# Research Objects and ROHub -
# A journey from theory to practical infrastructure

**Jose Manuel Gomez-Perez (Expert System)**
Raul Palma (PSNC)

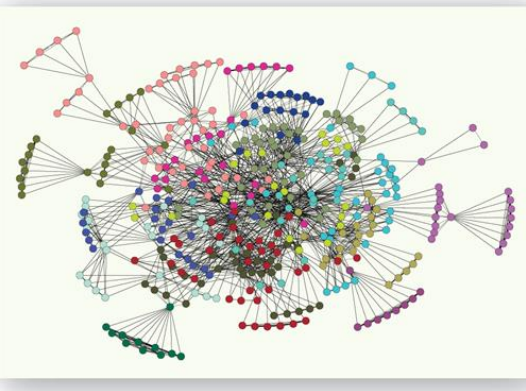EGU – 12th April 2018

# The Scientific Enterprise

**Collaborate**
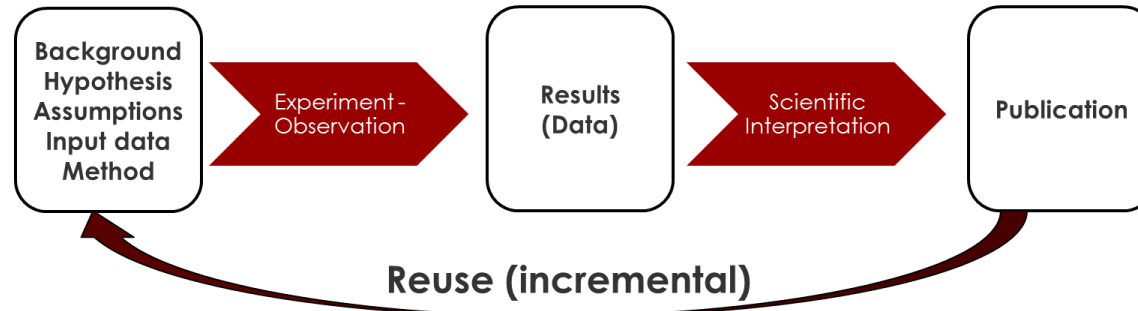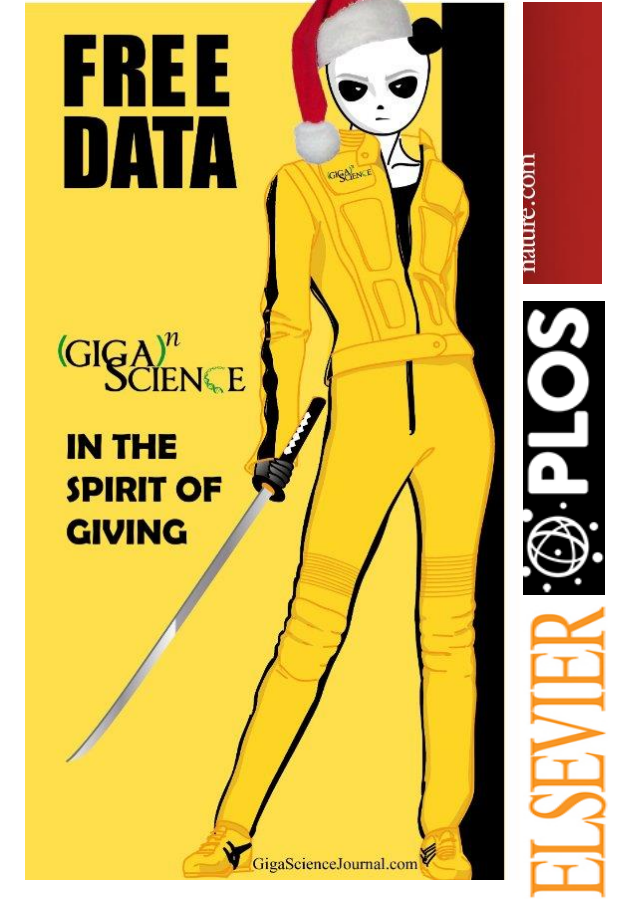
**Validate**

**Communicate**

Many landmark findings in preclinical oncology research are not reproducible, in part because of inadequate cell lines and animal models.

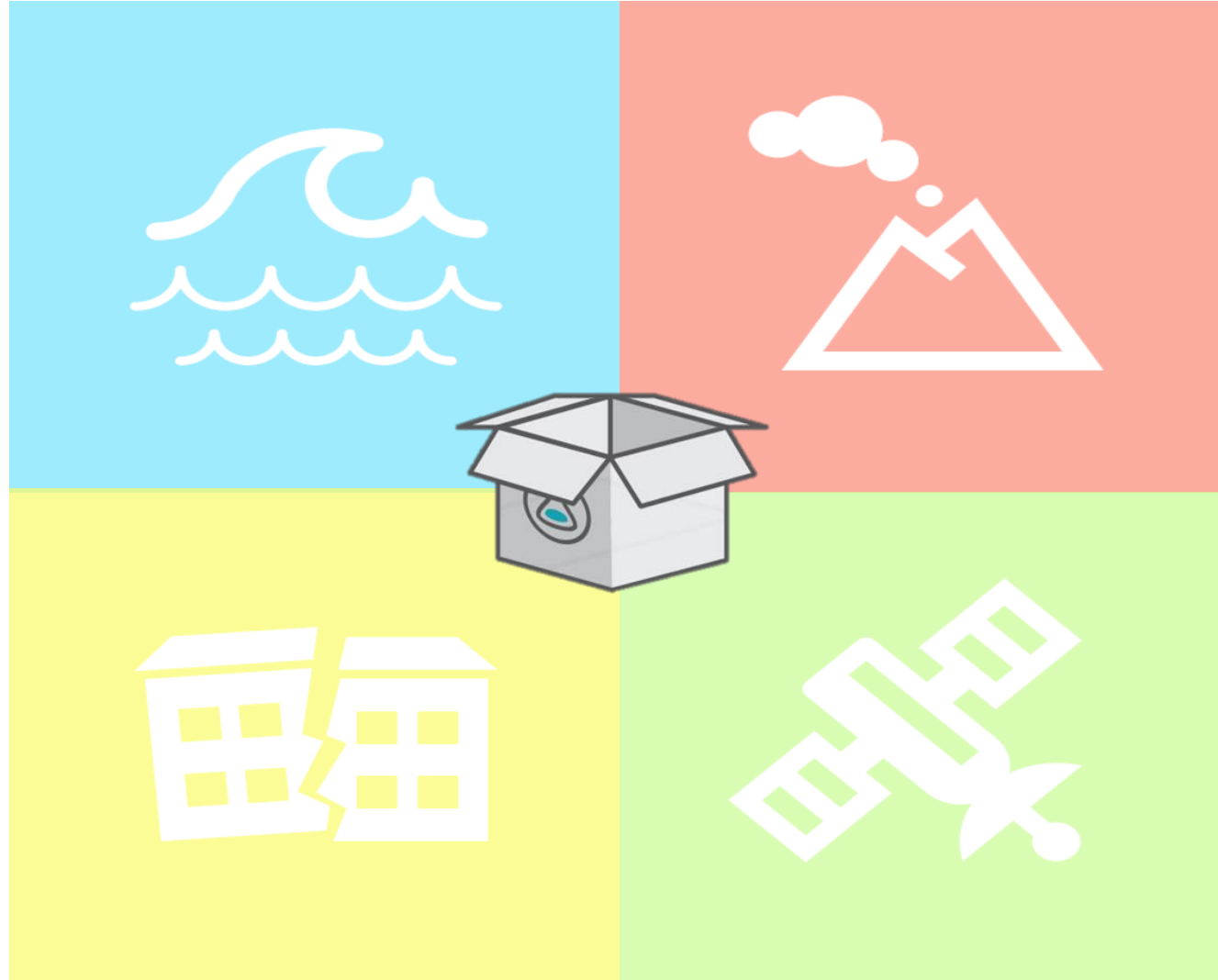## Raise standards for preclinical cancer research

| Background Hypothesis Assumptions Input data Method | → Experiment - Observation | Results (Data) | → Scientific Interpretation | Publication |

**Reuse (incremental)**

FREE DATA

(GIGA)ⁿ SCIENCE

IN THE SPIRIT OF GIVING

GigaScienceJournal.com

nature.com

PLOS

ELSEVIER

This project is co-funded by the European Union

# Research Objects in Earth Sciences

# Challenges

- **Long-term preservation**
  - Earth observation missions can cover a long timespan (+30 years)
  - Both data and models need to be preserved (e.g. harvesting of Bathymetry data for Sea Monitoring)
  - Long publication and documentation cycles

- **Sharing & attribution**
  - Reluctance of individual organisations and/or scientists to provide access to their data, methods and tools (IP issues, lack of time or resources, sensitivity of the resources involved, professional rivalry and competitiveness…)
  - Lack of data/methods citation mechanisms to give credit to and incentivize authors to share

- **Automation**
  - Long tail of software and computational resources
  - Limited adoption of scientific workflows for data orchestration and claim validation
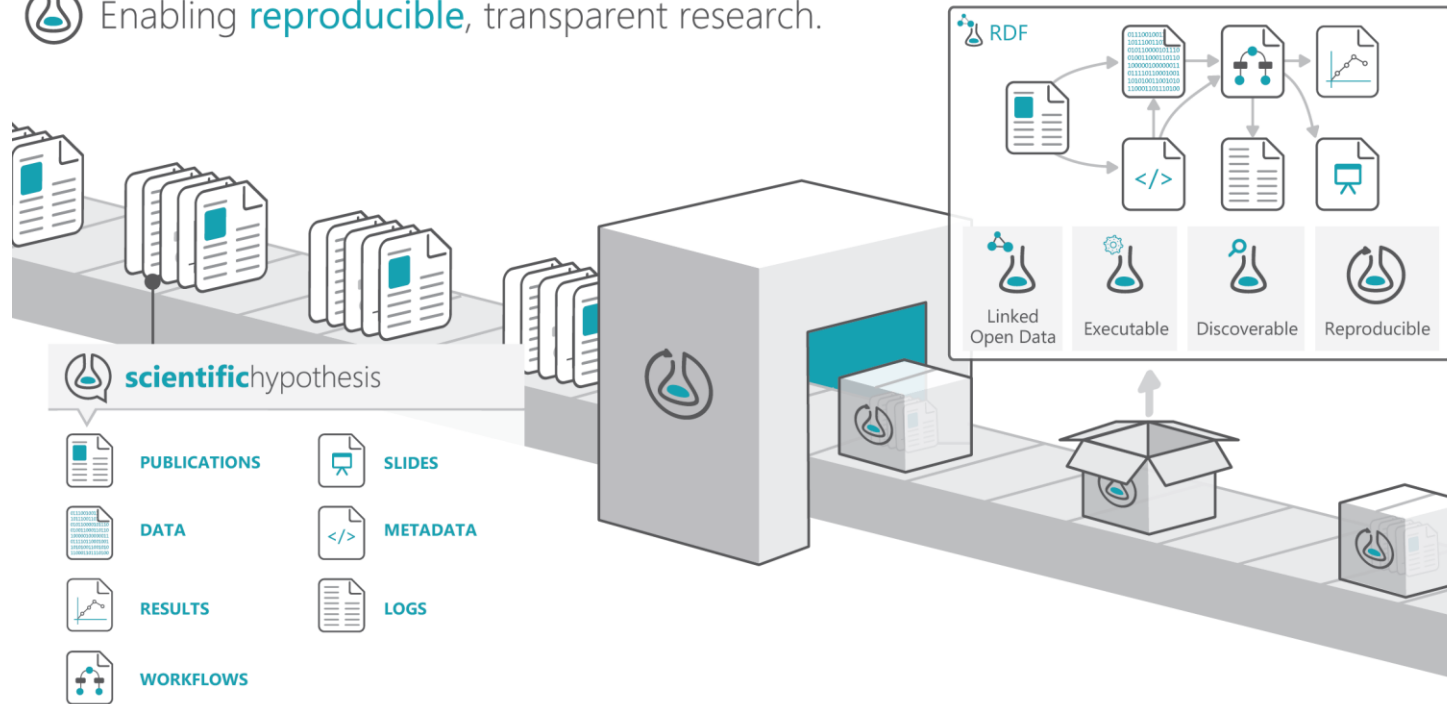
# Research Object – The Concept

A **Research Object** is an information artefact that contains and describes *everything* about your research, including how those things are related, in ways that are readable both by humans and machines



Enabling reproducible, transparent research.

scientific hypothesis

PUBLICATIONS  SLIDES

DATA  METADATA

RESULTS  LOGS

WORKFLOWS

RDF

Linked Open Data  Executable  Discoverable  Reproducible

http://www.researchobject.org

I. **Logically organize and describe in a single information unit** the resources, materials, methods and outcomes of an investigation

II. **Share** your research materials with other scientists at **discrete milestones of your investigation**. Uniquely identified

III. Enable **reproducibility** and **reuse** of scientific methods

IV. To be **recognized** and **cited**

V. **Preserve** results and **prevent decay**

VI. Provide **evidence** to findings claimed in **scholarly articles**

# Modeling Research Objects

- **Aggregation** (OAI-ORE) plus **annotation** (Annotation Ontology)

- Other vocabularies used in annotation bodies to provide information about resources, involving types, dependencies, and descriptions



**ro** (aggregation and annotation)
**wfdesc** (workflow description)
**wfprov** (workflow provenance)
**roevo** (evolution model)
**minim** (minimum information model)

- **Geospatial information**
- **Time-period coverage**
- **Data access policies**
- **Intellectual Property Rights**
- **General metadata (discipline, size, format, date…)**



- **Eight new types of research object, including:**
  - Workflow-centric, but also process and service-centric
  - Data-centric
  - Research product-centric
  - Documentation
  - Bibliographic



Available at: https://github.com/wf4ever/ro/tree/earth-science

# ROHub.org – The RO Management Platform



- Comprises both
  - Backend service (RODL) and API
  - A reference web client application (ROHub portal)

- Features include
  - Creation and preservation
  - Lifecycle and version mgmnt.
  - Change tracking
  - Quality and decay monitoring
  - Search, explore, reuse (fork)
  - Follow, subscribe and notifications
  - Likes and ratings

Palma R, Hołubowicz P, Corcho O, Gomez-Perez JM. ROHub - A Digital Library of Research Objects Supporting Scientists Towards Reproducible Science. In Presutti et al. (eds) Semantic Web Evaluation Challenge. SemWebEval 2014, Springer.

# Community Building and Content

**G**olden **E**xemplar **R**esearch **O**bjects



**Deep Sea Habitat Suitabilty Model**

In this research object we derive the MSFD indicator 1.5 (Habitat area) to assess the biological diversity descriptor. To do this in deep sea environment, the scientist (user) needs to implement a habitat suitability model.

**The Citizen science and jellyfish distribution**

A crowdsourcing app sponsored by Italian magazine and other different media provides scientific data to study jellyfish. CNR-ISMAR wants to fully exploit within the EVER-EST initiative the potential of the app to generate meaningful indicators in MSF perspective.

**Trend in the evolution of invasive jellyfish distribution**

Starting from Jellyfish sightings, we elaborate data to produce explicit geographical information concerning trends about the evolution and distribution of alien species according with MSF directive descriptors.

**Hazard Impact Model Development**

Research object to facilitate development of surface water flooding early warning systems and their impacts within the UK.

**Land monitoring Golden Exemplar**

Research object for the ingestion of satellite images acquired on land areas (with the support of information coming

**Volcano Source Modelling (VSM) - Application to Campi Flegrei (Italy)**

**IPWV map generation**

This research object contains the workflow which allows obtaining an integrated map of the precipitable water content over the Etna supersite, by using

**2013 Mount Etna Eruption (bibliographic Search)**

This is a bibliographic research object, which supports search of bibliographic

http://everest.expertsystemlab.com/home/#Golden%20Exemplars

- **20+ GEROs produced by the EVER-EST communities and collaboration with the USA National Ecological Observatory Network (NEON) and UNAVCO**



ROHub currently stores **3.099 research objects**, aggregating **84.593 resources** and **23.644 annotations**

**A**utomatically generated **B**ibliographic **R**esearch **O**bjects

- **~700 ABROs currently produced**
- **Semantically annotated**
- **Include grey literature, field reports, heterogeneous operational information…**



**INGV Reports**

List of daily and monthly reports by the Istituto Nazionale di Geofisica e Vulcanologia. 202 automatically generated Research Objects.

Learn more

**CNR Bibliographic References**

List of bibliographic references by the Consiglio Nazionale delle Ricerche and the Instituto di Scienze Marine. 209 automatically generated Research Objects.
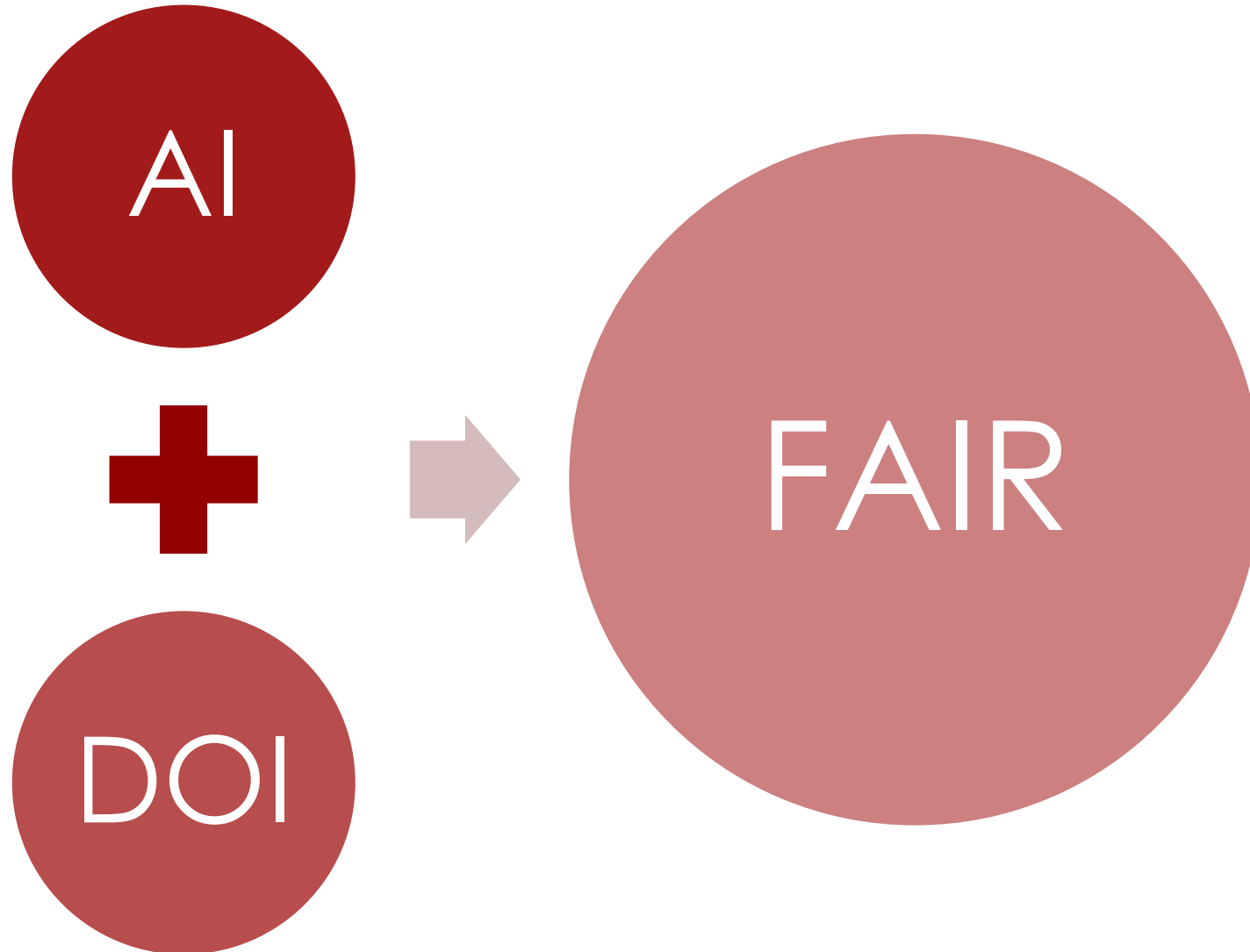
Learn more

**NHP Assessments**

List of daily hazard assessments by the Natural Hazard Partnership. 92 automatically generated Research Objects.

Learn more

# Two recent and relevant developments
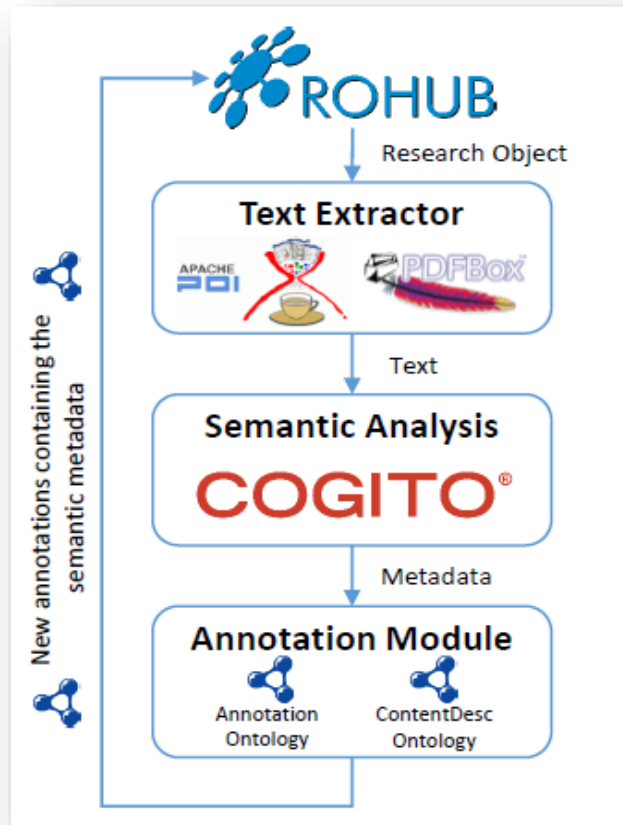
# The Metadata Chasm

- **Towards FAIR-ness, metadata is key**
  - **For scientists** (*"would this research object fit my investigation, as a whole or partially?"*)
  - **For machines,** through machine-readable annotations by search engines or recommendation systems
  - **For both:** To answer scientific questions

- **Research object metadata usually generated manually** (labor-intensive and scarce)

- **Metadata focused on lifecycle, structure and resource types rather than actual payload -** valuable knowledge sources like scientific papers, field notes or technical reports ignored

- **Related information hidden** and non-actionable for machine discovery, search or reasoning

- **Limited diffusion and reuse** of scientific outcomes

This project is co-funded by the European Union

Manual inspection of 2,500 research objects showed only a **third have a proper title**, with average length of 38 chars. Also, **usually short (138 chars) and non-descriptive descriptions**
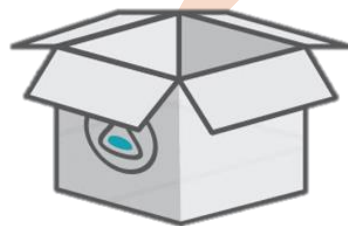
Where is the content metadata?

# Crossing the Chasm – Semantically Enriching Research Objects



- **Automated Semantic Annotation Natural Language Processing**

- **COGITO**, standard version w/o earth science extensions
  - Main entities include concepts, lemmas (canonical representation of a word) and relations (properties, hypernymy, polysemy, synonymy…)
  - ~300K syncons, ~400K lemmas, 80+ relation types (~2.8 million links)
  - Supports word-sense disambiguation based on word context

- Semantic enrichment annotates **the most significant concepts, domains, lemmas, noun phrases and named entities** in research object resources *(titles, descriptions, papers, bibliography…)*

Gomez-Perez, J. M., Palma, R., & Garcia-Silva, A. (2017). **Towards a Human-Machine Scientific Partnership Based on Semantically Rich Research Objects**. In *2017 IEEE 13th International Conference on e-Science (eScience)* (pp. 266–275). IEEE. https://doi.org/10.1109/eScience.2017.40

# Richer, Machine-Readable Metadata Enables Cool Apps...



Content

Other users

Me

Context of interest

Ranked recommendation

Information card

Automatically generated metadata

- Exploratory search and recommendation of research objects in scientific social networks

- Recommend by example

- Focused on the aggregated similarity of the context of interest with other items in the repository

- Based on metadata automatically generated from research object content and structure

- Extends findability and reusability

- Reduces cognitive load exploring scientific repositories

Rico M, Gomez-Perez JM, González R, Garrido A, Corcho O. (2017). Collaboration Spheres A Visual Metaphor to Share and Reuse Research Objects. *arXiv preprints https://arxiv.org/abs/1710.05604v1*

# DOI for Citation and Reuse

# The Research Object Journey



- Artificial Intelligence, NLP, Computer Vision
- Cross-modal content annotation (text, images, diagrams, tables, provenance…)
- **Towards a Digital Aristotle that can reason and answer science questions**

Content

Community

Platform

Model

Concept

# The research object of this talk



http://www.rohub.org/rodetails/EGU18_VRE_session_keynote-1/

# Don't publish. Release… often!

# Some key facts

- Number of Research Objects: 2,500+ (starting point 1,200)
    - Golden Exemplars: 2 or 3 per VRC
    - Generated automatically: 511

- Number of users: 162

- Activity events: 275K+
    - Last week > 47K ( approx. 6K per day )

- ROHUB storage size
    - solr:
        - named_objects: 868 / 912KB
        - notifications: 275620 / 193MB
        - ros: 2648 / 15MB
        - ros-private: 146 / 1.1MB
    - File system: 2.4GB
    - Tripple store: 1.1 GB
    - Database: 300MB

- ROHub access (previous week):

|  | all | success | ROs | resources |
|---|---|---|---|---|
| sum | 469081 | 466506 | 101935 | 367146 |
| avg | 58635 | 58313 | 12741 | 45893 |

*As of September 2017