# Apport des ressources Semantiques pour la gestion de données

Contact: Pascal.Neveu@inra.fr

# Data Challenge

**Context: more and more data!**

- Cheap storage capacity and high speed network
  e.g. **1 Gigabyte price** : $400K in 1980, $10K in 1990, $10 in 2000,
  **now less than $0.01**
- Many heterogeneous devices, simulations, machine learning, Internet
  data sources (Open, collaborative, etc) are available

## Make data valuable!

- **Knowledge discovery**
- **Decision support**
- **New services and Open Science**

- *Population treatment → individualized treatment*
- *When data did not quite match what we expect!*
- *Which theories/models are consistent and which ones are not*!
- *...*

**Need: A new generation of Information Systems**
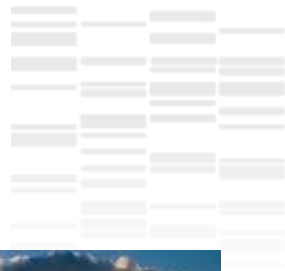
# Complex Data

Different scales

Genome    Organ    Plant    Field    Region

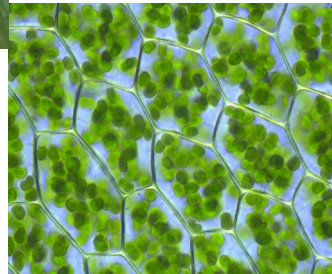# Complex Data

**Different interactions**

# Complex Data

**Different crops**

**Different stages and transformations**

# Complex Data

From various contexts

« omics » Platforms

## Various data complex types

Genomics

Composition and the structure of biopolymers
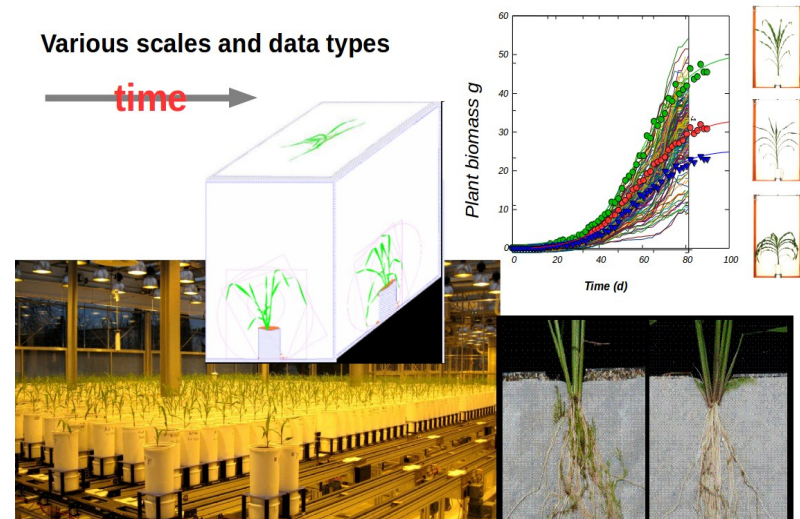
Quantification of metabolites and enzyme activities



Green house Platforms

Various scales and data types

time



Field Platforms

**Various scales and data types**
- Cell, organ, plant, population
- Images, hyperspectral, spectral, sensors, human readings...

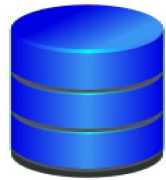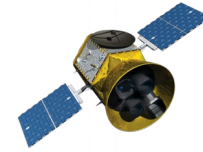time

Thousands of micro-plots



Farm Platforms

**Various scales and data types from thousands of farms**
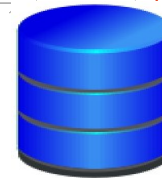- organ, plant, population, site
- Images, sensors, human readings...

time

# Complex Data

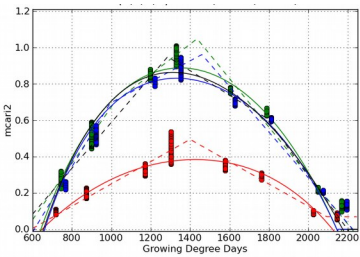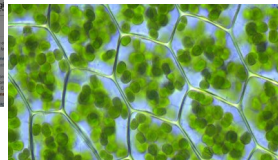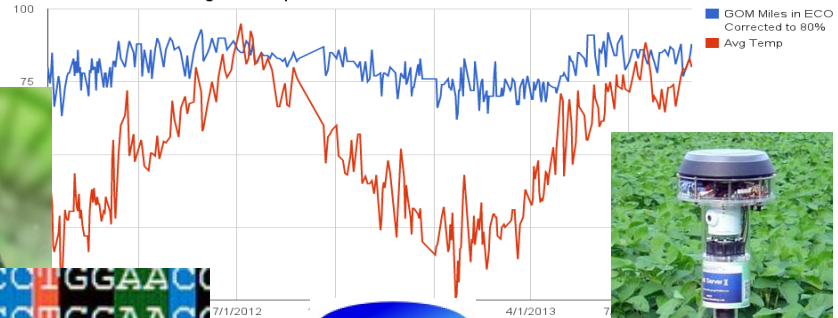# Complex Data

- **Orphan data → Worthless!**
- **Data has value if they are grouped...**
  **(re)analyse, meta-analyse, visualisation, etc**

# Context

Experiments or Observations

- Expensive, require a lot of resources and often very hard

- Cannot be reproduced

Scientific projets generate large and complex datasets

Strong needs of transparence and reproducibility

**But re-analyses meta-analyses and new analyses**

→ **impossible without information (metadata)**

**Findable: PID (globally unique)**, indexed in portals, standardized and relevant metadata

**Accessible:** open and standardized protocols (internet protocols), **licence rights**

**Interoperable (technology, syntax, semantic):** shared standardized formats, vocabularies and **formal languages for knowledge representation**

**Reusable: provenance**, relevant metadata for understanding

# Some common mistakes we do

- Metadata in file names (not standardized, very often not machine readable, reduces metadata quantity)
  → 2017-Paris-Syrah-irrig-goblet.csv
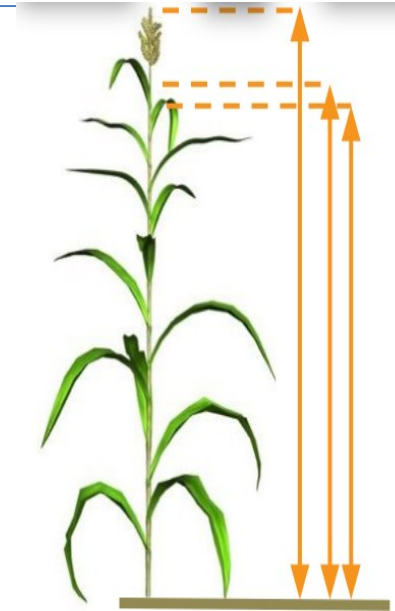  → 2017-St-Paul-Merlot-guyot-not_irrig.csv

- Same name for several variables

- Several names for same variable

- Sharing unstable variables

- Data are stored on personal computer

- Sofware parameters (calibration, etc) are lost

- Ambiguous ID

- X-Y position are lost

- Faults are not described

- no or few data links

- etc

Plot566 in 2016

Plot566 in 2017

**How to structure data ?**

# Structuration

**Data structure** enables a computer system to perform store, retrieve, process data and Implement good practices:
- ➔ Make **FAIR data**
- ➔ **Flexible**
- ➔ Ability to allow **understanding (and reproduce) data processing**
- ➔ Ability to enforce DMP and Open Science

Based on 2 **key elements**:

➔ **Identification and Naming convention**
- ➔ Objects: plants, plots, experiments, sensors, events, etc
- ➔ Persistent, unambiguous, resolvable, globally unique

➔ **Semantic and tagging (based on ontologie set)**
- ➔ Controlled vocabulary
- ➔ Formalized relationships between entities
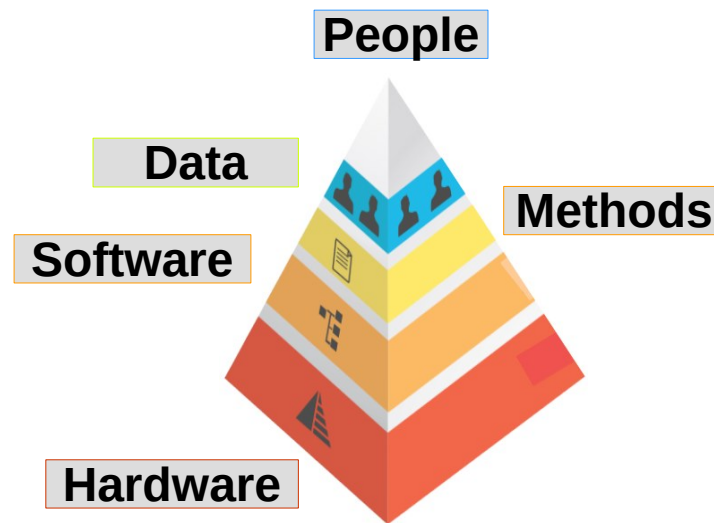- ➔ Data annotation and enrichment

# OpenSILEX

**Open source software set**

- Methods, tools, components to implement information systems for experimental data in agriculture and environment

  Information System: Organized system for the collection, organisation, storage, exchange and treatment of information

**People**

**Data**

**Methods**

**Software**

**Hardware**

# Structuration : PHIS approach

**PHIS is the Information System for Phenomics based on OpenSILEX**

**Scientific objects** (plant, plant organ, plot, etc.) are:
- Identified by **URI** standardized, unambiguous, shared, etc

**Events** (management, faults, meteo, etc)
- Identified by **URI**

Variables, Documents, Observations, Software are associated with these Objects and Events
- Identified by **URI**

**Organisation and linking of Objects and Events → done with a controlled semantic** (reference ontologies, vocabularies, thesaurus, taxonomies) and **application Ontologies**

# PHIS Identification

URI string used to identify a resource (Web standardized syntax)

→ **Standardized (easy integration in Web application)**

`http://subdomain.yourdomain.topdomain/path/identifier`

**Possibility to use prefix**:  m3p: <http://lepse.inra.fr./>



URI of plant :
mp3:arch/2014/pl/000000012

URI of pot :
mp3:arch/2001/pt/000001542

URI of cabin :
mp3:arch/2010/ca/cabine2

URI of camera :
mp3:arch/2011/ss/00003312

URI of image :
mp3:arch/2015/im/000000564

# PHIS Identification

URL identifies what exists on the Web;

**URI identifies, on the Web, what exists;**

IRI  identifies, on the Web, in any language, what exists.

URI of plant :
mp3:arch/2014/pl/000000012

URI of pot :
        mp3:arch/2001/pt/000001542

URI of cabin :
        mp3:arch/2010/ca/cabine2

URI of camera :
mp3:arch/2011/ss/00003312

URI of image :
        mp3:arch/2015/im/000000564

# PHIS Identification

URI string used to identify a resource (Web standardized syntax)

- ➢ **Standardized**

- ➢ **Unambiguous, globally unique**

- ➢ **Resolvable (actionable, dereferencable)**

- ➢ **Persistent (services: B2HANDLE, PIC, PURL, etc)**

URI → the scientific responsible and generated by tools

Resource identifications: standardized, unique, unambiguous

URI of plant :
mp3:arch/2014/pl/000000012

URI of pot :
mp3:arch/2001/pt/000001542

URI of cabin :
mp3:arch/2010/ca/cabine2

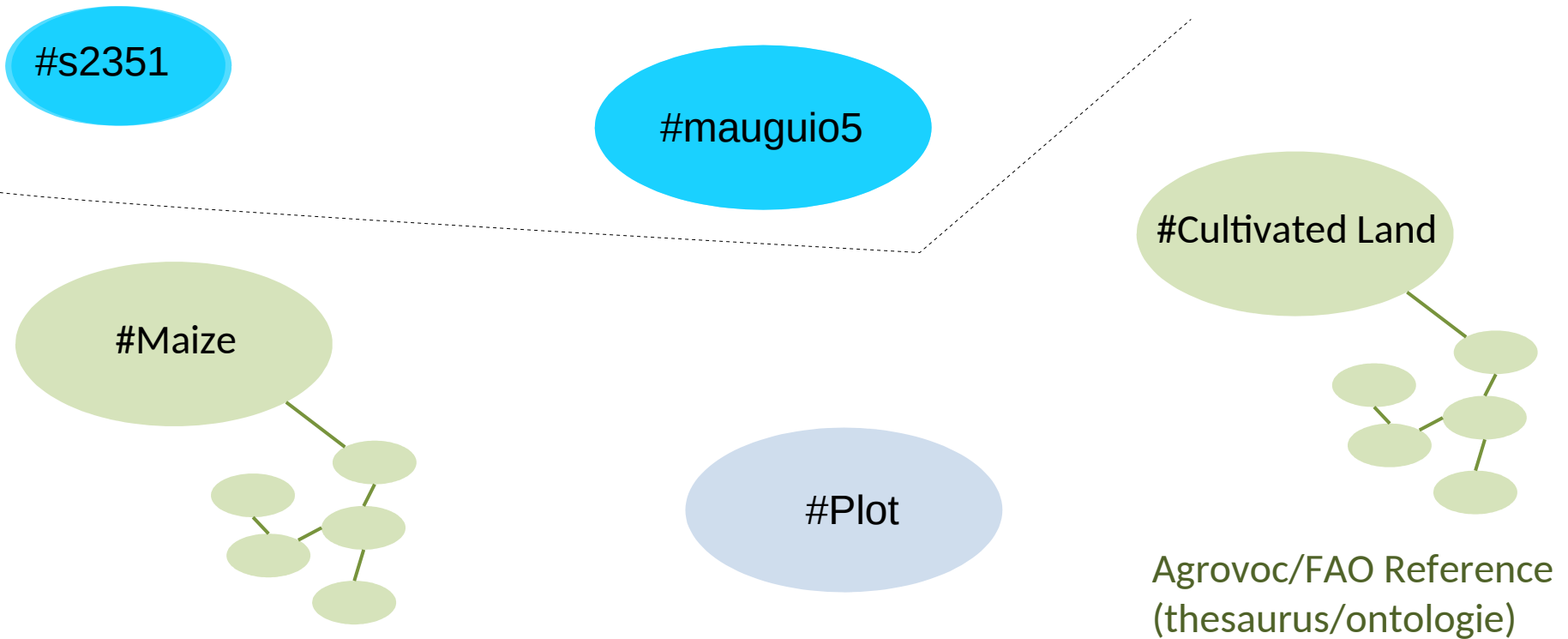URI of camera :
mp3:arch/2011/ss/00003312
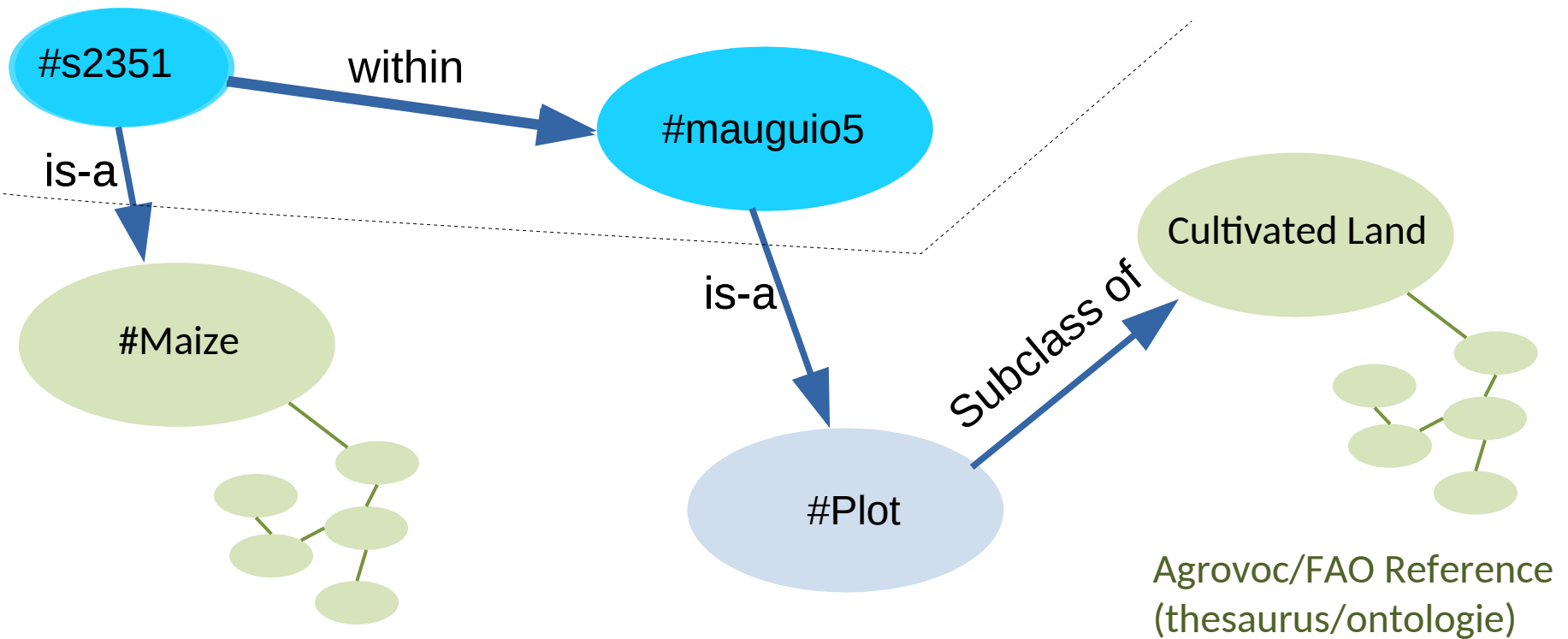
URI of image :
mp3:arch/2015/im/000000564

# PHIS

■ **Metadata / ontologies provide the meaning of data**
→ **Link each data element to a controlled, shared, vocabulary and machine readable**
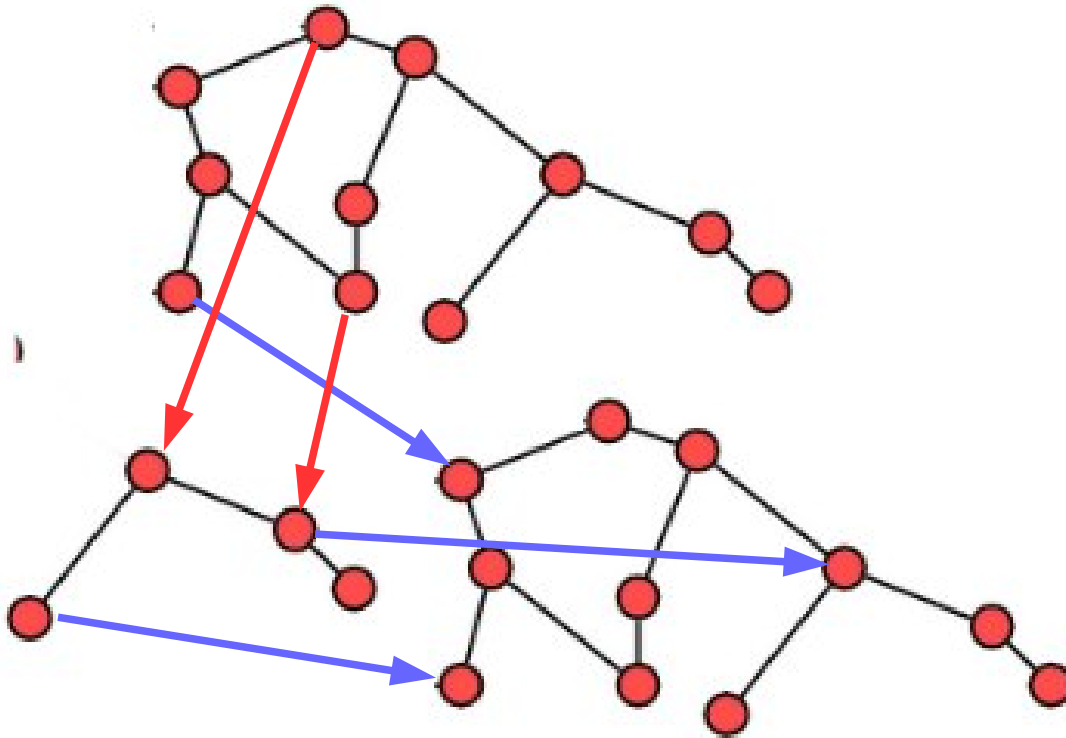→ **Structure the data in a graph**

#s2351

#mauguio5

#Cultivated Land

#Maize

#Plot

Agrovoc/FAO Reference
(thesaurus/ontologie)

# PHIS

- **Metadata / ontologies provide the meaning of data**
    - → **Link each data element to a controlled, shared, vocabulary and machine readable**
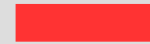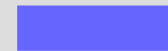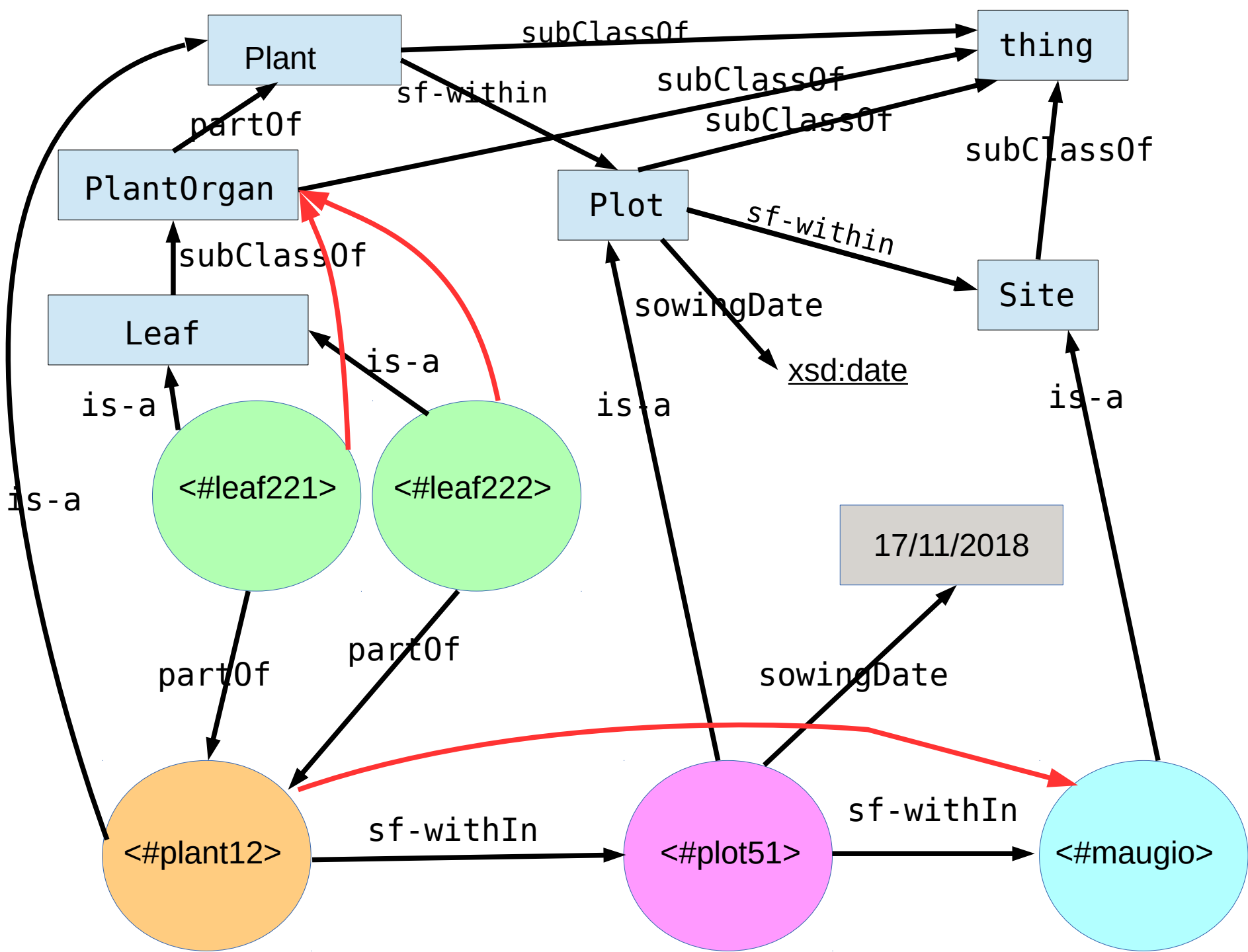    - → **Structure the data in a graph**

#s2351 —within→ #mauguio5

#s2351 —is-a→ #Maize

#mauguio5 —is-a→ #Plot

#Plot —Subclass of→ Cultivated Land

Agrovoc/FAO Reference
(thesaurus/ontologie)

# PHIS

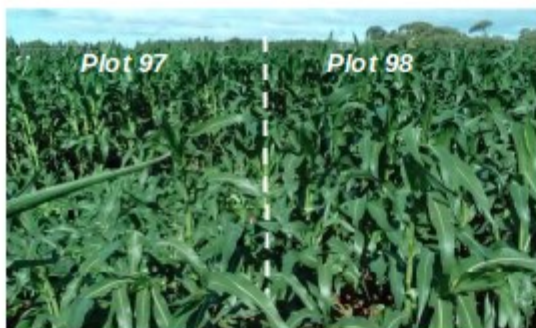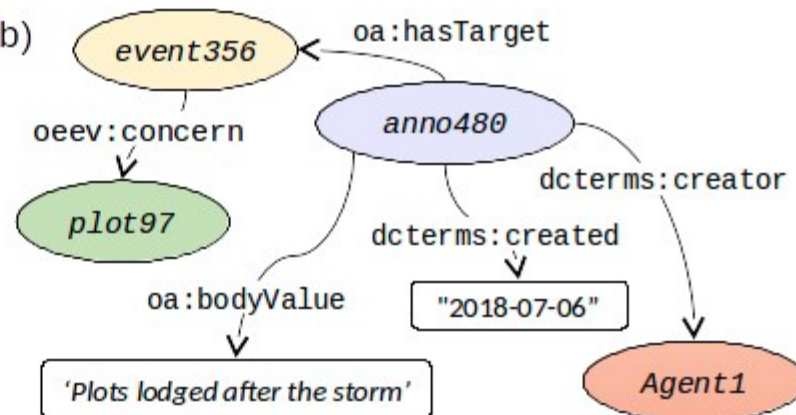**Reference ontologies
See AgroPortal**

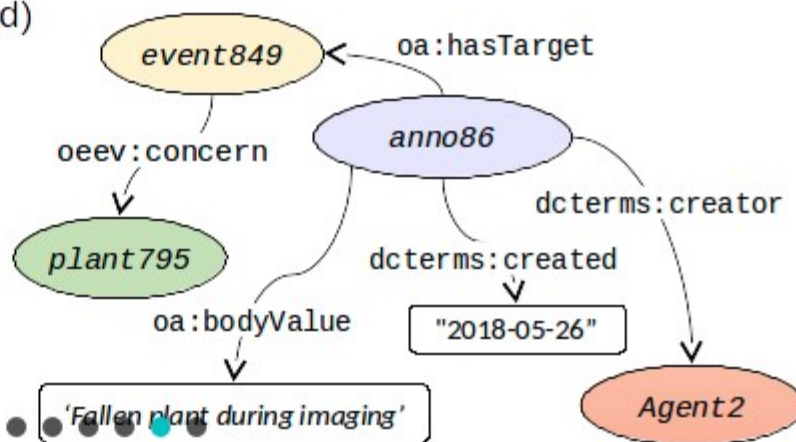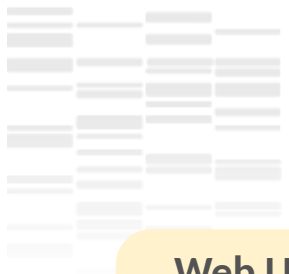**Application ontologies**

# PHIS



(a)

(b) event356 — oa:hasTarget — anno480
event356 — oeev:concern → plot97
anno480 — dcterms:creator → Agent1
anno480 — dcterms:created → "2018-07-06"
anno480 — oa:bodyValue → 'Plots lodged after the storm'
Plot 97 Plot 98

(c)

(d) event849 — oa:hasTarget — anno86
event849 — oeev:concern → plant795
anno86 — dcterms:creator → Agent2
anno86 — dcterms:created → "2018-05-26"
anno86 — oa:bodyValue → 'Fallen plant during imaging'

# OpenSILEX - PHIS

What do <u>we recommended</u> in global context

- Use URI for unambiguous name (in global context)

- Actionable URI for an accessible description of variable
  Description can be read by machine and human

- Try to reuse existing variable if available

- Use standardized/shared representation schema for formalisation
  of new variable (and share it)

  In PHIS

**Variable = Entity + Quality/Quantity + Method + Unit**

**PlantHeight = plant + hauteur + ruler + cm**

# OpenSILEX - PHIS

# Trait – Provenance

# PHIS

✔ **Allows management of huge and complex data**

✔ **Enables and facilitates cloud computing (data center, EGI)**

→ **distributed computing, distributed storage, backup**

✔ **Flexible design**

✔ **International identification (URI and DOI)**

✔ **Semantic management (ontologies, standardized vocabularies)**

✔ **Open technologies , Web APIs and portal interoperability**

✔ **Provenance and reproducibility for data processing**

✔ **Over 10 instances of PHIS** for various installations **(field and greenhouse)**

✔ **Phenoarch instance → Over 700 Tb of data over 10 plant species**

✔ **Other implementations of OpenSilex : WEIS, SunAGRI, SIUE**

✔ **Open Software** - support and development (MISTEA team)

# PHIS and OpenSilex

➢**PHIS** demonstration
- http://phis.inra.fr/
- http://www.opensilex.org/opensilex-sandbox/web/
- Research paper:
  https://nph.onlinelibrary.wiley.com/doi/full/10.1111/nph.15385

➢How to contribute to OpenSILEX?
- Github repository: https://github.com/OpenSILEX/
- Developer documentation: https://opensilex.github.io/docs-community-dev/

➢User documentation of the version in development:
- https://opensilex.github.io/phis-docs-community/
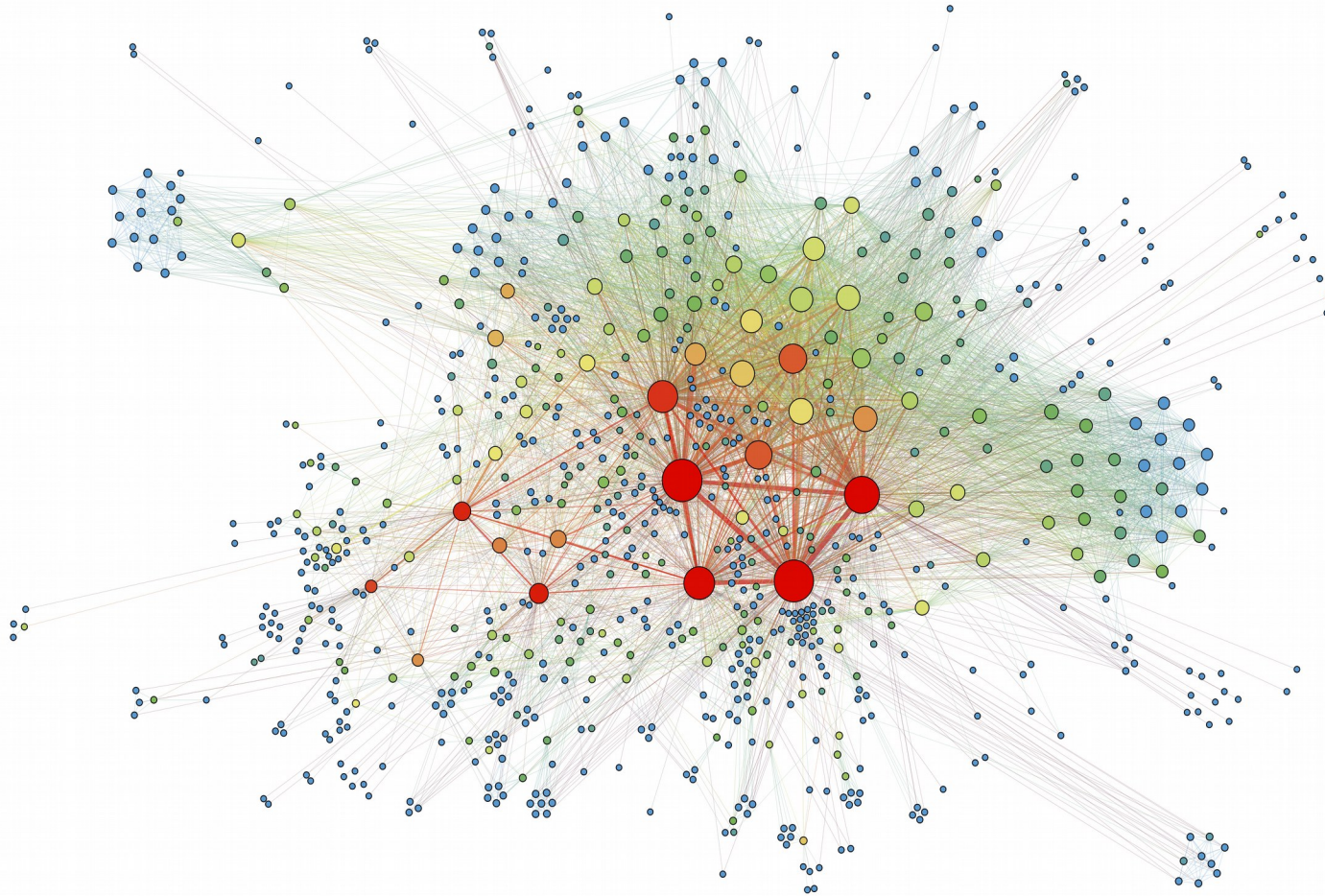
# I use spreadsheet! What I can do?



**you can provide extra information known as metadata about CSV files using a JSON metadata file**

2018ObservationSite.csv **then call the metadata file**
2018ObservationSite.csv-metadata.json

**More information:** http://w3c.github.io/csvw/primer/

- Merge duplicate resource IDs

- Reuse local unique

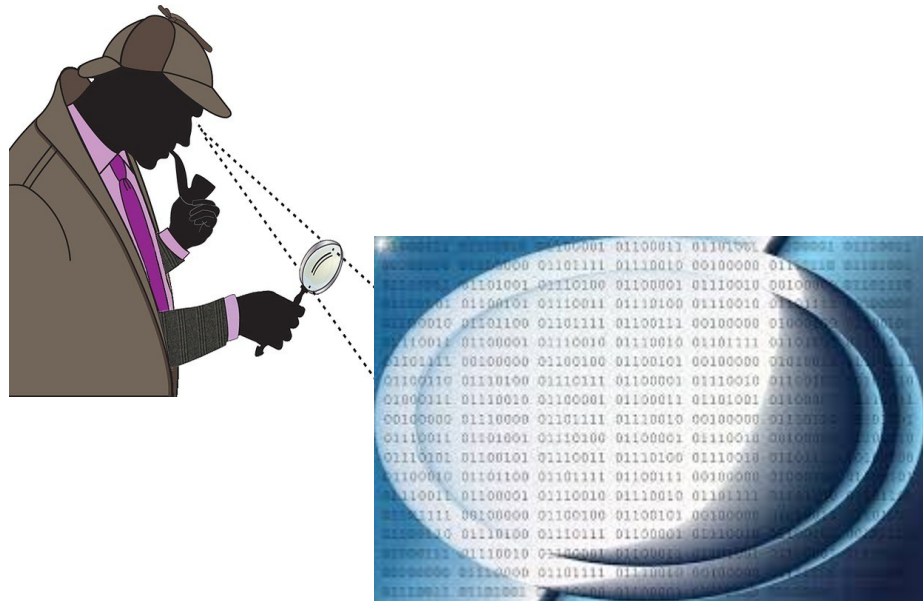- Reuse local unique

# DATA in a Semantic GRAPH

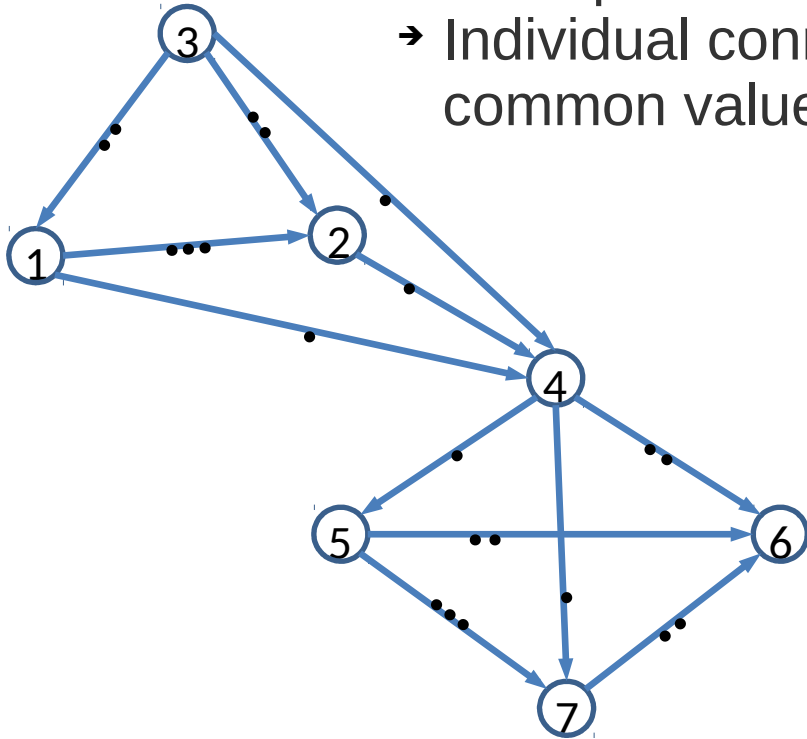# Analysis of Graph Structured Data

How to structure data ?

How to analyse data ?

# Data analytics – Overview

Example: Set of variables
➜ Individual connections depend on the number of common values

# Data analytics – Overview

Objective: find combinations of variables (RDF properties) in heterogeneous large dataset that discriminate a resource.

- Key : a set of properties (variables) that uniquely identify individuals

- Idea: Use key Interpret numerical data in a symbolic way

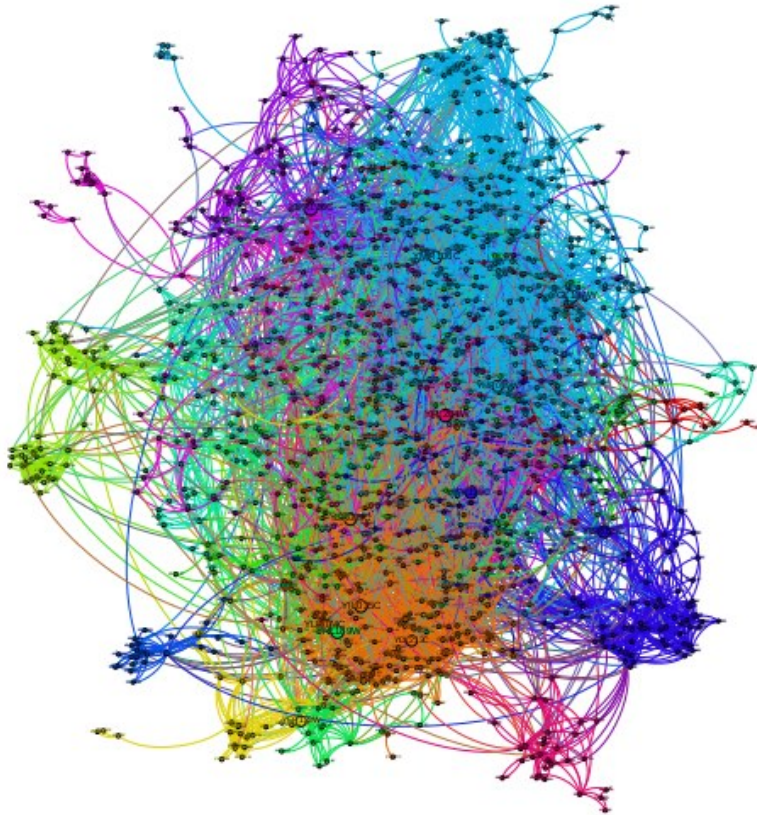- Automatically discover keys and evaluate their quality

- ***How do we discriminate the wines??***

| Quantiles | |
|---|---|
| | PH |
| Wine1 | 3.15 |
| Wine2 | 3.22 |
| Wine3 | 3.23 |
| Wine4 | 3.24 |
| Wine5 | 3.56 |
| Wine6 | 3.68 |

# Data analytics – Overview

## Uncovering Latent structures in networks



Interaction network

Discovering features in large network
Accounting for meta-informations :

Examples :
- Multiples networks
- Time series of networks
- Networks in space
-

**Model-based Statistical Model**:
Stochastic Block Model

**Idea** : Modelling the connection as a function of meta-information and some unknown latent structure.

*Mariadassou et al. ('10) Ann. App. Stat.*
*Barbillon et al.('16) JRSSA*
*ECONET ANR-18-CE02-0010.*

# Conclusion

- ✔ **Go to FAIR data Not Only for machine: allows teams and communities a better formalization of data**

- ✔ **Complex Big Data must be structured in graph in order to be able to FAIR, integrate and process it.**

- ✔ **Data mining and statistics methods are under development or already available to address volume and complexity issues**

- ✔ **Even if traditional methods remain accessible, the use of graph analysis can offer new ways**

# Thank you for your attention