# Introduction to Ontology Semantics & tools for data annotation

Harold Duruflé

# Introduction to Ontology
# Semantics & tools for data annotation

Harold Duruflé

# Problem: Several words for retrieving the same data


(SUNRISE project)

**Sunflower**

*Helianthus spp.*

*Helianthus annuus* **L.**
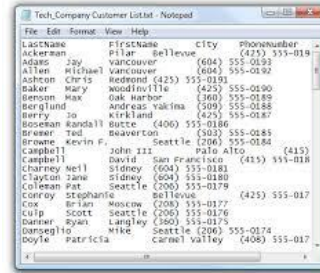
**Tournesol**

**Helianthus**

→ **Semantic heterogeneity (concepts)**

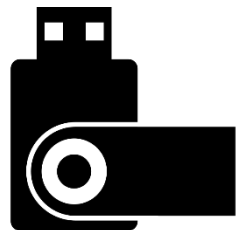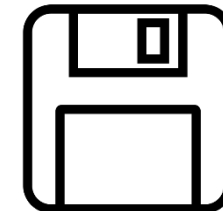# Problem: Data is multi-format and everywhere

Databases

Tabulated files

Publication / Reports

text files

Excel files

**→ Structural and syntactic heterogeneity**

# Several domains



SEMANTICS - THE WAY TO RECONCILE POINTS OF VIEW AND DATA
THE EXAMPLE OF "RICE"

→ **Domain heterogeneity**

4

# Different types of semantic resources



Stronger Semantics (Abstraction)

(e.g. Relations, concepts)

Ontologies

Taxonomies

Web Ontology Language (OWL)

Controlled vocabularies

Concept Maps

RDF

Thesaurus

UML

Glossaries

XML

Word/HTML

Weaker Semantics

Time

(Adapted to Semantic Spectrum & Network Inference)

# Example:



Level of Abstraction | Problem Solving | Yield Improvement
| Conversation | Pathways
| Sentences | Structural motif
| Concepts | Protein
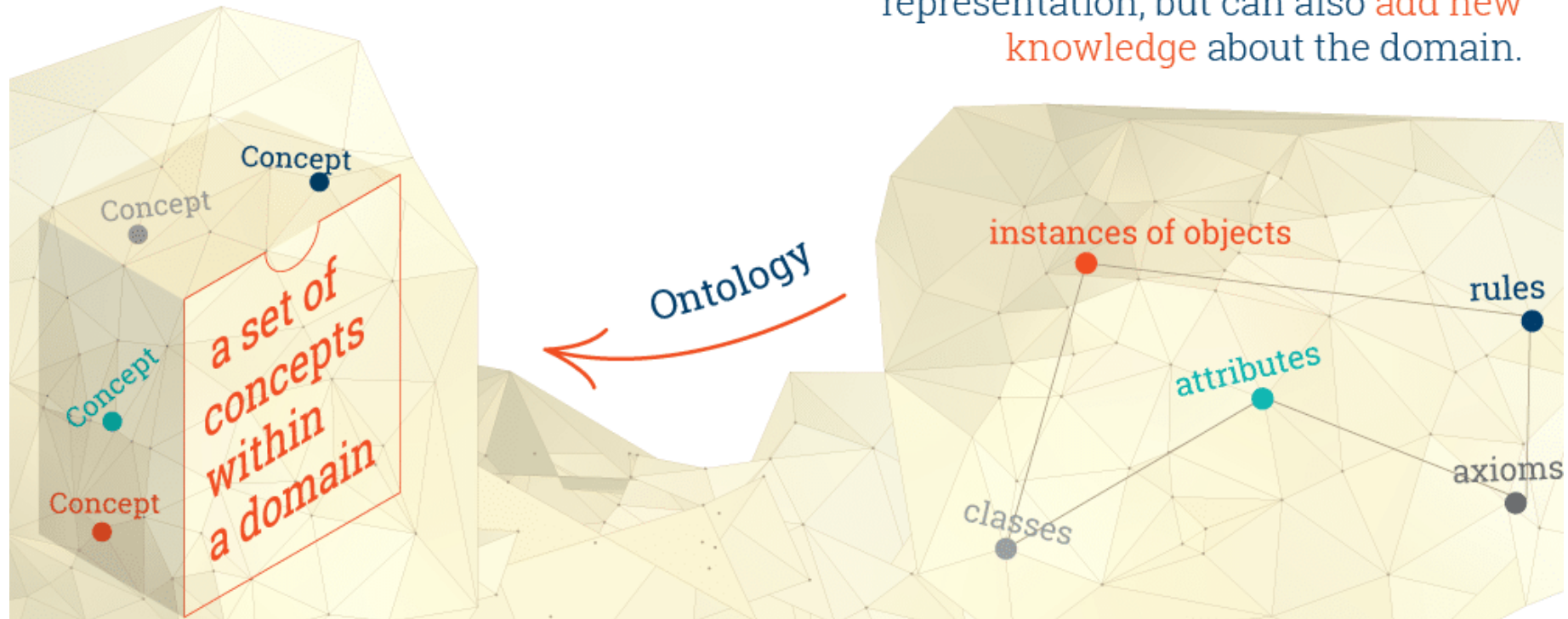| Words | Sequence
| Sound | DNA

# What is an ontology?

- Provide a shared vocabulary for a domain (all the terms)
- Provide textual definitions that describe the intended meaning of the terms in vocabularies
- Provide standard identifiers for concepts describing a given domain
- Provide machine-readable axioms and definitions that enable computational access to some aspects of the meaning of classes and relations – logical representation of human knowledge

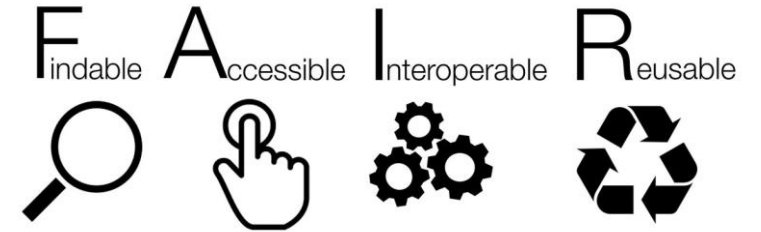→ **Facilitate data publication / data access and analysis**

# What is an ontology?



Ontologies do not only introduce a sharable and reusable knowledge representation, but can also add new knowledge about the domain.

Concept

Concept

Ontology

a set of concepts within a domain

Concept

Concept

instances of objects

rules

attributes

classes

axioms

https://www.ontotext.com/knowledgehub/fundamentals/what-are-ontologies/

# FAIR principles



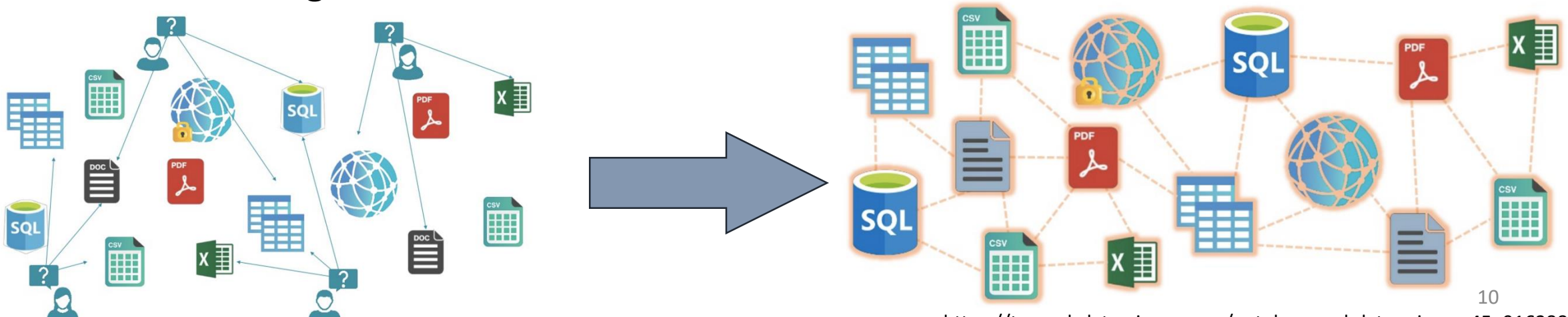FAIR: Findable, Accessible, Interoperable, Reusable

(Wilikinson, 2016 nature DOI: DOI:10.1038/sdata.2016.18)

**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.
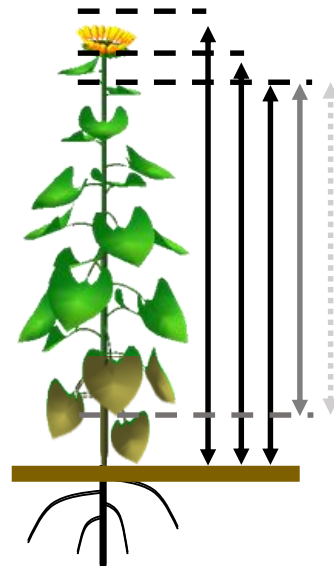
# Why biologist have adopted ontologies?

- To provide canonical representation of scientific knowledge

- To annotate experimental data to enable interpretation, comparison, and discovery across databases (Example: GO 👁 )

- To facilitate knowledge-based applications for
  - Decision support
  - Natural language-processing
  - Data integration

https://towardsdatascience.com/ontology-and-data-science-45e916288cc5

# Example of the crop ontology

- No naming convention for variables and methods of measurement which are heterogeneous

- Trait & Variable definitions and measurement are not similar between farmers, breeders, agronomists, modelers,…

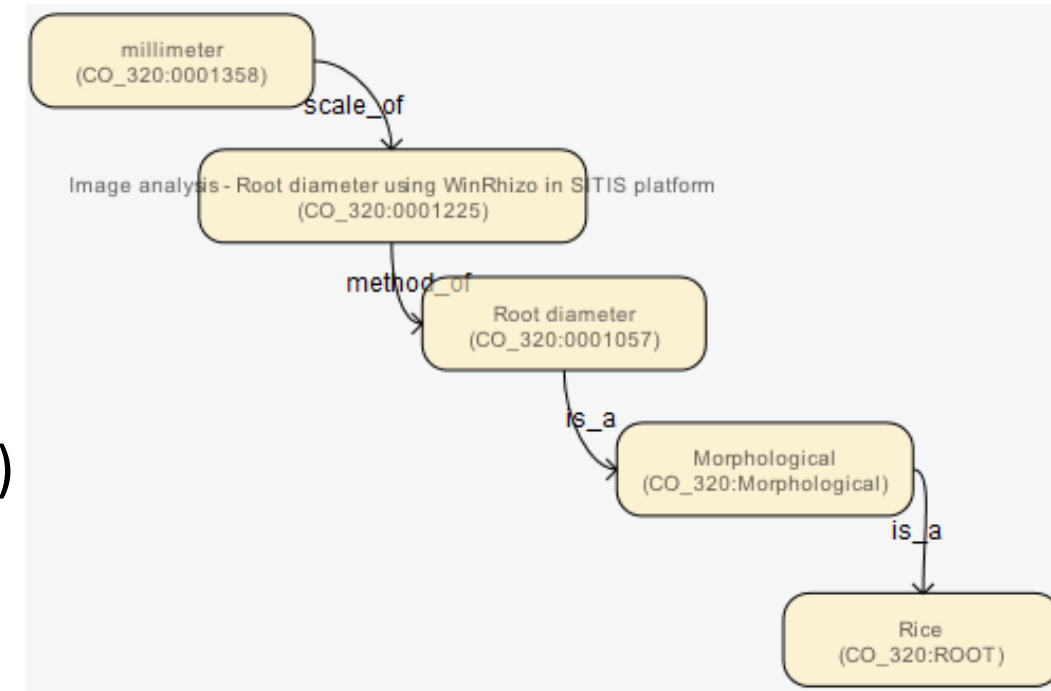Plant Height

→ One trait = x traits…

# Example of the crop ontology

Annotation must explain:

1/ What is the observation about? = **TRAIT**
     (e.g. Plant Height, Color of grain)

2/ How is the trait observed? = **METHOD**
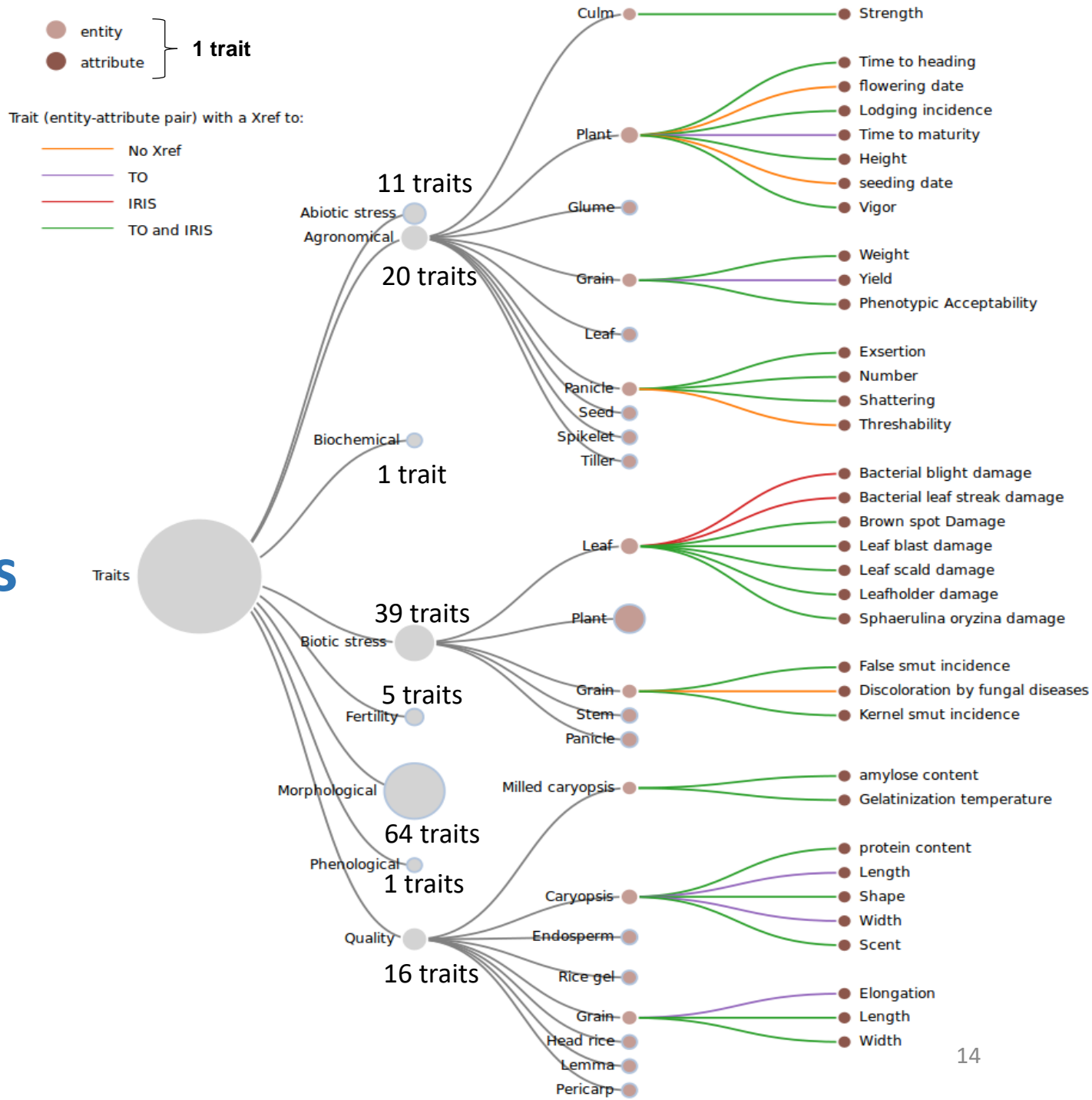     (e.g. Measuring, Estimated visually, Calculated)

3/ How is the trait observation expressed? = **SCALE**
     (e.g. cm, short/medium, white/black)

| Identifier | CO_320:0001057 |
|---|---|
| Trait description | Diameter of a cross-section of the root |
| Attribute | Diameter |

| Identifier | CO_320:0001225 |
|---|---|
| Method class | Measurement |
| Method description | Scan roots from soil depths of 0-45 cm using CI-600 scanner system associated with the WinRhizo software and calculate root diameter based on image Winrhyzo analysis. WinRHIZO uses a non-statistical method for measuring root morphology. It calculates total root length from a one pixel thinned image by multiplying the number of pixels by pixel size, and calculates average diameter by dividing the projected area of the imaged object by the total length. |
| Method name | Image analysis - Root diameter using WinRhizo in SITIS platform |

| Identifier | CO_320:0001358 |
|---|---|
| Scale Xref | UO:0000016 |
| Scale class | Numerical |
| Scale name | millimeter |

# Example of the crop ontology

Annotation must explain:

1/ What is the observation about? = **TRAIT**
        (e.g. Plant Height, Color of grain)



2/ How is the trait observed? = **METHOD**
        (e.g. Measuring, Estimated visually, Calculated)

3/ How is the trait observation expressed? = **SCALE**
        (e.g. cm, short/medium, white/black)

# Example of an ontology



Rice Traits

# Example of link between ontologies

**fruit color trait** (TO:0002617) across various species



Fruit Color

Kernel Color

Grain Color

Achene Color

Berry Color

Pod Color

→ **Ontologies can link to data from multiple species**

# Crop Ontology Workflow



**Breeding Management System**
Breeders, Data Manager & Scientists

Plants' traits ontology

Crop Ontology
for agricultural data

Defined and organized variables
In Trait Dictionary

Fieldbook creation

Ease data collection in the field thanks to defined variables

Published online
Possibility to download variables
Possibility to upload variables

Database
Store annotated data, and allow their interpretation and harmonization

**Data Analysis**

# Thank you

- Bioversity international CGIAR
  - **Elizabeth Arnaud**
  - **Marie Angélique Laporte**

  - **Cyril Pommier**

## HOW STANDARDS PROLIFERATE:
(SEE: A/C CHARGERS, CHARACTER ENCODINGS, INSTANT MESSAGING, ETC)

SITUATION: THERE ARE 14 COMPETING STANDARDS.

14?! RIDICULOUS! WE NEED TO DEVELOP ONE UNIVERSAL STANDARD THAT COVERS EVERYONE'S USE CASES.
YEAH!

SOON:
SITUATION: THERE ARE 15 COMPETING STANDARDS.

(xkcd.com/927)

17

## 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

(Wilikinson, 2016 nature DOI: DOI:10.1038/sdata.2016.18)