



Semantic Data Science

Markus Stocker

TIB Leibniz Information Centre for Science and Technology

@envinf

<https://orcid.org/0000-0001-5492-3212>



Objective

Make the case that **scientists** and research infrastructures **are elements of scientific knowledge infrastructures** and that there is an urgent need in such infrastructures to **ensure that the meaning of data is made explicit**.



Schedule

Time	Content
08:00 - 09:00	Part I: The conceptual and technology framework
09:00 - 09:30	Part II: A knowledge infrastructure in atmospheric physics as a case in point
09:30 - 10:15	Part III: Hands-on assignment
10:15 - 10:30	Part IV: Discussion

Part I

The conceptual and technology framework

—

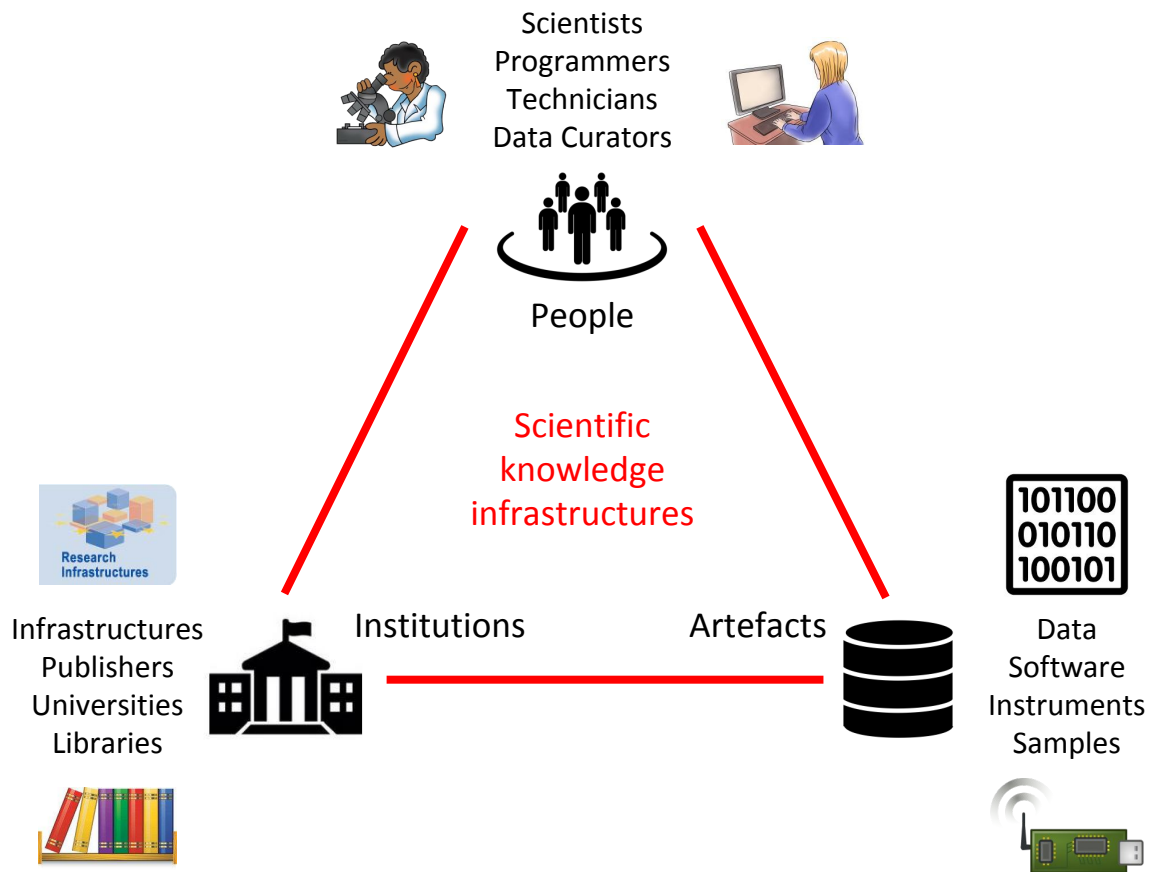
The conceptual framework



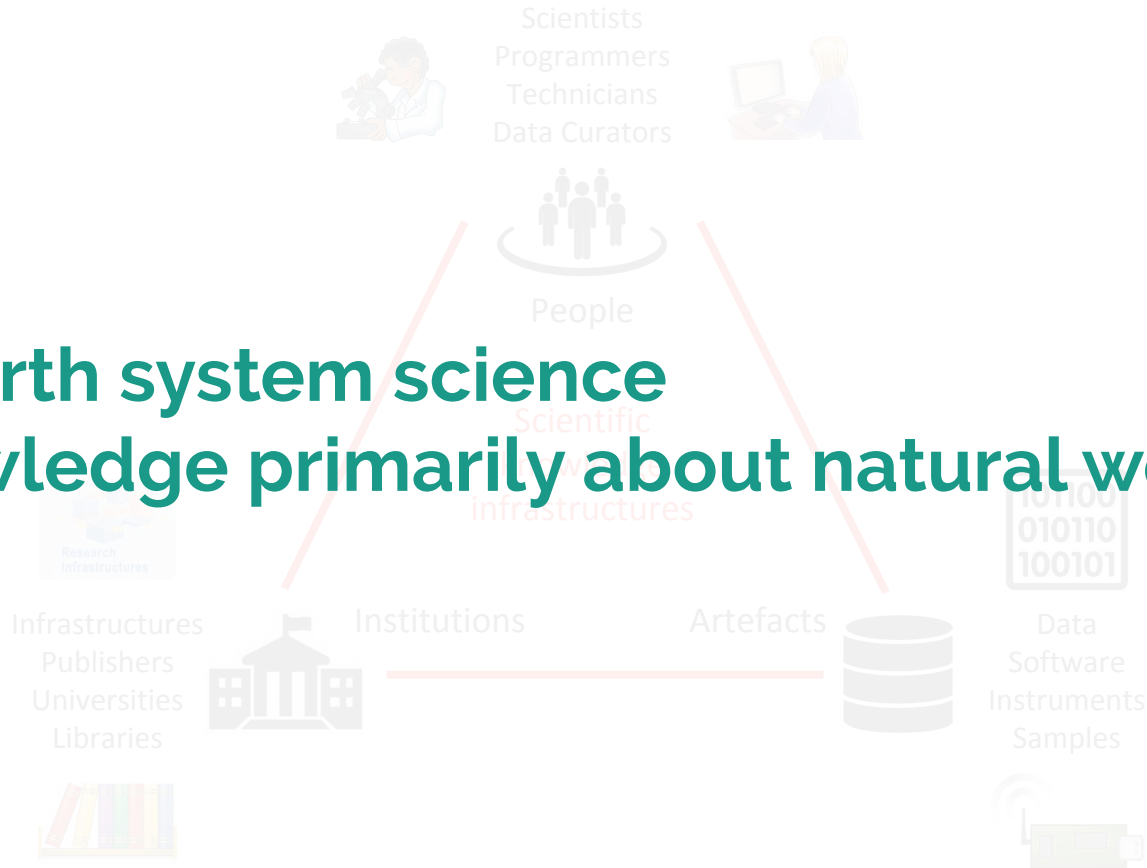
Knowledge infrastructures

Robust networks of people, artifacts, and institutions that generate, share, and maintain specific knowledge about the human and natural worlds.

-- Paul Edwards (2010). *A Vast Machine: Computer Models, Climate Data, and the Politics of Global Warming*, p. 17

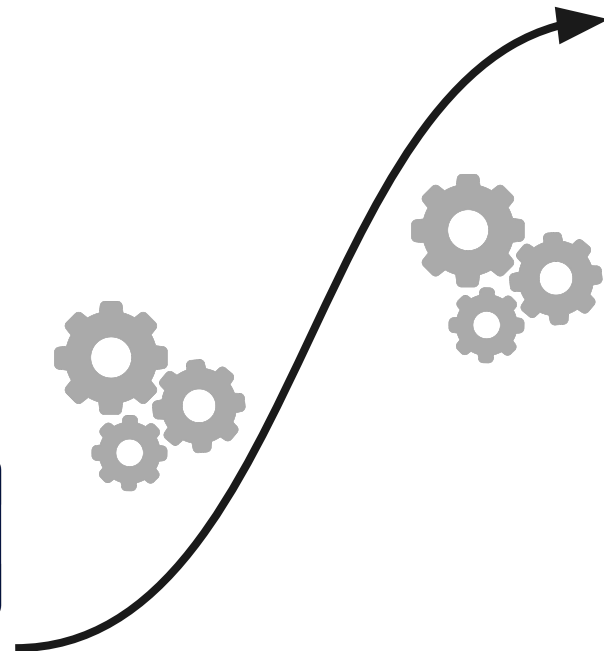


In Earth system science knowledge primarily about natural worlds

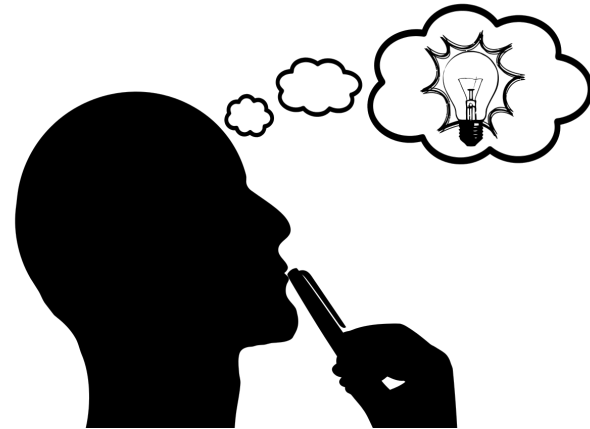


Data

111
001011
01100110
10011011011
01101100
101101
111

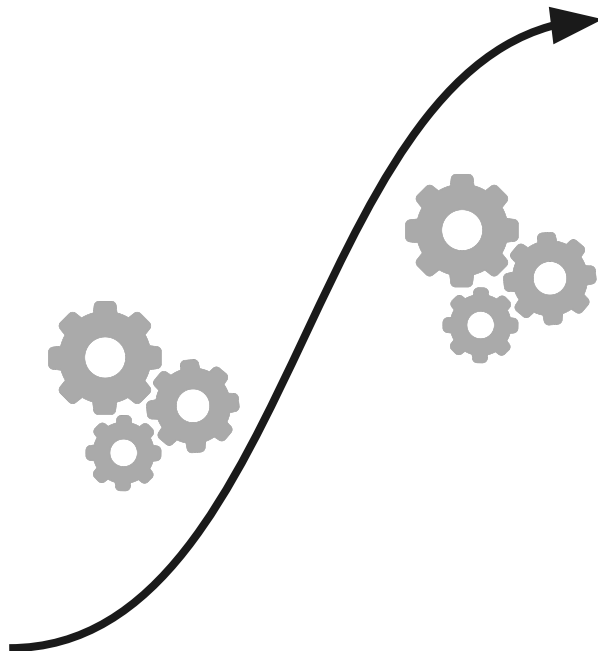
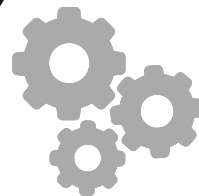
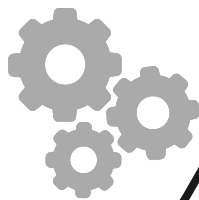


Knowledge



Observations

111
001011
01100110
10011011011
01101100
101101
111



Knowledge

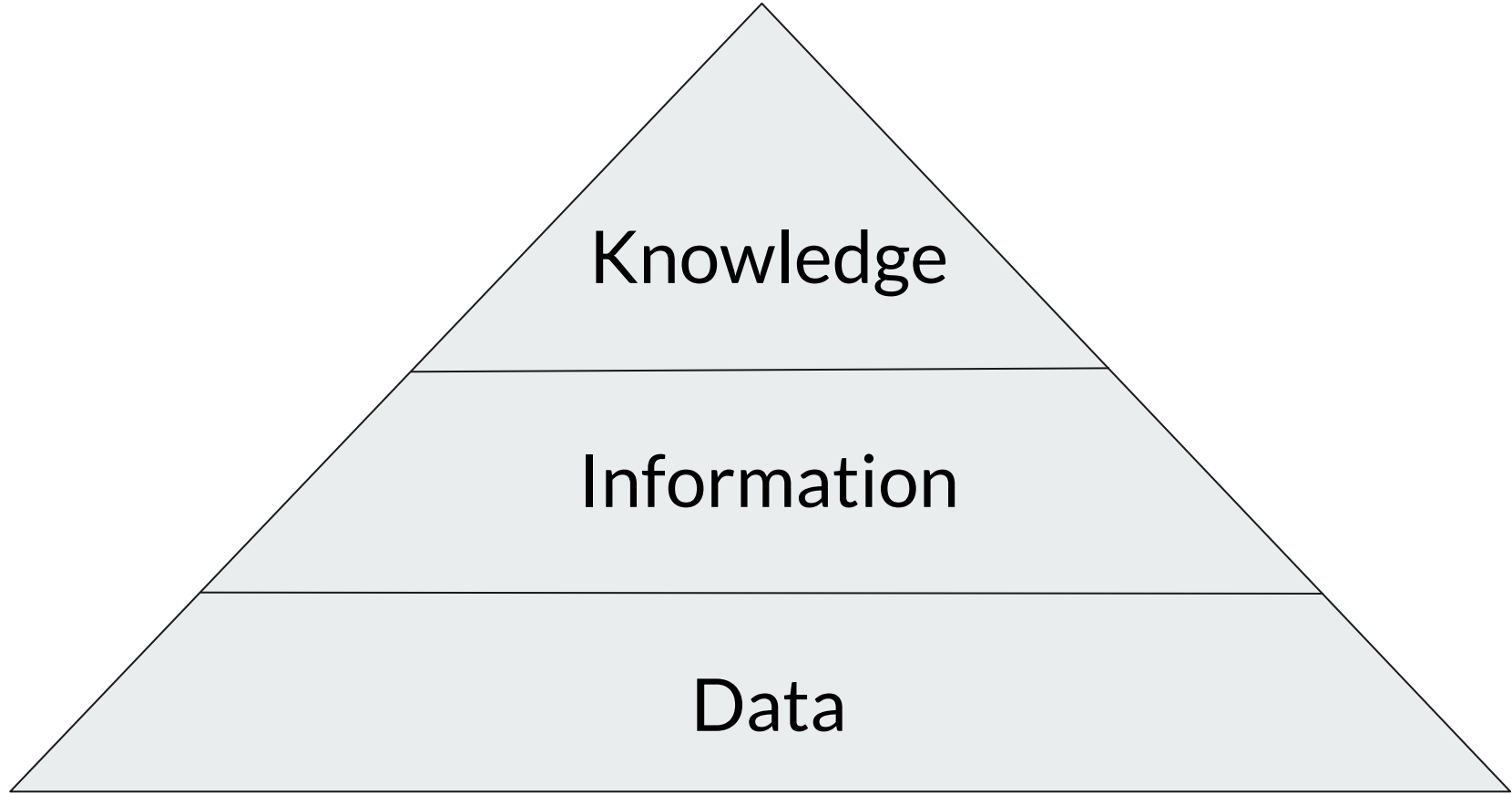


The background of the entire image is a photograph of a sunset or sunrise over a body of water. The sky is filled with wispy clouds, illuminated from below with warm orange and yellow light, transitioning to a deep blue at the top. The water in the foreground is calm, reflecting the colors of the sky. In the distance, a dark, silhouetted shoreline is visible. The ICOS logo is positioned in the top left corner, consisting of the letters 'ICOS' in a large, bold, black sans-serif font. To the right of 'ICOS' is a vertical line, followed by three small colored dots (red, blue, red) and the text 'INTEGRATED CARBON OBSERVATION SYSTEM' in a smaller, black, all-caps sans-serif font.

ICOS

INTEGRATED
CARBON
OBSERVATION
SYSTEM

**Knowledge
through
observations**



It's more complicated

21

21°C

The temperature in my living
room is 21°C

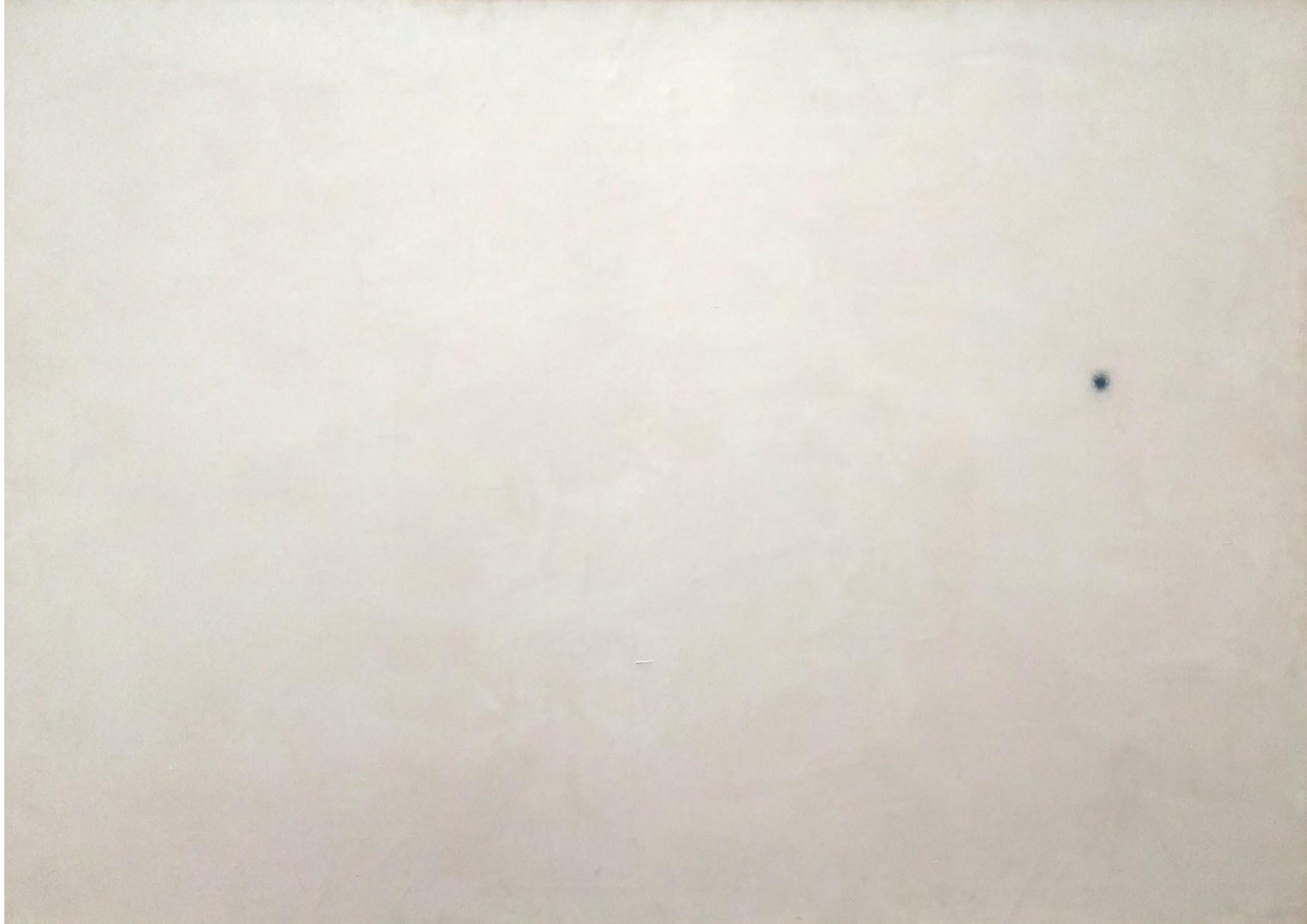
My thermometer observes
that the temperature in my
living room is 21°C

**This is an observation
and it is information**

**So, what are data, information,
knowledge?**

Datum

Joan Miró
Landscape
(1968)



Datum

A datum is a difference within some context

Floridi, L. (2011). The Philosophy of Information. Oxford University Press.



Primary and derivative data

- Primary data are the principal data stored, for example in a database
 - For instance, numerical values resulting from observation activities
 - Measurement data acquired from sensor networks
- Derivative data are data that are extracted from some (primary) data
 - Here, primary data used as indirect sources
 - About things other than those directly addressed by the primary data themselves



Information

- An item σ is an instance of information if
 - σ consists of n data, $n \geq 1$
 - the data are well formed
 - the well-formed data are meaningful
 - the meaningful data are truthful

Floridi, L. (2011). The Philosophy of Information. Oxford University Press.



Data interpretation

- Activity carried out by an interpreter through which data becomes information
- Data are uninterpreted symbols with no meaning for the system concerned
- Interpretation occurs within a real-world context and for a particular purpose
- The interpreter thus determines the contextual meaning of data

Aamodt, A. and Nygård, M. 1995. Different roles and mutual dependencies of data, information, and knowledge – An AI perspective on their integration. *Data & Knowledge Engineering*, 16(3): 191–222.
[https://doi.org/10.1016/0169-023X\(95\)00017-M](https://doi.org/10.1016/0169-023X(95)00017-M)



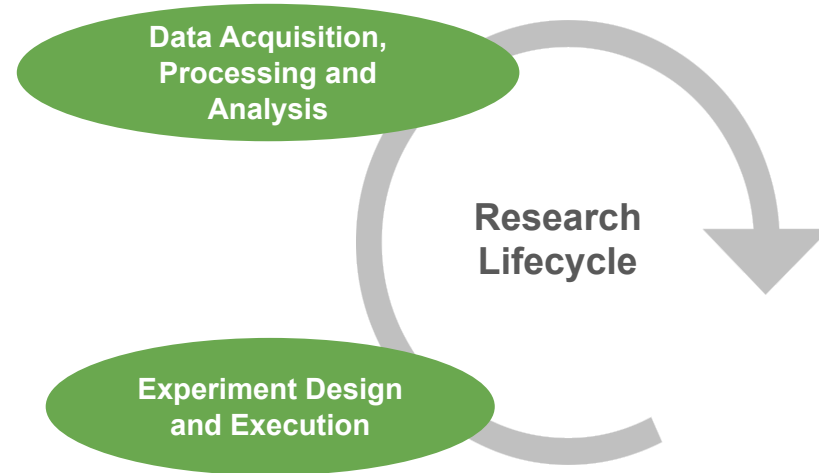
Knowledge

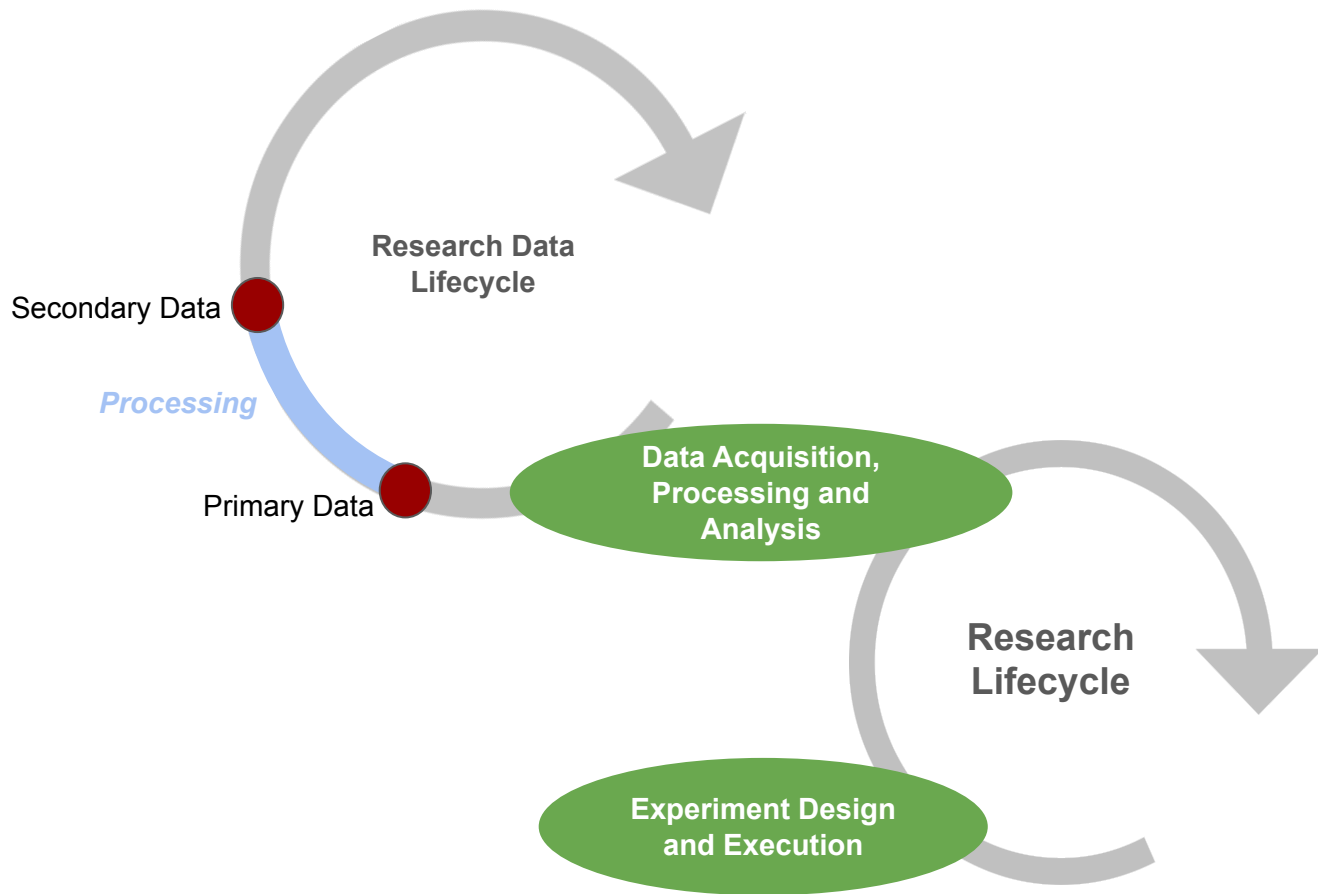
- Learned information
- Information incorporated in an agent's reasoning resources
- Made ready for use within decision processes
- Output of learning processes
- Tacit knowledge is made explicit through externalization

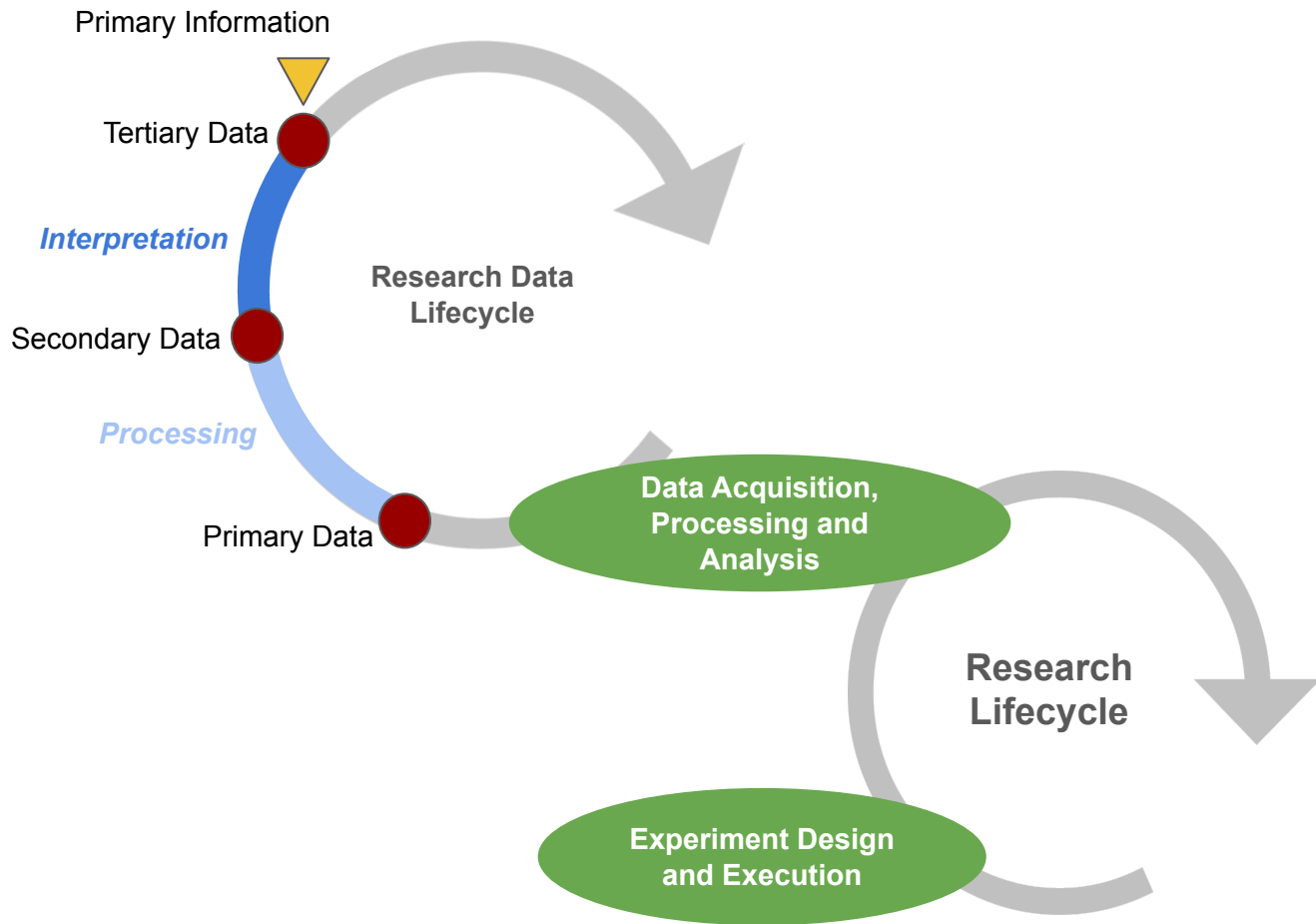
Aamodt, A. and Nygård, M. 1995. Different roles and mutual dependencies of data, information, and knowledge – An AI perspective on their integration. Data & Knowledge Engineering, 16(3): 191–222.
[https://doi.org/10.1016/0169-023X\(95\)00017-M](https://doi.org/10.1016/0169-023X(95)00017-M)

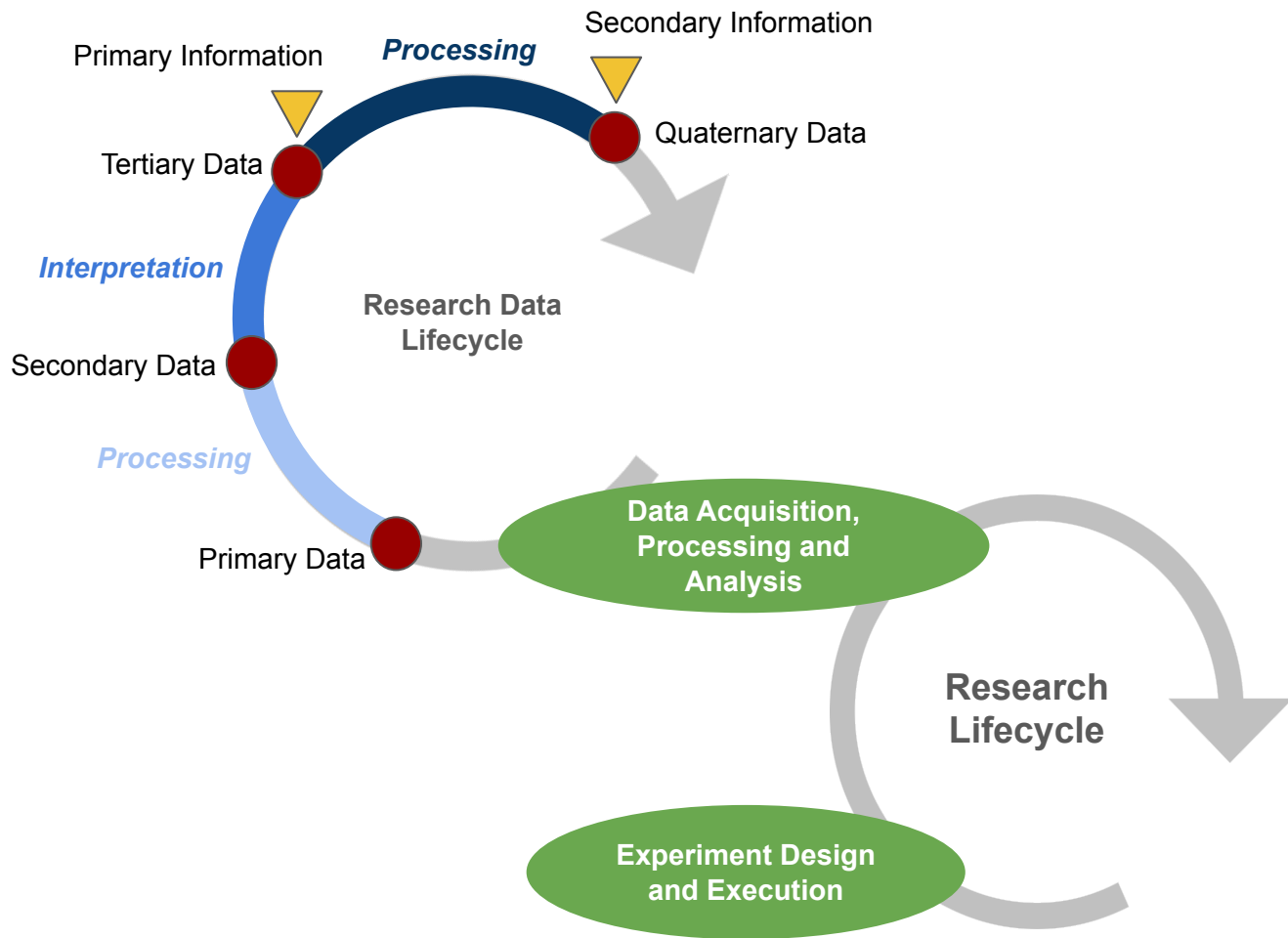
Chyi Lee, C. and Yang, J. 2000. Knowledge value chain. Journal of Management Development, 19(9), 783–794.
<https://doi.org/10.1108/02621710010378228>

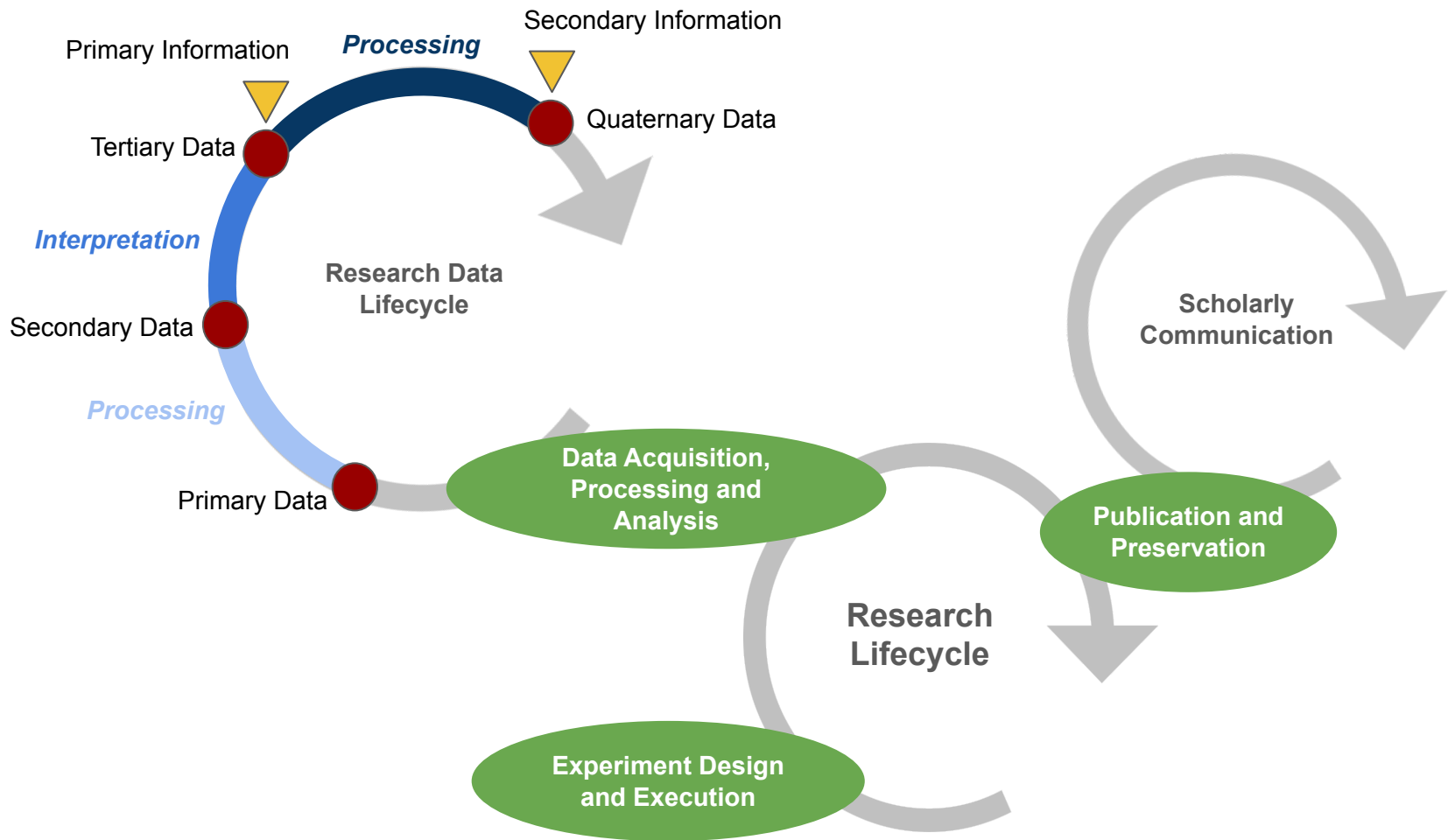
**How does this apply to
research?**

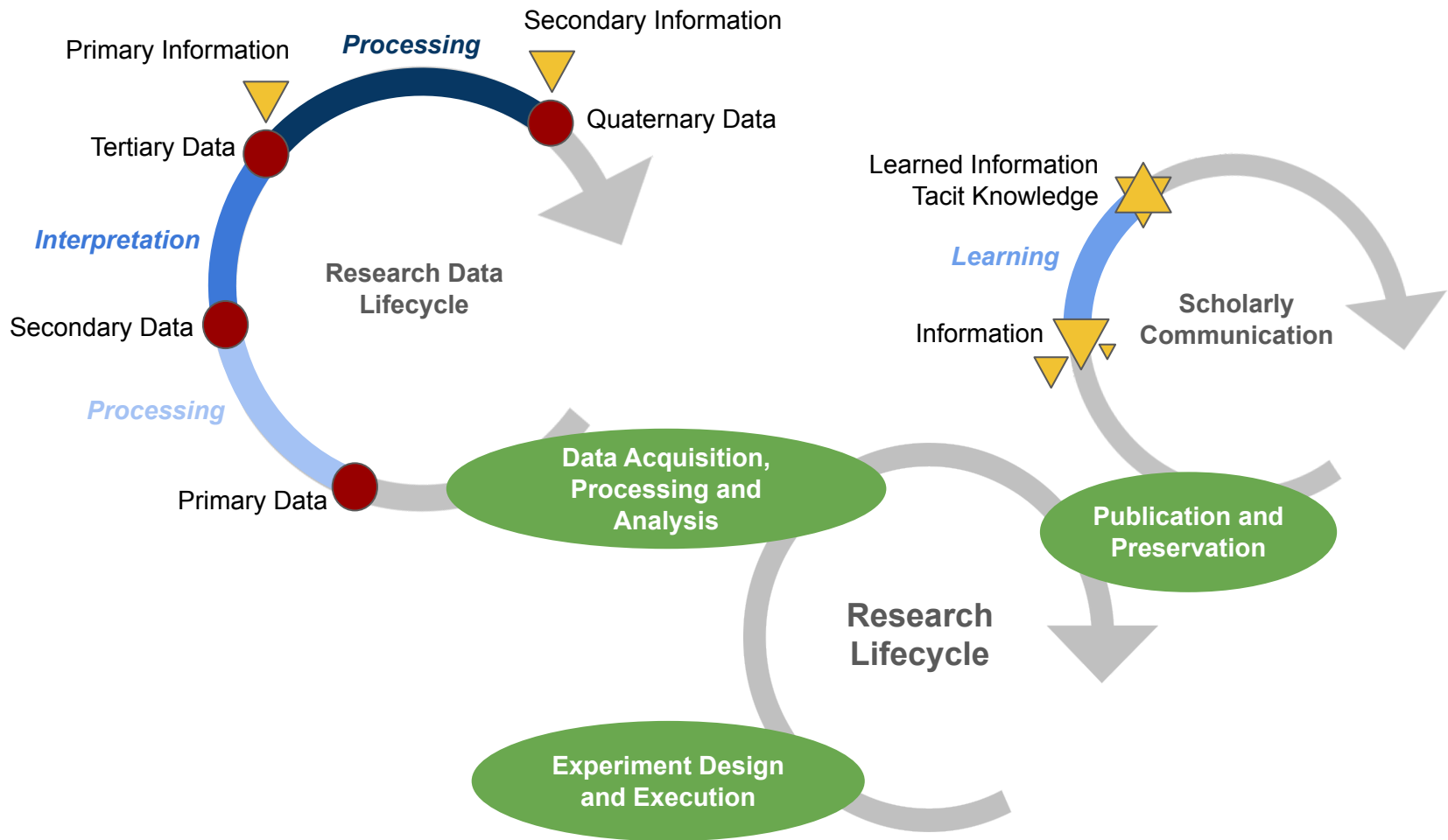


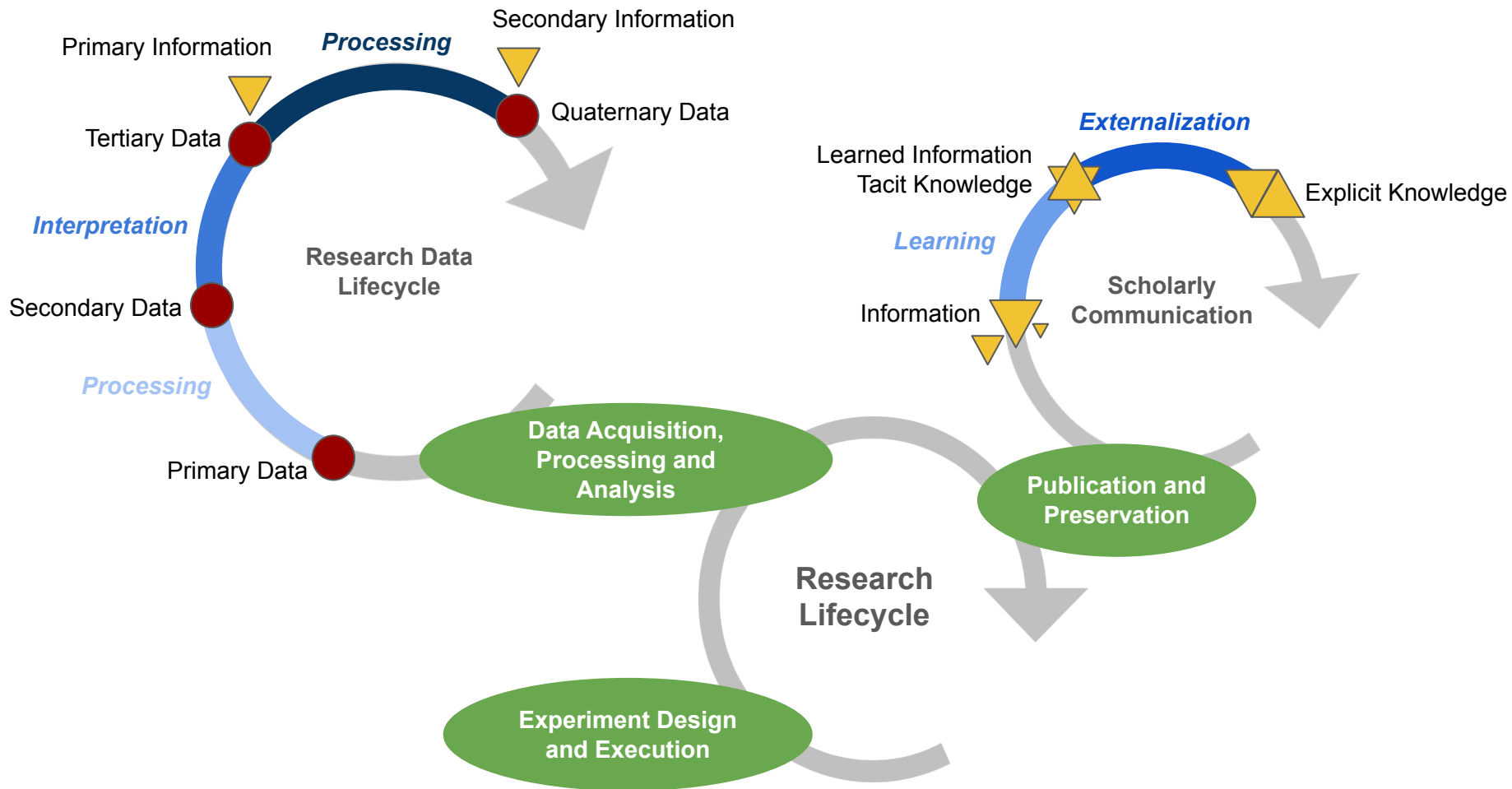
















But for machines ...



... are **data**!

**What the data mean is
(sometimes more, sometimes
less) clear to humans but for
machines meaning is implicit
(not expressly stated)**

**It is important for knowledge
infrastructures to ensure
meaning is explicit and formal**

Information and knowledge-based systems

Not “just” databases

**So, how do we represent
meaning in information
systems?**

—

The technology framework

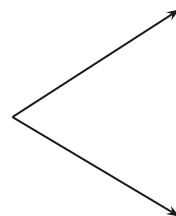
21°C

Temperature [°C]
21

21°C



Quantity Value



Numerical value

Unit of measure

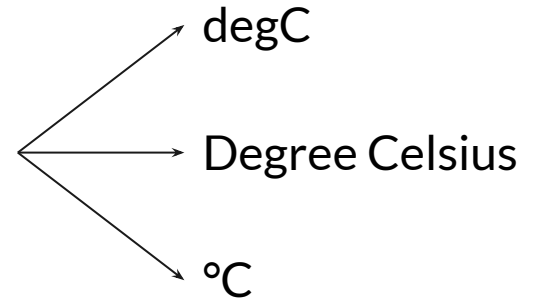
Numerical value: 21

Unit of measure:

http://qudt.org/vocab/unit/DEG_C



Resource





Resource Description Framework (RDF)

- A W3C Recommendation and technology developed within the Semantic Web Activity
- A Statement-centric data model
- A Statement is a triple structure consisting of a subject, predicate and object

(a quantity value, has unit, a unit)

Subject Predicate Object



Resource Description Framework (RDF)

- Subject, predicate and object are Resources identified by Uniform Resource Identifier (URI)
- Object can be a literal value such as a number or a string
- URIs are typically prefixed, e.g. “qudt:” for <http://qudt.org/schema/qudt/>
- The resulting (RDF) data can be processed, e.g. querying using SPARQL

(http://qudt.org/vocab/unit/DEG_C,
 <http://www.w3.org/1999/02/22-rdf-syntax-ns#type>,
 <http://qudt.org/schema/qudt/Unit>)

Subject

Predicate

Object

(unit:DEG_C, rdf:type, qudt:Unit)

With prefixes

My thermometer observes
that the temperature in my
living room is 21°C

:obs1 rdf:type sosa:Observation .
:obs1 sosa:madeBySensor :thermometer1 .
:obs1 sosa:observedProperty :temperature .
:obs1 sosa:hasFeatureOfInterest :room1 .
:obs1 sosa:hasResult :qv1 .
:qv1 qudt:numericValue "21" .
:qv1 qudt:unit unit:DEG_C .

My thermometer observes that the
temperature in my living room is 21°C

**The meaning of data is now
explicit and formal**

**Meaning is formalized using
unambiguous and identified
terms of shared vocabularies**



Vocabularies

- Concepts and relations (qudt:Unit, sosa:Observation, rdf:type, etc.) are constituents of vocabularies
- Also known as terminologies, controlled vocabularies, thesauri, taxonomies, ontologies
- They are most useful if sociotechnical systems widely share them
- This is key to attain *semantic* data interoperability
- There exist numerous languages to specify vocabularies (e.g. RDFS, OWL)
- These enable stating that a Resource is a concept or a relation, a sub concept of another concept, ...

<http://www.w3.org/ns/sosa/Observation>

Lookup by human

Lookup by machine

4.3.2.2 `sosa:Observation`

IRI: <http://www.w3.org/ns/sosa/Observation>

a OWL Class

Observation - Act of carrying out an (Observation) `Procedure` to estimate or calculate a value of a property of a `FeatureOfInterest`. Links to a `Sensor` to describe what made the `Observation` and how; links to an `ObservableProperty` to describe what the result is an estimate of, and to a `FeatureOfInterest` to detail what that property was associated with.

Example

The activity of estimating the intensity of an Earthquake using the Mercalli intensity scale is an `Observation` as is measuring the moment magnitude, i.e., the energy released by said earthquake.

Restrictions

`sosa:madeBySensor` **EXACTLY 1**
`sosa:madeBySensor` **ONLY** `sosa:Sensor`
`sosa:usedProcedure` **ONLY** `sosa:Procedure`
`sosa:hasFeatureOfInterest` **EXACTLY 1**
`sosa:hasFeatureOfInterest` **ONLY** `sosa:FeatureOfInterest`
`sosa:observedProperty` **EXACTLY 1**
`sosa:observedProperty` **ONLY** `sosa:ObservableProperty`
`ssn:wasOriginatedBy` **EXACTLY 1**
`ssn:wasOriginatedBy` **ONLY** `ssn:Stimulus`
`sosa:phenomenonTime` **EXACTLY 1**
`sosa:hasResult` **MIN 1**
`sosa:hasResult` **ONLY** `sosa:Result`
`sosa:resultTime` **EXACTLY 1**

```
sosa:Observation rdf:type owl:Class .  
sosa:Observation rdfs:label "Observation".  
  
...
```

How do I find the vocabulary I need?

EMBL-EBI

Services | Research | Training | About us

Ontology Lookup Service

Home | Ontologies | Documentation | About

Welcome to the EMBL-EBI Ontology Lookup Service.

Search OLS...

Examples: [diabetes](#), [GO:0098743](#)

Looking for a pa...

About OLS

The Ontology Lookup Service (OLS) is a repository for biomedical ontologies that aims to provide a single point of access to the latest ontology versions. You can browse the ontologies through the website as well as programmatically

Related Tools

In addition to OLS the SPOT team also provides the [OxO](#), [Zooma](#) and [Webulous](#) services. [OxO](#) provides cross-ontology mappings between terms from different ontologies. [Zooma](#) is a service to assist in

Contact

For feedback, or

QUDT

Home

About

Contact

Join

Resources

Search for a class

Enter a class, e.g. Shape, Trait, etc...

Advanced Search

Find a semantic resource (ontology, thesaurus, etc.)

Start entering ontology name, e.g. [PhytoTraits](#), then choose from list

Browse Ontologies

Ecoportal Statistics

Ontologies	6
Classes	137

Ontology Visits (June 2019)

WORMS

World Register of Marine Species

Quick search...

QUADT

QUADT.org is a 501(c)(3) not-for-profit organization founded to provide semantic specifications for units of measure, quantity kind, dimensions and data types. QUADT is an advocate for the development and implementation of standards to quantify data expressed in RDF and JSON. Our mission is to improve interoperability of data and the specification of information structures through industry standards for Units of Measure, Quantity Kinds, Dimensions and Data Types.

QUADT.org is a member of the World Wide Web Consortium (W3C)

Why QUADT.org

QUADT.org exists to make the QUADT Ontologies, derived models and vocabularies available to the public. Originally, QUADT models were developed for the NASA Exploration Initiatives Ontology Models (NEXIOM) project, a Constellation Program initiative at the AMES Research Center (ARC).

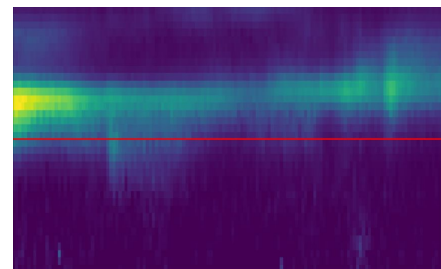
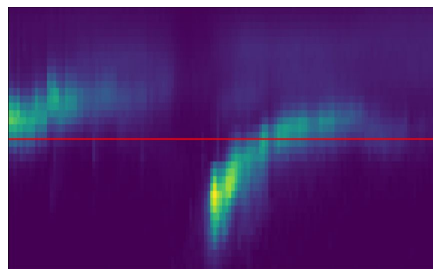
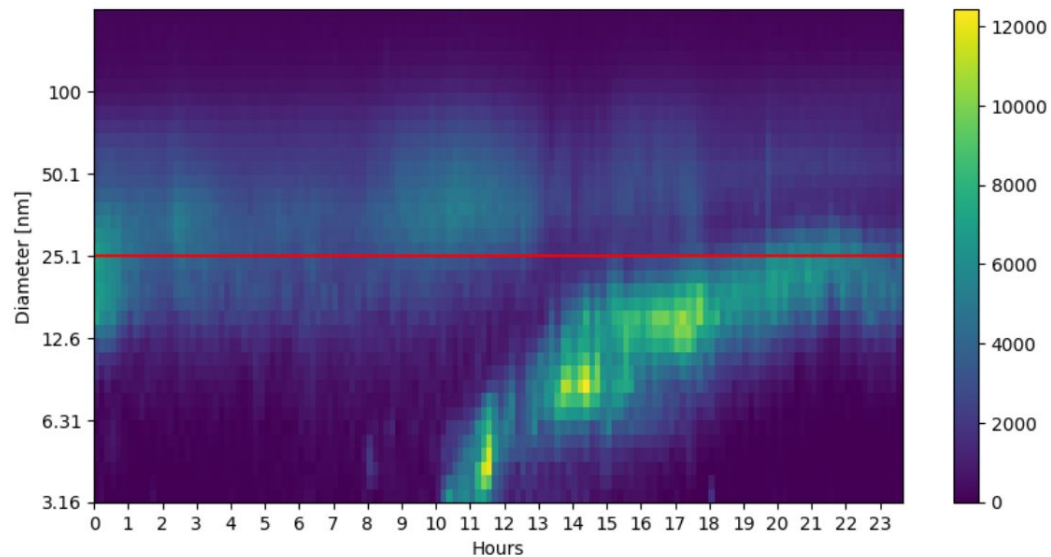
An au



Part II

**A knowledge infrastructure in
atmospheric physics as a case in
point**





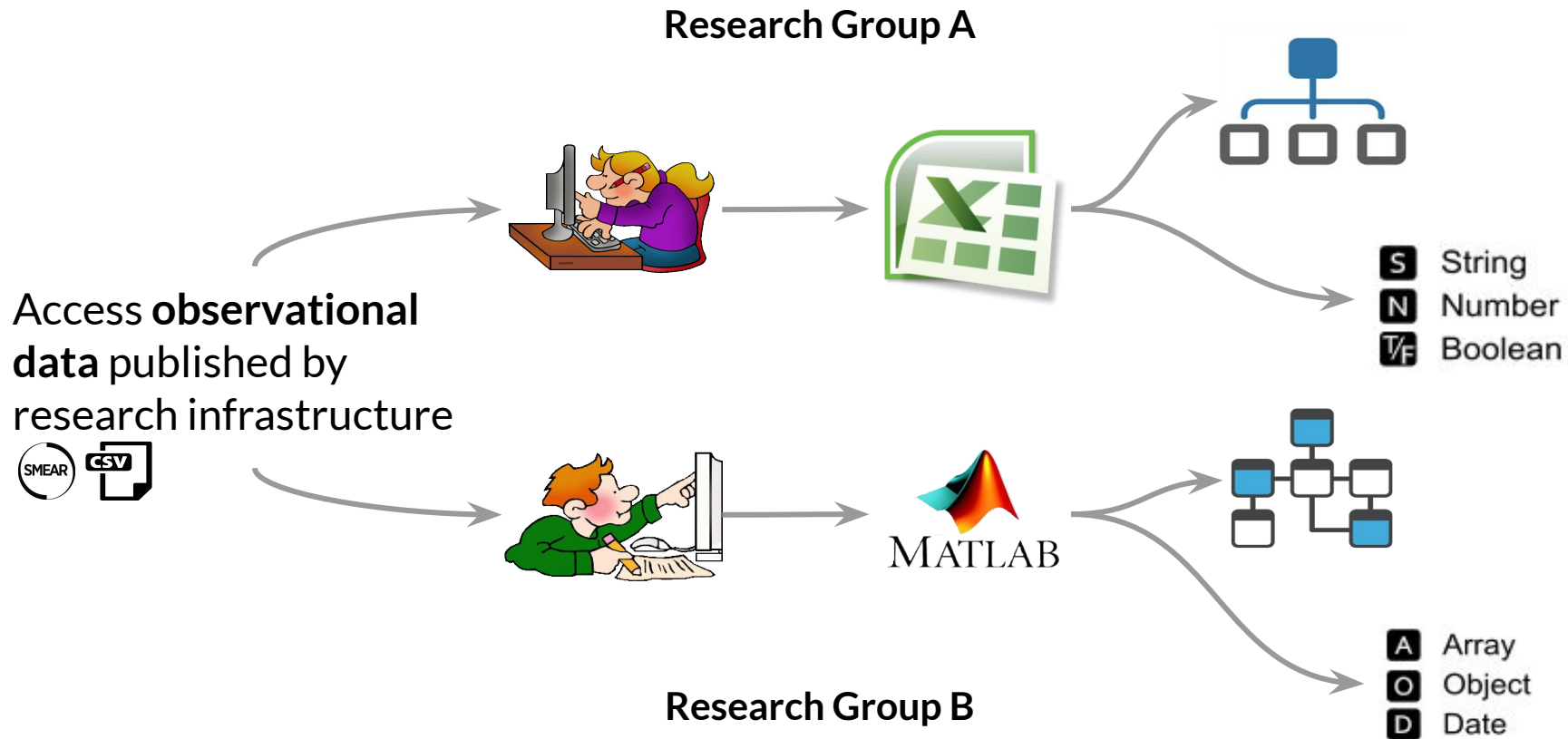
**Measured is particle size distribution,
the primary observational data.**

**Described are new particle formation
events, the derived data (actually,
information)**

—

Let's look at how two research groups implement data analysis.

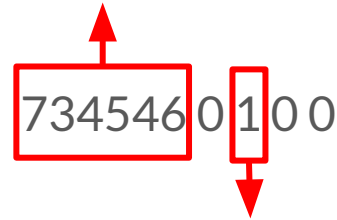
Is the meaning of data explicit, unambiguous, formal?





734545 0 0 0 1
734546 0 1 0 0
734547 0 0 0 1
734548 0 0 0 1
734549 0 0 0 1
734550 0 0 0 1
734551 0 0 0 1
734552 0 0 1 0

MATLAB datenum



Class Ib event

Encoding of data describing NPFE
by Research Group A



2011-07-04,2

2011-07-04,2

2011-07-05,4

2011-07-06,4

2011-07-07,4

2011-07-08,4

2011-07-09,NE

2011-07-10,BD

Encoding of data describing NPFE
by Research Group B



Classifications

Research Group A

Class I: parameters derived with good confidence

Class Ia: very clear and strong NPF

Class Ib: other Class I events

Class II: parameter derivation not possible

Non Event

Undefined

Bad Data

Partly Bad Data

dal Maso et al. (2005) BOREAL ENV. RES. Vol. 10:323-336

Research Group B

Class 0: Undefined

Class 1: Strong NPF

Class 2: Intermediate NPF

Class 3: Weak NPF

Class 4: ?

NE: Non Event

BD: Bad Data

Hamed et al. (2007) Atmos. Chem. Phys. 7:355-376



Problem: Heterogeneous Derived Data

734545 0 0 0 1	2011-07-04,2
734546 0 1 0 0	2011-07-04,2
734547 0 0 0 1	2011-07-05,4
734548 0 0 0 1	2011-07-06,4
734549 0 0 0 1	2011-07-07,4
734550 0 0 0 1	2011-07-08,4
734551 0 0 0 1	2011-07-09,NE
734552 0 0 1 0	2011-07-10,BD

Problem: Implicit Semantics



734545 0 0 0 1
734546 0 1 0 0
734547 0 0 0 1
734548 0 0 0 1
734549 0 0 0 1
734550 0 0 0 1
734551 0 0 0 1
734552 0 0 1 0

Can we improve this?

Part III

Hands-on exercise

GitHub, Inc. (US) | <https://github.com/markusstocker/carbon-workshop> 110% ... ☆

markusstocker / carbon-workshop Watch 0 Star 0 Fork 0

Code Issues 0 Pull requests 0 Projects 0 Wiki Security Insights Settings

Hands-on assignment for the Semantic Data Science lecture at the TERENO-NEON Carbon Workshop 2019 Edit

Manage topics

25 commits 1 branch 0 releases 1 contributor

Branch: master New pull request Create new file Upload files Find File Clone or download

markusstocker Update README.md Latest commit 7952305 1 minute ago

.gitignore fixed json 7 days ago

<https://github.com/markusstocker/carbon-workshop>

data.csv-metadata.json added name 6 days ago

requirements.txt updated 5 days ago

README.md

TERENO-NEON Carbon Workshop 2019

This hands-on assignment for the Semantic Data Science lecture at the TERENO-NEON Carbon Workshop 2019 demonstrates how research data can be described for their meaning, here using the [CSV on the Web](#) W3C Recommendation. To complete the assignment, you should fork this repository and then execute the assignment [Jupyter Notebook](#) (assignment.ipynb) using [binder](#).

The material is licensed [Creative Commons Attribution 4.0 International \(CC BY 4.0\)](#).

[launch binder](#)



Aims


- Building on the presented case in atmospheric physics
- Interpret observational data and generate derivative data about events, in CSV
- Inspect the corresponding “metadata” that specifies data semantics
- Convert the CSV data into RDF as a format with explicit and formal data semantics
- Understand the resulting RDF data
- Formulate SPARQL queries on RDF to demonstrate a kind of RDF data processing




Steps

- Fork the repository <https://github.com/markusstocker/carbon-workshop>
- Use <http://mybinder.org> to launch the forked repository
- Wait until Jupyter Notebook (Lab) is launched
- Open assignment.ipynb and follow the instructions


 markusstocker / carbon-workshop


 Watch ▾ 0

 Star 0

 Fork 0

 Code


 Issues 0

 Pull requests 0

 Projects 0

 Wiki

 Security

 Insights

 Settings

Hands-on assignment for the Semantic Data Science lecture at the TERENO-NEON Carbon Workshop 2019

Edit

[Manage topics](#)

 25 commits

 1 branch

 0 releases

 1 contributor

Branch: master ▾

New pull request

Create new file

Upload files

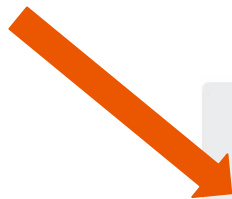
Find File

Clone or download ▾



Turn a Git repo into a collection of interactive notebooks

Have a repository full of Jupyter notebooks? With Binder, open those notebooks in an executable environment, making your code immediately reproducible by anyone, anywhere.



Build and launch a repository

GitHub repository name or URL

GitHub ▾

Git branch, tag, or commit

Path to a notebook file (optional)

File ▾

launch

Copy the URL below and share your Binder with others:



Copy the text below, then paste into your README to show a binder badge:



NPFE Data with CSV on the Web

This notebook exemplifies how data about [new particle formation events](#) can be described using [CSV on the Web](#).

We interpret particle size distribution data (primary data) as measured by an observation system of the [SMEAR](#) research infrastructure in order to detect the occurrence of new particle formation events on selected days in [Hyytiälä](#), Finland. Detected events are then described, whereby we generate secondary (derivative) data about events. These data are stored to disk in CSV format.

In addition, we use CSV on the Web to describe the secondary CSV data using common, shared terminology that is unambiguously identified and described on the web. This makes the secondary CSV data more interoperable, reusable and understandable to machines. For instance, we can use the description to transform the CSV data into RDF and then leverage SPARQL to query this data.

Before we start, we need to load required Python modules as well a few user-defined functions.

```
In [ ]: import requests, io, os, pandas as pd, numpy as np
        from urllib.parse import urlencode
        from pytz import timezone
        from datetime import datetime, timedelta
        from csvlib import CSVWConverter
        from matplotlib import pyplot as plt
        from rdflib.plugins.sparql.results.csvresults import CSVResultSerializer

        def fetch(date):
            time_from = timezone('Europe/Helsinki').localize(datetime.strptime(date, '%Y-%m-%d'))
            time_to = time_from + timedelta(days=1)

            query = {
                'table': 'HYV_DMPS', 'quality': 'ANY', 'averaging': 'NONE', 'type': 'NONE',
                'from': str(time_from), 'to': str(time_to), 'variables': 'd316e1,d355e1,d398e1,\
                'd447e1,d501e1,d562e1,d631e1,d708e1,d794e1,d891e1,d100e2,d112e2,d126e2,d141e2,d158e2,\
                'd178e2,d200e2,d224e2,d251e2,d282e2,d316e2,d355e2,d398e2,d447e2,d501e2,d562e2,d631e2,\
                'd708e2,d794e2,d891e2,d100e3,d112e3,d126e3,d141e3,d158e3,d178e3,d200e3'
            }

            url = 'https://avaa.tdata.fi/smea-services/smeardata.jsp?' + urlencode(query)
            response = requests.post(url)

            return pd.read_csv(io.StringIO(response.text))

        def plot(data):
            d = data.copy(deep=True)
            d = d.iloc[:, 6:].values
            m = len(d)
            n = len(d[0])
            x = range(0, m)
            y = range(0, n)
            x, y = np.meshgrid(x, y)
            z = np.transpose(np.array([row[1:] for row in d]).astype(np.float))
            plt.figure(figsize=(10, 5), dpi=100)
            plt.pcolormesh(x, y, z)
```

```

1 {
2   "@context": "http://www.w3.org/ns/csvw",
3   "url": "data.csv",
4   "tableSchema": {
5     "aboutUrl": "http://avaa.tdata.fi/{date}",
6     "columns": [{
7       "name": "date",
8       "propertyUrl": "http://purl.obolibrary.org/obo/STATO_0000093",
9       "dc:description": "The ISO date at which the event occurred",
10      "datatype": {
11        "base": "date",
12        "format": "yyyy-MM-dd"
13      }
14    }, {
15      "name": "start",
16      "propertyUrl": "http://purl.obolibrary.org/obo/RO_0002537",
17      "dc:description": "The ISO event start time",
18      "datatype": {
19        "base": "time",
20        "format": "HH:mm"
21      }
22    }, {
23      "name": "end",
24      "propertyUrl": "http://purl.obolibrary.org/obo/RO_0002538",
25      "dc:description": "The ISO event end time",
26      "datatype": {
27        "base": "time",
28        "format": "HH:mm"
29      }
30    }, {
31      "name": "class",
32      "propertyUrl": "http://purl.obolibrary.org/obo/OBI_0000999",
33      "dc:description": "The event classification according to the dal Maso et al. scheme",
34      "datatype": {
35        "base": "string",
36        "format": "Ia|Ib|II"
37      }
38    }, {
39      "name": "type",
40      "virtual": true,
41      "propertyUrl": "rdf:type",
42      "valueUrl": "http://purl.obolibrary.org/obo/ENVO_01001372"
43    }
44  ]
45 }
46 }

```

date	start	end	class
2007-05-05	12:00	14:00	lb
2007-04-15	11:00	12:00	la
2013-04-04	10:00	11:30	la



```
<http://avaa.tdata.fi/2007-04-15> a ns1:ENVO_01001372 ;
    ns1:OBI_0000999 <http://avaa.tdata.fi/class/Ia> ;
    ns1:RO_0002537 "11:00:00"^^xsd:time ;
    ns1:RO_0002538 "12:00:00"^^xsd:time ;
    ns1:STATO_0000093 "2007-04-15"^^xsd:date .
```

```
<http://avaa.tdata.fi/2007-05-05> a ns1:ENVO_01001372 ;
    ns1:OBI_0000999 <http://avaa.tdata.fi/class/Ib> ;
    ns1:RO_0002537 "12:00:00"^^xsd:time ;
    ns1:RO_0002538 "14:00:00"^^xsd:time ;
    ns1:STATO_0000093 "2007-05-05"^^xsd:date .
```

```
<http://avaa.tdata.fi/2013-04-04> a ns1:ENVO_01001372 ;
    ns1:OBI_0000999 <http://avaa.tdata.fi/class/Ia> ;
    ns1:RO_0002537 "10:00:00"^^xsd:time ;
    ns1:RO_0002538 "11:30:00"^^xsd:time ;
    ns1:STATO_0000093 "2013-04-04"^^xsd:date .
```

Part IV

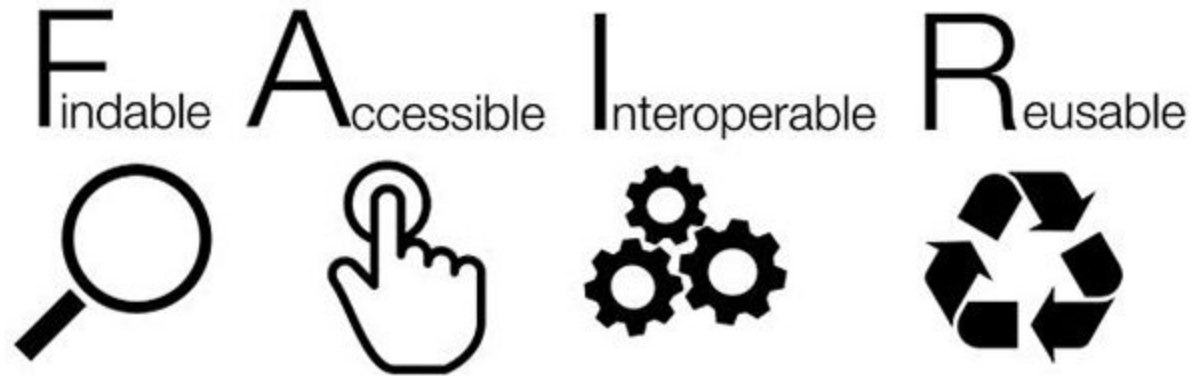
Discussion

How are data semantics explicit and formal here?

```
<http://avaa.tdata.fi/2007-04-15> a ns1:ENVO_01001372 ;  
    ns1:OBI_0000999 <http://avaa.tdata.fi/class/Ia> ;  
    ns1:RO_0002537 "11:00:00"^^xsd:time ;  
    ns1:RO_0002538 "12:00:00"^^xsd:time ;  
    ns1:STATO_0000093 "2007-04-15"^^xsd:date .
```

```
<http://avaa.tdata.fi/2007-05-05> a ns1:ENVO_01001372 ;  
    ns1:OBI_0000999 <http://avaa.tdata.fi/class/Ib> ;  
    ns1:RO_0002537 "12:00:00"^^xsd:time ;  
    ns1:RO_0002538 "14:00:00"^^xsd:time ;  
    ns1:STATO_0000093 "2007-05-05"^^xsd:date .
```

```
<http://avaa.tdata.fi/2013-04-04> a ns1:ENVO_01001372 ;  
    ns1:OBI_0000999 <http://avaa.tdata.fi/class/Ia> ;  
    ns1:RO_0002537 "10:00:00"^^xsd:time ;  
    ns1:RO_0002538 "11:30:00"^^xsd:time ;  
    ns1:STATO_0000093 "2013-04-04"^^xsd:date .
```



Wilkinson, M. D. *et al.* (2016) The FAIR Guiding Principles for scientific data management and stewardship. *Sci. Data* 3:160018
<https://doi.org/10.1038/sdata.2016.18>

Do you publish your (CSV) research data?

Are you going to use these technologies?