve as a high level summary for dataset discovery.</p> <p>For example, if variableMeasured is: molecule concentration, measurementTechnique could be: "mass spectrometry" or "nmr spectroscopy" or "colorimetry" or "immunofluorescence". If the variableMeasured is "depression rating", the measurementTechnique could be "Zung Scale" or "HAM-D" or "Beck Depression Inventory".</p> <p>If there are several variableMeasured properties recorded for some given data object, use a PropertyValue for each variableMeasured and attach the corresponding measurementTechnique. The value can also be from an enumeration, organized as a MeasurementMethodEnumeration.</p>
<b>variableMeasured</b>	Property or PropertyValue or StatisticalVariable or Text	The variableMeasured property can indicate (repeated as necessary) the variables that are measured in some dataset, either described as text or as pairs of identifier and description using PropertyValue, or more explicitly as a StatisticalVariable.
<b>Properties from CreativeWork</b>		
<b>about</b>	Thing	The subject matter of the content. Inverse property: subjectOf
<b>abstract</b>	Text	An abstract is a short description that summarizes a CreativeWork.



# Dataset Search



## Wikidata

[Explore at www.wikidata.org](http://www.wikidata.org)[Explore at m.wikidata.nym.sk](http://m.wikidata.nym.sk)[Explore at unblocked.to](http://unblocked.to)[Explore at the Datahub](#)[Explore at data.wu.ac.at](http://data.wu.ac.at)

6 scholarly articles cite this data set ([View in Google Scholar](#))

### Data set provided by

[Wikimedia Foundation](#)

### Licence

[CC0 1.0 Universal Public Domain Dedication](#)

Licence information was derived automatically

### Available download formats from providers

sparql endpoint, full rdf turtle dump, full json dump, simplified ("truthy") rdf n-triples dump

### Description

Wikidata offers a wide range of general data about our universe as well as links to other databases. The data is published under the CC0 "Public domain dedication" license. It can be edited by anyone and is maintained by Wikidata's editor community.

```
{
  "@context": "http://schema.org",
  "@id": "http://datahub.io/dataset/e2427d01-3ef6-40f7-a148-f0bf28176184",
  "@type": "Dataset",
  "dateModified": "2015-11-11T10:59:55.856901",
  "datePublished": "2015-11-11T10:18:28.306264",
  "description": "The free knowledge base anyone can edit\r\nhttps://wikidata.org",
  "distribution": [
    {
      "@id": "http://datahub.io/dataset/e2427d01-3ef6-40f7-a148-f0bf28176184/resource/8fabf551-3d9e-4e23-93cb-0ef03ce45fea",
      "@type": "DataDownload",
      "contentUrl": "https://www.wikidata.org/w/api.php",
      "datePublished": "2015-11-11T10:21:42.576076",
      "description": "The MediaWiki action API is a web service that provides convenient access to wiki features, data, and meta-data over HTTP, via a URL at api.php. ",
      "fileFormat": "api/json",
      "license": "http://www.opendefinition.org/licenses/cc-zero",
      "name": "Wikidata API"
    }
  ],
}
```

```
17 .ui.card .meta {
18   font-size: 0.8em;
19 }
20 </style>
21 </head>
22 <body>
23   <header id="header" class="ui vertical segment">
24     <h2 class="ui header">
25       <div class="content">
26         Open Data Portal Watch
27         <div class="sub header">Mapping and export of Schema.org metadata descriptions for over 250 Open Data p
28       </div>
29       <a href="http://data.wu.ac.at/"></a>
30     </h2>
31   </header>
32   <nav id="mainnav" class="ui fluid secondary stackable pointing icon menu background">
33     <div class="ui item white"><a class="ui button" href="http://data.wu.ac.at/schema/">Data Portal List</a></div>
34     <div class="ui item white"><a class="ui primary button" href="http://data.wu.ac.at/odgraphsearch/">Spatio-Temp
35     <div class="ui item white"><a class="ui primary button" href="http://data.wu.ac.at/portalwatch/">Data Portal Mo
36     <!--<a class="link item"><i class="icon download"></i>Data</a>-->
37     <div class="right menu">
38       <div class="ui item white"><a href="http://data.wu.ac.at/portalwatch/about" class="ui primary button">About
39     </div>
40   </nav>
41   <!-- CONTENT -->
42   <section class="ui text">
43     <div class="ui container">
44       <div class="ui two column grid">
45         <div class="row">
46           <div class="column">
47             <div class="ui fluid card">
48               <div class="content">
49                 <div class="header">
50                   Dataset Information
51                 </div>
52                 <div class="description">
53                   <table class="ui very basic red compact table">
54                     <tbody>
55                       <tr>
56                         <td class="right aligned">Title</td>
57                         <td><a href="http://datahub.io/dataset/e2427d01-3ef6-40f7-a148-f0bf28176184">
58                       </tr>
59                       <tr>
60                         <td class="right aligned">Description</td>
61                         <td><!-- The free knowledge base anyone can edit -->
```

Dataset

All (1) ▾

Dataset		0 ERRORS 3 WARNINGS ^
ID: http://datahub.io/dataset/e2427d01-3ef6-40f7-a148-f0bf28176184		
@type	Dataset	
@id	http://datahub.io/dataset/e2427d01-3ef6-40f7-a148-f0bf28176184	
dateModified	2015-11-11T10:59:55	
datePublished	2015-11-11T10:18:28	
description	The free knowledge base anyone can edit https://wikidata.org	
keywords	Wikimedia	
keywords	wikibase	
keywords	Wikidata	
keywords	open data	
name	Wikidata	
distribution		
@type	DataDownload	
@id	http://datahub.io/dataset/e2427d01-3ef6-40f7-a148-f0bf28176184/resource/8fabf551-3d9e-4e23-93cb-0ef03ce45fea	
contentUrl	https://www.wikidata.org/w/api.php	
datePublished	2015-11-11T10:21:42	
description	The MediaWiki action API is a web service that provides convenient access to wiki features, data, and meta-data over HTTP, via a URL at api.php.	
fileFormat	api/json	
license	http://www.opendefinition.org/licenses/cc-zero	
name	Wikidata API	
encodingFormat	The encodingFormat field is recommended. Please provide a value if available.	
distribution		



Dev Wikidata

General Information

Free knowledge database project hosted by Wikimedia and edited by volunteers.

Homepage <http://wikidata.org/>

Developed in [Worldwide](#)

Created in 2012

FAIRsharing is a repository of standards and databases. It exposes resource metadata using schema.org

```
<script type="application/ld+json">
{
  "@context": "http://schema.org",
  "@id": "https://doi.org/10.25504/FAIRsharing.6s749p",
  "@type": "Dataset",
  "alternateName": "Wikidata",
  "citation": {
    "@type": "CreativeWork",
    "identifier": "doi:10.25504/FAIRsharing.6s749p"
  },
  "datePublished": "2016-09-06 16:10:14",
  "description": "Free knowledge database project hosted by Wikimedia and edited by volunteers.",
  "identifier": "10.25504/FAIRsharing.6s749p",
  "license": "https://creativecommons.org/licenses/by-sa/4.0/. Please link to https://fairsharing.org and https://fairsharing.org/static/img/home/svg/FAIRsharing-logo.svg for attribution.",
  "name": "Wikidata",
  "url": "https://doi.org/10.25504/FAIRsharing.6s749p"
}
```

# Publishing RDF - Linked Data Principles

1. Use **Uniform Resource Identifiers (URIs/URLs)** as identifiers for things
2. Use **HTTP URIs**, so that people can look up those entities
3. When someone looks up a URI, provide **useful information**, using Semantic Web standards
4. Include **links** to other URIs, so that they can discover more things

# FAIR Knowledge Graphs

## 1. Build your knowledge graph using existing standards

- a. Assign a **unique identifier (URI)** to every entity (kg, types, relations, instances)
- b. Format the data using a **data standard** (e.g. RDF, nanopublications)
- c. Capture the **provenance** and **context** of each assertion, and of the graph itself
- d. Where possible, make **links** to other published resources

## 2. Create high quality metadata to document your KG

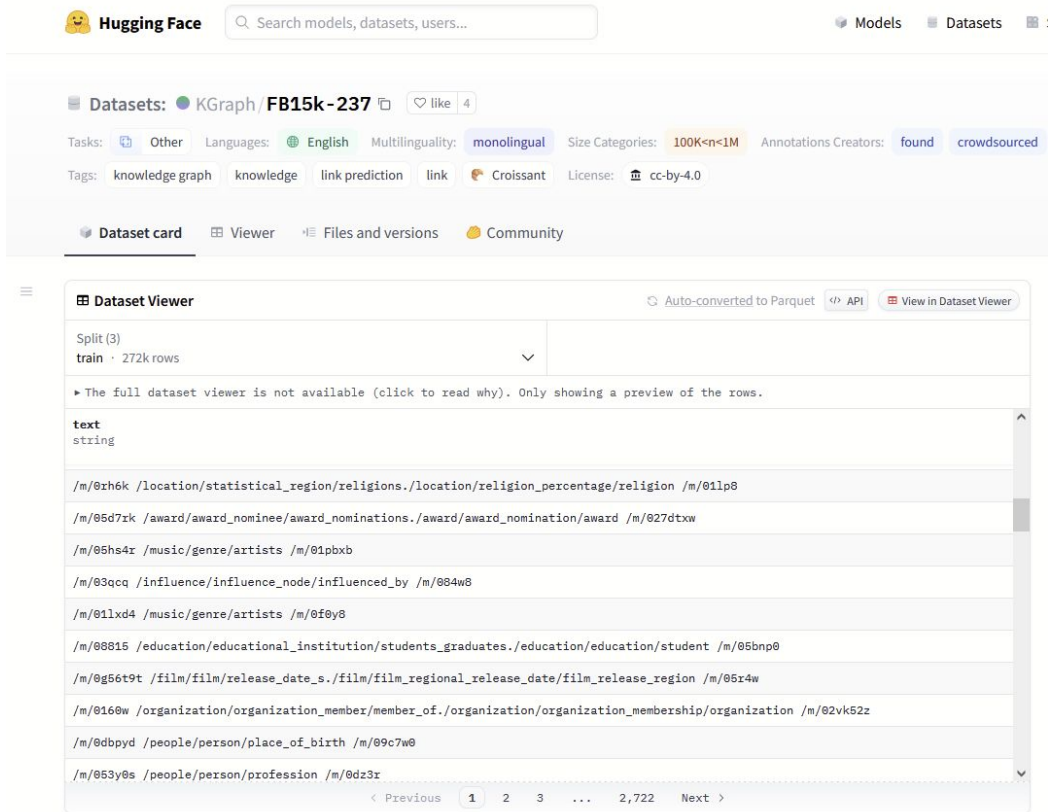
- a. **use** metadata schemas (schema.org, HCLS)
- b. **add KG specific metadata** (eg. number of triples, vocabs used, etc)
- c. **validate** the metadata against the schema

## 3. Make your KG available to others

- a. **Version** your KG, and **publish the data and metadata** to an archive
- b. Create programmatic access points via APIs

# Publish to FAIR Repositories

Hugging Face publishes structured metadata and has APIs to access the content. Also can support large files through git LFS.



The screenshot shows the Hugging Face interface for the dataset **KGGraph FB15k-237**. At the top, there's a search bar and navigation links for Models, Datasets, and a user profile. The dataset page includes a header with the dataset name, a 'like' button, and a '4' indicating the number of likes. Below this, there are filters for Tasks (Other), Languages (English), Multilinguality (monolingual), Size Categories (100K<n<1M), Annotations Creators (found, crowdsourced), Tags (knowledge graph, knowledge, link prediction, link), and License (cc-by-4.0). The main content area is titled 'Dataset card' and has tabs for Viewer, Files and versions, and Community. The 'Dataset Viewer' tab is active, showing a preview of the dataset. It indicates the dataset is split into 3 parts (train, dev, test) with 272k rows. A message states: 'The full dataset viewer is not available (click to read why). Only showing a preview of the rows.' Below this, a table shows a preview of the data with columns for 'text' and 'string'. The table contains several rows of triples, such as '/m/@zh6k /location/statistical\_region/religions./location/religion\_percentage/religion /m/@1lp8' and '/m/@5d7rk /award/award\_nominee/award\_nominations./award/award\_nomination/award /m/@27dtwx'. At the bottom of the table, there are navigation links for 'Previous', '1', '2', '3', '...', '2,722', and 'Next'.

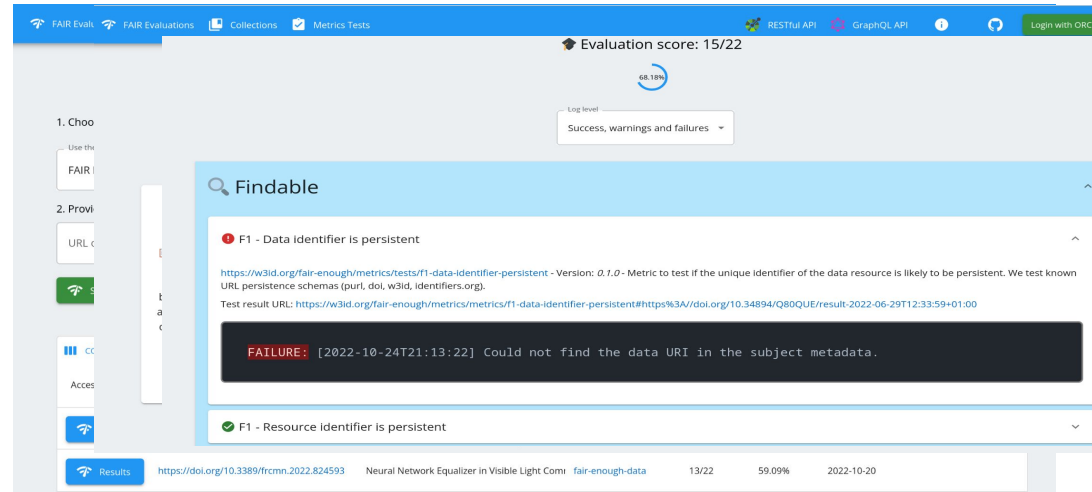
## Dataset Card for FB15k-237

### Dataset Summary

FB15k-237 is a link prediction dataset created from FB15k. While FB15k consists of 1,345 relations, 14,951 entities, and 592,213 triples, many triples are inverses that cause leakage from the training to testing and validation splits. FB15k-237 was created by Toutanova and Chen (2015) to ensure that the testing and evaluation datasets do not have inverse relation test leakage. In summary, FB15k-237 dataset contains 310,079 triples with 14,505 entities and 237 relation types.

# Quality of FAIR Implementation

- **FAIR Enough** is a system to perform automated assessment of the technical quality of the FAIRness implementation.
- FAST and can be used with different metric collections
- Keeps track of past assessments to monitor status
- Extensible via service based framework (can use FAIR Evaluator harvester and metrics)
- Open source and Docker deployable



45

<https://fair-enough.semanticscience.org>





Hugging Face

Search models, datasets, users...

Models

Datasets

Spaces

Docs

Solutions

Pricing



meta-llama / **Llama-2-70b-hf** like 369



Text Generation



PyTorch



Safetensors



Transformers



English

llama

facebook

meta

llama-2



text-generation-inference



arxiv:2307.09288



Model card



Files and versions



Train



Deploy



Use in Transformers



### Access Llama 2 on Hugging Face

This is a form to enable access to Llama 2 on Hugging Face after you have been granted access from Meta. Please visit the [Meta website](#) and accept our license terms and acceptable use policy before submitting this form. Requests will be processed in 1-2 days.

Downloads last month

168,220



Safetensors ⓘ

Model size

69B params

Tensor type

F16 · F32



<https://huggingface.co/meta-llama/Llama-2-70b-hf>

Identifier of this evaluation: <https://w3id.org/fair-enough/evaluations/2d2f90ff23f8ee901295b5786a950e4ddb0b16>

Evaluated in 3s with the [fair-enough-metadata](#) collection on the 2023-07-26

By <https://orcid.org/0000-0003-4727-9435>

See other evaluations for <https://huggingface.co/meta-llama/Llama-2-70b-hf>

#### Extracted metadata

Title: meta-llama/Llama-2-70b-hf · Hugging Face



Evaluation score: 11/16

