

# 15 principles in 20 min

Carlos Utrilla Guerrero



F A I R

# FAIR explained in 2 min



<https://www.youtube.com/watch?v=5OeCrQE3HhE&t=5s>

# Challenge:

Understand 15 principles in 20 min

My suggestions, must-read:

- 1- [The FAIR Guiding Principles for scientific data management and stewardship](#) 2016
- 2- [FAIR Principles Interpretations and Implementation Considerations](#), 2020
- 3- [Cloudy, increasingly FAIR: revising the FAIR Data Guiding Principles for the EOSC](#) 2017



Australian Research Data Commons

Importantly, it is our intent that **the principles apply not only to ‘data’ in the conventional sense, but also to the algorithms, tools, and workflows** that led to that data.

**All scholarly digital research objects**—from data to analytical pipelines—**benefit from** application of **these principles**, since all components of the research process must be available to ensure transparency, reproducibility, and reusability.

Wilkinson et al. 2016

But today, our main focus is on data.

What is "F" stand for?

F  
indable





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

Globally unique and persistent identifier for one  
resource in internet



<https://doi.org/10.34894/FQGRKC>



# Git and Dataverse repository

August 20, 2019

maastrichtlawtech/case-law-explorer v1.1

Pedro V; Kody Moodley

Scripts for automatic citation extraction from Rechtspraak.nl court decisions.

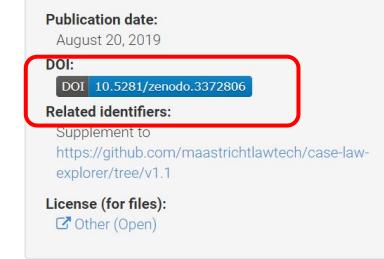
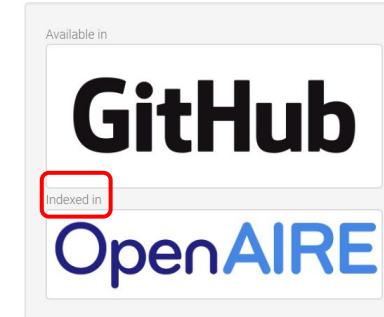
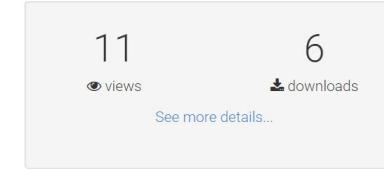
Preview

**case-law-explorer-v1.1.zip**

- maastrichtlawtech-case-law-explorer-c760a80
  - .gitignore
  - README.md
  - data\_extraction**
    - citations**
      - extract\_citations\_unix.sh
      - extract\_citations\_win.bat
      - rechtspraak\_citations\_extractor.py
      - requirements.txt
    - data\_model**
      - rechtspraakdatamodel.pdf
      - rechtspraakdatamodel.vsdx

Files (253.5 kB)	
Name	Size
maastrichtlawtech/case-law-explorer-v1.1.zip	253.5 kB
md5:2fe54512749425d4c056d26fe48cb336	

[Preview](#) [Download](#)



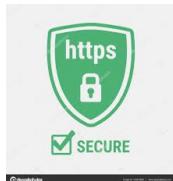
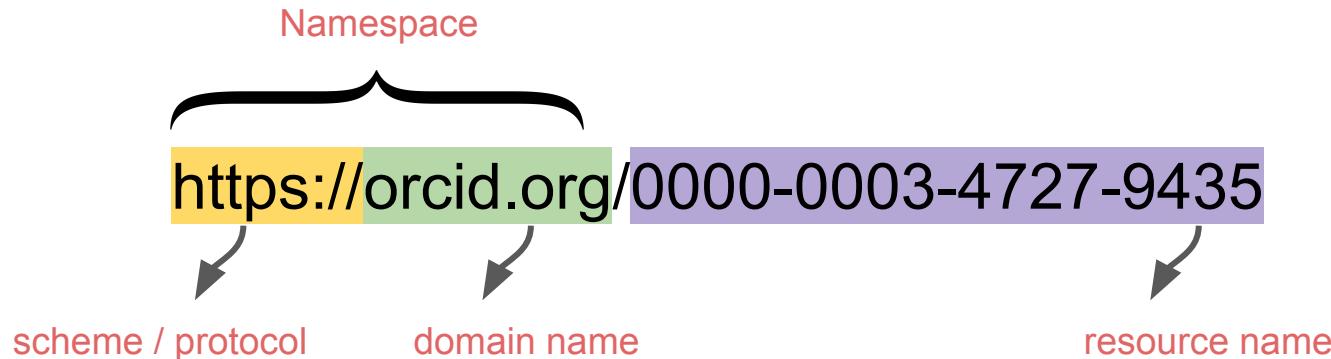
Globally unique and persistent identifier for one particular person on the planet earth

F  
indable

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



enables **unambiguous** identification of resources  
of interest.



resource name  
Michel Dumontier

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

 Print view 

Websites & Social Links  
Lab web site  
Institute of Data Science  
Maastricht University profile page  
intro-nanopub

Keywords  
data science, biomedical informatics, drug discovery, semantic web, ontology, intelligent systems, k

Other IDs  
Scopus Author ID: 6701759312



# F2: Data are described with rich metadata

# Metadata and data

Data	Metadata
	<p><b>Collection:</b> Land North and South of Blundell's Road, Tiverton</p> <p><b>Title:</b> General shot of trench 11, looking East</p> <p><b>File name:</b> General_Shot_of_Trench_11_Looking_East</p> <p><b>Created:</b> 2014-12-09</p> <p><b>Copyright:</b> Cotswold Archaeology</p> <p><b>Creator:</b> Cotswold Archaeology</p> <p><b>Coordinates:</b> OSGB 298080E 113000N</p> <p><b>Location:</b> Blundell's Road; Tiverton; Devon; England</p> <p><b>Subject:</b> Trench</p> <p><b>Period:</b> Late Neolithic; Middle Bronze Age; Roman</p> <p><b>Data Type:</b> Raster Image</p> <p><b>File Type:</b> JPEG File Interchange Format</p> <p><b>Software:</b> Canon AOS (version 6)</p> <p><b>Checksum:</b> 38e27eba5c21ded56b51b8db8dbe0adc</p>

F indable 

# Metadata and data

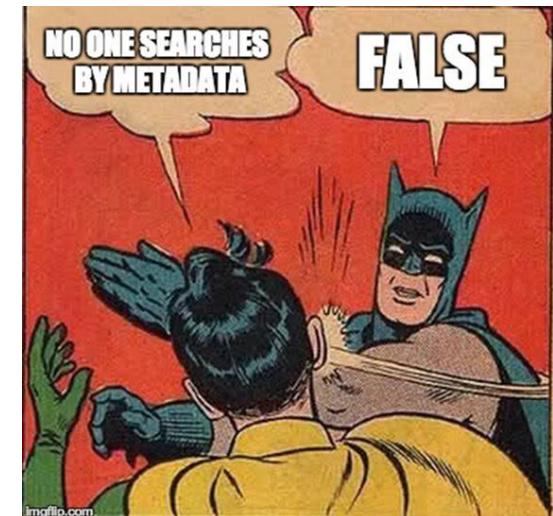
Data	Metadata
	<p><b>Collection:</b> Land North and South of Blundell's Road, Tiverton</p> <p><b>Title:</b> General shot of trench 11, looking East</p> <p><b>File name:</b> General_Shot_of_Trench_11_Looking_East</p> <p><b>Created:</b> 2014-12-09</p> <p><b>Copyright:</b> Cotswold Archaeology</p> <p><b>Creator:</b> Cotswold Archaeology</p> <p><b>Coordinates:</b> OSGB 298080E 113000N</p> <p><b>Location:</b> Blundell's Road; Tiverton; Devon; England</p> <p><b>Subject:</b> Trench</p> <p><b>Period:</b> Late Neolithic; Middle Bronze Age; Roman</p> <p><b>Data Type:</b> Raster Image</p> <p><b>File Type:</b> JPEG File Interchange Format</p> <p><b>Software:</b> Canon AOS (version 6)</p> <p><b>Checksum:</b> 38e27eba5c21ded56b51b8db8dbe0adc</p>

Findable 

Good metadata helps others to identify the your digital resources quickly via repositories for example

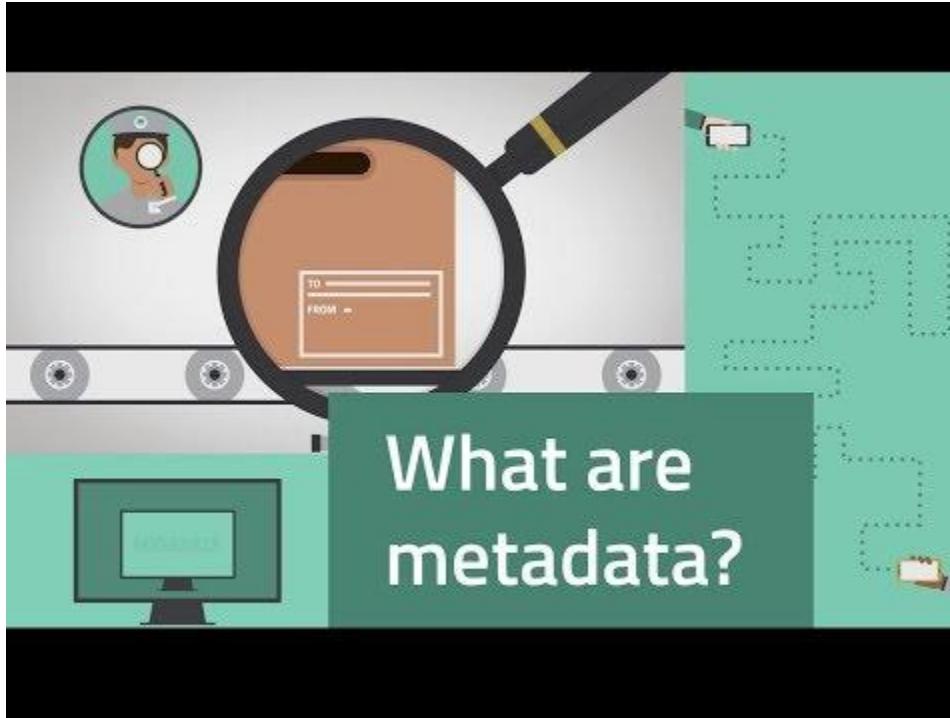
# High quality with 300 Metadata fields? Any one got time for that?

- 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?



<https://www.pinterest.com/pin/298645019032711046/>

# Why are metadata (and why are they so important?)



[https://www.youtube.com/watch?v=-4\\_MFhi4GpU](https://www.youtube.com/watch?v=-4_MFhi4GpU) (2 min)

F  
indable

# Barriers of Metadata

- Sounds like a ton of work - too time consuming
- We never learned how to share data
- Many initiatives seem to have very little buy-in from researchers
- Not very enticing to enter tons of metadata into a form in the vague hope that someday this will become easier to find for someone else
- We need some reward for the time and effort spent
- Increased utility of your data + methods

Preparing data is too time-consuming





# F3: Metadata specify the data identifier

# Because good research need good data following best practices



Citation Metadata ▾

Dataset Persistent ID ⓘ doi:10.34894/FQGRKC

Publication Date ⓘ 2020-12-08

Title ⓘ Estimated number of patients who contact the general practitioner for the first time with COVID-19-like symptoms

Alternative Title ⓘ Aantal patiënten met een eerste contact vanwege COVID-19-achtige klachten bij de huisarts

Alternative URL ⓘ <https://www.nivel.nl/nl/nivel-zorgregistraties-eerste-lijn/monitor-cijfers-covid-19-achtige-klachten-huisartsenpraktijken>

Author ⓘ M. Hooiveld (Nivel) - ORCID: 0000-0002-5513-1740

Contact ⓘ Use email button above to contact.  
Directie (Nivel)

Description ⓘ A representative sample of around 350 GP practices across the Netherlands provide data once a week to the Nivel (Netherlands Institute for Health Services Research) Primary Care Database. They record reported and observed symptoms and diagnoses of the consulting patients they see. Using this data, Nivel determines how many patients saw their GP in the past week for the first time for symptoms that could indicate COVID-19. Nivel determines this on the basis of the diagnostic codes reported by GPs ('acute upper respiratory system infection', 'other respiratory infection(s)', 'influenza', 'pneumonia', 'other viral disease(s)', 'other infectious disease', 'fever', 'shortness of breath', 'coughing') in combination with textual information provided that could indicate COVID-19. To improve this indicator's accuracy, Nivel also recalculates the data for the preceding weeks, so they include information that only became available at a later time. The Nivel Primary Care Database provides this information once a week (Thursdays) and delivers this information in JSON format to the Ministry of Health, Welfare and Sport. The source file gives the numbers per week (Monday through Sunday).

Subject ⓘ Medicine, Health and Life Sciences

Related Publication ⓘ Hooiveld M, Heik K, Heins M, Hendriksen J, Bolt E, Weesie Y, Spreeuwenberg P, Korevaar J. Cijfers COVID-19-achtige klachten in huisartsenpraktijken. Nivel Zorgregistraties Eerste Lijn, Utrecht: Nivel, 2020. url: <https://www.nivel.nl/sites/default/files/bestanden/1003896.pdf>

Notes ⓘ This dataset is updated every week.

Depositor ⓘ Hooiveld, Mariette

Deposit Date ⓘ 2020-12-07

Date of Collection ⓘ Start: 2020-03-02

Kind of Data ⓘ Clinical aggregated data

Software ⓘ Stata, Version: 15.1



## Mandatory Properties

Mandatory and Recommended properties and sub-properties are especially valuable for information seekers and added-service providers, such as indexers. The DataCite I Working Group members strongly urge the inclusion of metadata identified as Required for the purpose of achieving greater exposure for the resource's metadata record, an underlying research itself.

Mandatory Properties
Identifier
Creator
Title
Publisher
PublicationYear
ResourceType



## Recommended Properties

Recommended Properties
Subject
Contributor
Date
RelatedIdentifier
Description
GeoLocation

## Optional Properties

Optional Properties
Language
AlternateIdentifier
Size
Format
Version
Rights
FundingReference



# F4: (Meta)data are registered or indexed in a searchable resource

**F3: Metadata clearly and explicitly include the identifier of the data they describe**  
***(basically it states data and metadata should be in different files)***

Supports manual exploration - twitter vaccination dataset example



Google



## Dataset Search

Search for Data Sets



Try [coronavirus covid-19](#) or [education outcomes site:data.gov](#).

[Learn more](#) about Dataset Search.

Think: How does your discipline share data?



*Illustration credit: Ainsley Seago.*

<https://www.nature.com/sdata/policies/repositories>

# Kept safe in a secure environment



Following a year-long public discussion and building on existing community consensus<sup>1</sup>, several stakeholders, representing various segments of the digital repository community, have collaboratively developed and endorsed a set of guiding principles to demonstrate digital repository trustworthiness. Transparency, Responsibility, User focus, Sustainability and Technology: the TRUST Principles provide a common framework to facilitate discussion and implementation of best practice in digital preservation by all stakeholders.

<https://www.nature.com/articles/s41597-020-0486-7>



What is "A" stand for?

A  
ccessible



Here the focus is on the '*ability*' and '*conditions about compliance*'



**A1:** (meta) data are retrievable by their identifier using a standard communications protocols.

**A1.1:** the protocol is open, free, and universally implementable.

**A1.2:** the protocol allows for an authentication and authorisation procedure, where necessary.

*'Ability'* = there should be no additional barriers to retrieve data

A.I.I. Several open, free and universally implementable protocols framework but HTTP protocol widely supported by community



HTTP - Hypertext Transfer Protocol

Accessible



Swagger. Supported by SMARTBEAR

Why Swagger? Tools Resources

Search Sign In Try Free

SwaggerHub  
Swagger Inspector  
Open Source Tools  
Specification  
What Is OpenAPI?  
Basic Structure

```
1. securityDefinitions:  
2.   BasicAuth:  
3.     type: basic  
4.   ApikeyAuth:  
5.     type: apiKey  
6.     in: header  
7.     name: X-API-Key  
8. OAuth2:  
9.   type: oauth2  
10.  flow: accessCode  
11.  authorizationUrl: https://example.com/oauth/authorize  
12.  tokenUrl: https://example.com/oauth/token  
13.  scopes:  
14.    read: Grants read access  
15.    write: Grants write access  
16.    admin: Grants read and write access to administrative information
```

Each security scheme can be of type:

- **basic** for Basic authentication
- **apiKey** for an API key
- **oauth2** for OAuth 2

# Why do we care about machine accessibility?

A  
ccessible 

- Teach computers how communicate and quickly access to data related to the pandemic - EU Open source solutions to help medical staff, public administrations, citizens in daily activities.



PUBLIC HEALTH

Mapping COVID-19

By Lauren Gardner, January 23, 2020

## Covid19 GraphQL API

<https://covid19-graphql.now.sh>

 Deploy

Data is pulled directly from <https://github.com/pomber/covid19>, which is a JSON representation of <https://github.com/CSSEGISandData/COVID-19>. All data is up to date.

Example query

```
query {  
  # time series data  
  results (countries: ["US", "Canada"], date: { lt: "3/10/2020" }) {  
    country {  
      name  
      date  
      confirmed  
      deaths  
      recovered  
      growthRate  
    }  
    # by country  
    country(name: "US") {  
      name  
      mostRecent {  
        date(format: "yyyy-MM-dd")  
        confirmed  
      }  
    }  
  }  
}
```

# Why do we care about machine accessibility?

A  
ccessible 



ELSEVIER

## What's an API? 5 things you need to know to stay current

APIs are having a growing impact on our lives and work – even if you're not a techie

By Alice Atkinson-Bonasio Posted on 24 September 2014

*"Our users are researchers, and generally they spend a lot of time finding the right content, and then organising and formatting their data. We want to support our dev community so they use the Mendeley API to integrate our product not only with Elsevier's data and offerings, but those of other publishers – and any product out there that suits their needs and makes their lives easier."*

# Why do we care about machine accessibility?

A  
ccessible 



ELSEVIER

Interoperability is about ensuring that applications, tools and data sets from different providers can work together. The Mendeley API represents our commitment to interoperability with any tools that researchers need.

# Why do APIs matter to me?

A  
ccessible 

Academic research

## Take your research further with Twitter data

From social science to computer science, Twitter data can advance research objectives on topics as diverse as the global conversations happening on Twitter.

# When to create an API?



- Collaboration and sharing your model with data scientist and software developers.
- Your data set is large, making download via FTP complicated.
- Your research community only need access to a part of the data at any one time.
- Your data changes or is updated frequently.
- Capture and share time-series data.
- Unstructured data to convert into structured format (e.g. *Linguistics*).
- [Nice article for historian programming: why to use API as a researcher?](#)  
[And Europeana developer for cultural heritage.](#)



# *Why open, free, universally implementable?*

A  
ccessible 



AI.2 for FAIR != "open"



A  
ccessible

Most of the **ESA Earth observation datasets** are available on the Internet free of charge.  
Access is granted after user registration, which also provides the detailed content of the  
free datasets.



# Sensitive Data Community of Practice



A  
ccessible

- Government Data
- Health Data
- Indigenous Data
- Commercial Data



# Defense-security high protected data, but still FAIR



Accessible 

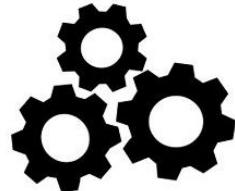
API OPTIONS	RF
Request-driven API	RESTful HTTPS (JSON) ✓
	RESTful HTTPS (XML) ✓
Event-driven API	HTTPS Websocket (JSON) ✓
Export Formats	JSON, XML, SHP, KML ✓



API = Application Programming Interface

What is "I" stand for?

I  
nteroperable

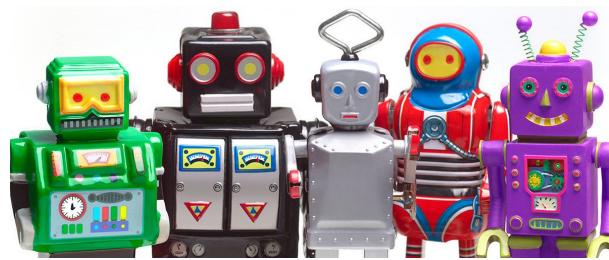




# I1: (Meta)data use formal, accessible, shared, and broadly applicable language for knowledge representation

# Here we start empathizing with machines, and before!!

**Be inclusive with others “common understanding”**



**inclusive** = empower these agents by providing an environment that is understandable for them and easily make use of our data

# Which is the primary goal of FAIR?

**Step towards a common understanding of the data by machines, which is a prerequisite for a functional internet of FAIR Data - Interpretations 2020**



Suppose a machine visits two datasets in which field "temperature" is present



Date	Temp
01/01/2021	191
02/01/2021	195
03/01/2021	200
04/01/2021	199
05/01/2021	180



Machines and also humans need **contextual** information to determine things rather than assume

Date ( <u>DD/MM/YY</u> )	Temp ( <u>Kelvin</u> )
01/01/2021	191
02/01/2021	195
03/01/2021	200
04/01/2021	199
05/01/2021	180

```
{
  "coord": {
    "lon": -122.08,
    "lat": 37.39
  },
  "weather": [
    {
      "id": 800,
      "main": "Clear",
      "description": "clear sky",
      "icon": "01d"
    }
  ],
  "base": "stations",
  "main": {
    "temp": 282.55,
    "feels_like": 281.86,
    "temp_min": 280.37,
    "temp_max": 284.26,
    "pressure": 1023,
    "humidity": 100
  },
  "visibility": 16093,
  "wind": {
    "speed": 1.5,
    "deg": 350
  },
  "clouds": {
    "all": 1
  },
  "dt": 1560350645,
  "sys": {
    "type": 1,
    "id": 5122,
    "message": 0.0139,
    "country": "US",
    "sunrise": 1560343627,
    "sunset": 1560396563
  },
  "timezone": -25200,
  "id": 420006353,
  "name": "Mountain View",
  "cod": 200
}
```

## GeoJSON



Even better for  
machine-readable i.e.  
value-key pair"

Interoperable



### Fields in API response

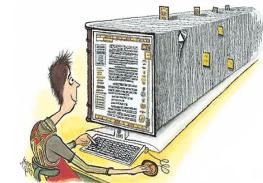
- coord
  - coord.lon City geo location, longitude
  - coord.lat City geo location, latitude
- weather (more info Weather condition codes)
  - weather.id Weather condition id
  - weather.main Group of weather parameters (Rain, Snow, Extreme etc.)
  - weather.description Weather condition within the group. You can get the output in your language. [Learn more](#)
  - weather.icon Weather icon id
- base Internal parameter
- main
  - main.temp Temperature. Unit Default: Kelvin, Metric: Celsius, Imperial: Fahrenheit.

FAIR speaks to the ability of data to be reused by a generic agent, rather than a community-specific agent.



*“The most widely-accepted choice is the Resource Description Framework (RDF) which is the W3C’s recommendation for how to represent knowledge on the Web in a machine-accessible format” .  
(FAIR principles Interpretation, 2021)*

[https://www.mitpressjournals.org/doi/full/10.1162/dint\\_r\\_00024](https://www.mitpressjournals.org/doi/full/10.1162/dint_r_00024)





# I2: (Meta)data use vocabularies that follow the FAIR principles

Terminology you use, for instance, the units of measure, classifications, and relationship definitions are themselves FAIR



Insufficient information to enable a machine understand the meaning and context



Date	Thing	Number
01/01/2020	Kiwi	1
02/01/2020	Kiwi	2
03/01/2020	Kiwi	3
04/01/2020	Kiwi	2
05/01/2020	Kiwi	1

# Shared vocabularies



What do you mean by “Kiwi”?



<https://en.wikipedia.org/wiki/Kiwifruit>



[https://en.wikipedia.org/wiki/Kiwi\\_\(bird\)](https://en.wikipedia.org/wiki/Kiwi_(bird))

Choose the one most suitable for your dataset



Date (DD/MM/YYYY)	Species	Number
01/01/2020	<a href="#"><u>Kiwi</u></a>	1
02/01/2020	<a href="#"><u>Kiwi</u></a>	2
03/01/2020	<a href="#"><u>Kiwi</u></a>	3
04/01/2020	<a href="#"><u>Kiwi</u></a>	2
05/01/2020	<a href="#"><u>Kiwi</u></a>	1



# I3: (Meta)data include qualified references to other (meta)data

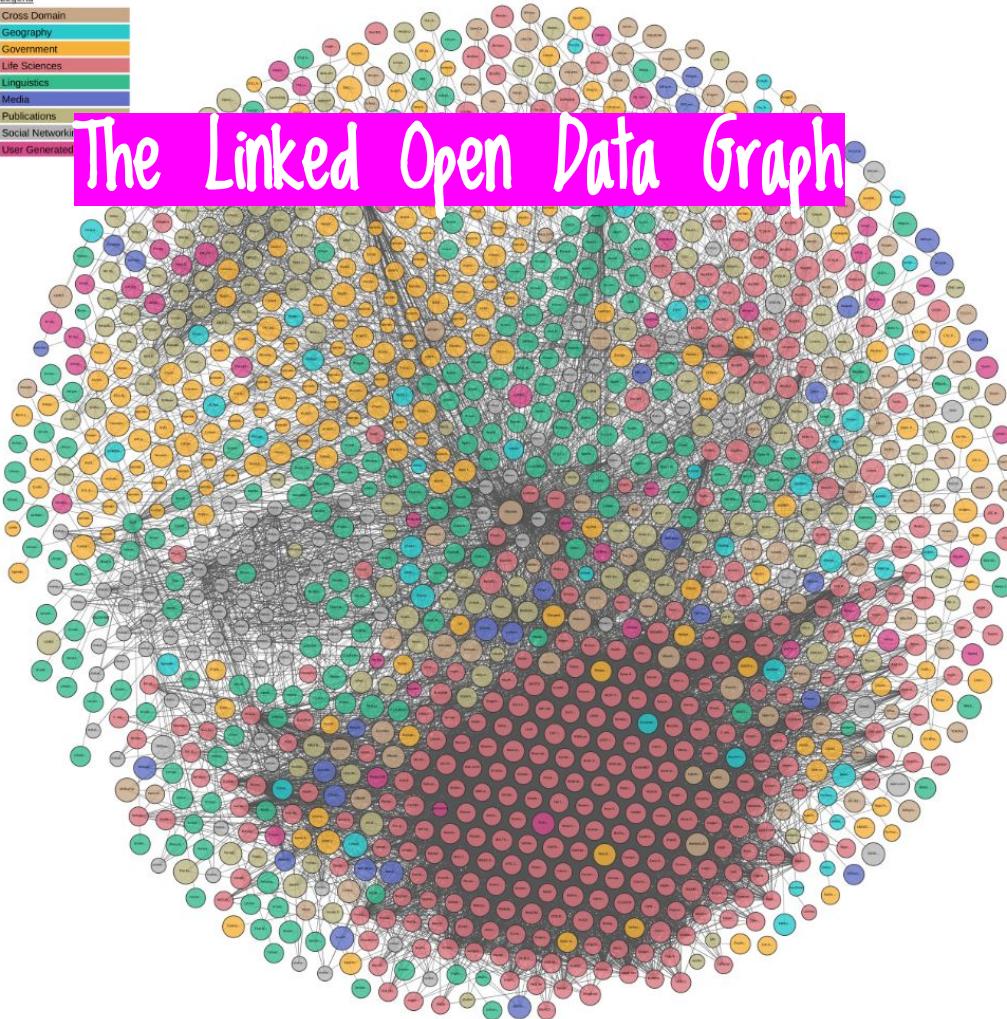


## Data silos – not made for sharing

"We must do what is necessary to ensure that the **knowledge** representing a resource is **connected** to that of other resources to create a meaningfully **interlinked** network of data and services "- Interpretation FAIR 2020

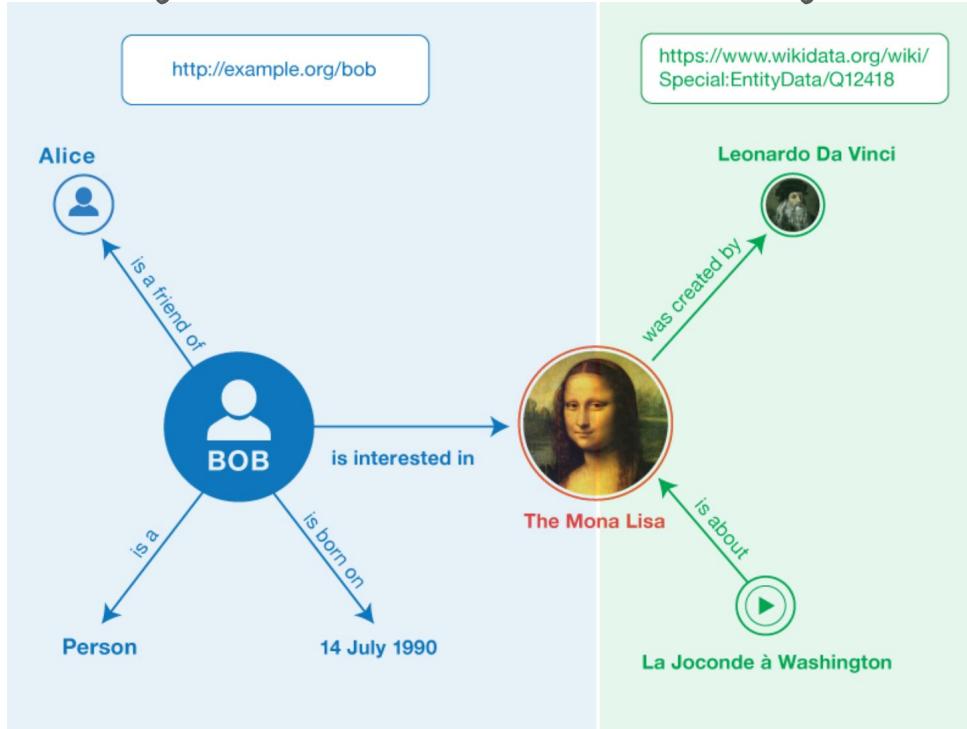
Legend
Cross Domain
Geography
Government
Life Sciences
Linguistics
Media
Publications
Social Networks
User Generated

# The Linked Open Data Graph



1,224 datasets  
with 16,113 links  
(as of June 2018)

# Subjects, Predicates and Objects



**RDF Data Model:** A set of **triples** that consists in **three components**:



The triple expresses a relationship between the subject and the object. The subject and object represent nodes in the RDF graph and the predicate represents an edge / relation between these nodes

# Knowledge Graph



Interoperable 

What is "R" stand for?

R<sub>Reusable</sub> 

R1.1. (Meta)data are richly described with a plurality of accurate and relevant attributes

R<sub>Reusable</sub> 

# wait..It's not similar to Findable Principle (F2)?



We chose the term '*plurality*' to indicate that the metadata author should be as generous as possible in providing metadata, even including information that may seem irrelevant. - GO FAIR

**F2:** Repository and indexed metadata.

**R1:** Is your data valid for my study? Give me instructions (i.e. readme) to assess if your data is appropriate for me.



# R1.1. (Meta)data are related with clear and accessible data usage license



## Licenses & Standards

### About Open Source Licenses

Open source licenses are licenses that comply with the [Open Source Definition](#) — in brief, they allow software to be freely used, modified, and shared. To be approved by the Open Source Initiative (also known as the OSI), a license must go through the [Open Source Initiative's license review process](#).

### Popular Licenses

The following OSI-approved licenses are popular, widely used, or have strong communities:

- Apache License 2.0
- BSD 3-Clause "New" or "Revised" license
- BSD 2-Clause "Simplified" or "FreeBSD" license
- GNU General Public License (GPL)
- GNU Library or "Lesser" General Public License (LGPL)
- MIT license
- Mozilla Public License 2.0
- Common Development and Distribution License
- Eclipse Public License version 2.0

<https://opensource.org/licenses>




**creative commons**

**Choose an open source license**

An open source license protects contributors and users. Businesses and savvy developers won't touch a project without this protection.

Which of the following best describes your situation?



I need to work in a community.

Use the license preferred by the community you're contributing to or dependent on. It's the right thing to do.

If you have a dependency that doesn't have a license, ask its maintainers to add a license.



I want it simple and permissive.

The [MIT License](#) is short and to the point. It lets people do almost anything they want with your code, including changing and distributing closed source versions.

[Ansible](#), [Bash](#), and [GIMP](#) use the MIT License.



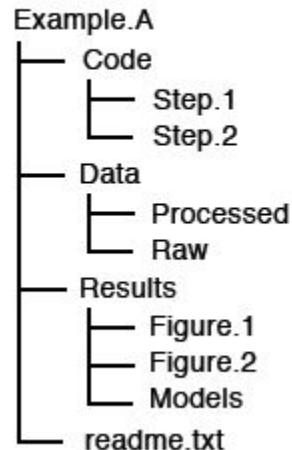
I care about sharing improvements.

The [GNU GPLv3](#) also lets people do almost anything they want with your project, except distributing closed source versions.

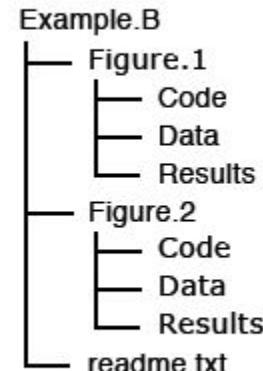
# R1.2. (Meta)data are associated with detailed provenance



**A) Organized by File type**



**B) Organized by Analysis**



# Facilitate reproducibility by reproducible research things with community standards



There's a README.md as part of the [git repository](#) that can give you some background. A Institute generic version of this can be found in a branch on the github.

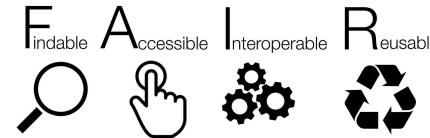
## Lessons

---

- [Lesson 1 - Documentation](#)
- [Lesson 2 - Naming Conventions](#)
- [Lesson 3 - Folder Structure](#)
- [Lesson 4 - Automation](#)
- [Lesson 5 - Version control](#)
- [Lesson 6 - Cloud Backups](#)
- [Lesson 7 - Computer Security](#)
- [Lesson 8 - Separating Identified Variables](#)
- [Lesson 9 - Permanent Identifiers for your Published Results](#)



# Summary FAIR principles



**Findable**: Globally unique, persistent identifier (e.g., DOI), rich metadata, indexed for search

**Accessible**: Retrievable using standardised, open protocol (e.g., HTTPS). Metadata stay accessible when data is removed.

**Interoperable**: Metadata use a formal, accessible, shared, broadly applicable language/vocabulary, references to other metadata

**Re-usable**: Accurate and relevant attributes, provenance and data usage licence is clear, meet domain-relevant community standards

# Why go FAIR? [write 3 pros]

- To be able to use your data it should be **findable** and **accessible**
- Make it easier for me to **use my data for a new purpose**
- For effective **reuse** since all relevant aspects are reported and it has a detailed license and provenance
- Make it easier/possible for people to **verify my work**
- Ensure the data are available in the future, especially after I stopped caring about them
- Making data interoperable allows for better analysis
- Satisfy expectations around data management from institutions

# Why practice FAIR objects? [write 3 pros]

## The idealist

- **Shoulders of giants!**
- Validates scientific knowledge
- Allows others to build on your findings
- Improved transparency
- Increased transfer of knowledge
- Increased utility of your data + methods

## The pragmatist

- Increased efficiency
- Reduces false leads based on irreproducible findings
- Data sharing citation advantage (Piwowar 2013)
- “It takes some effort to organize your research to be reproducible... the **principal beneficiary is generally the author herself.**” - Schwab & Claerbout

Poll

Go to slido.com

Enter code: FAIR-UM

# FAIR for your dataset

- 1. Build and collect your dataset using existing standards:**
  - a. Assign a meaningful variable name to every column
  - b. Format the data using your community standards (e.g. csv, json, rdf)
  - c. Capture the provenance and context of your dataset
- 2. Create high quality metadata to document your dataset:**
  - a. Use vocabularies and standards
  - b. Use specific metadata terms where necessary
- 3. Make your dataset available to others**
  - a. Version your dataset, and publish the data and metadata e.g. DataverseNL
  - b. Create access points via APIs if necessary