

Input-Output Data Science of Paris-Saclay

The data platform

The linked data platform

The linked open data platform

Karima Rafes – BorderCloud
Karima.rafes@bordercloud.com

Cécile Germain – LRI, Université Paris Sud, CNRS, INRIA
cecile.germain@lri.fr

What must to have the Data Scientists in 2017 ?

A **collaborative** Data Science IT Platform for :

- Sharing
- Linking
- Reproducing
- Starting to imagine the futures strategies for the management of Scientific Information

Design principles

- **Respectful, volunteer**: centralization or intrusiveness is excluded
 - Paris-Saclay is cross section of institutional science: massive heterogeneity in scale, organizations, data policies, resources,...
 - CDS is a scientific collaboration, not an IT service
 - We « build on top » of existing practices and repositories
- **Collaborative, agile**: Scientific knowledge is a continuous process
 - Empower grassroots scientists with the tools to build **and re-build** the domain vision
 - Following mainstream approaches and standards

Data Science IT Platform, today

Input-Output Data Science of Paris-Saclay

<https://io.datascience-paris-saclay.fr>

For the friends
IODS

- Sharing : **Map** the existing data repositories operated by CDS partners
 - Better knowledge of the available resources.
 - Helping collaborations to get off the ground.
- Linking : **Interoperable** data
 - Provide tools, usable across disciplines
 - Build a Linked (Open) Data **platform** to publish these.
 - Linking data between domains: **collective knowledge building**
- Infrastructure for **reproducibility**.
 - Ability to test algorithms on new or extended datasets
 - Provide tools to reuse the queries of interoperable data

CDS Datasets



- 34 datasets from CDS visible
- Objective for CDS2: be reasonably representative of the datasets in the CDS2 perimeter
 - Proprietary or sensitive excluded for now
 - Remember: you keep your data in-house
- Objectives of the demo
 - Iods basics: how to publish a dataset - much easier than you imagine
 - A few clicks
 - Any format, eg excel, text, community, sql, nosql, rdf...
 - Iods advanced: linking data
- Help available: karima.rafes@gmail.com

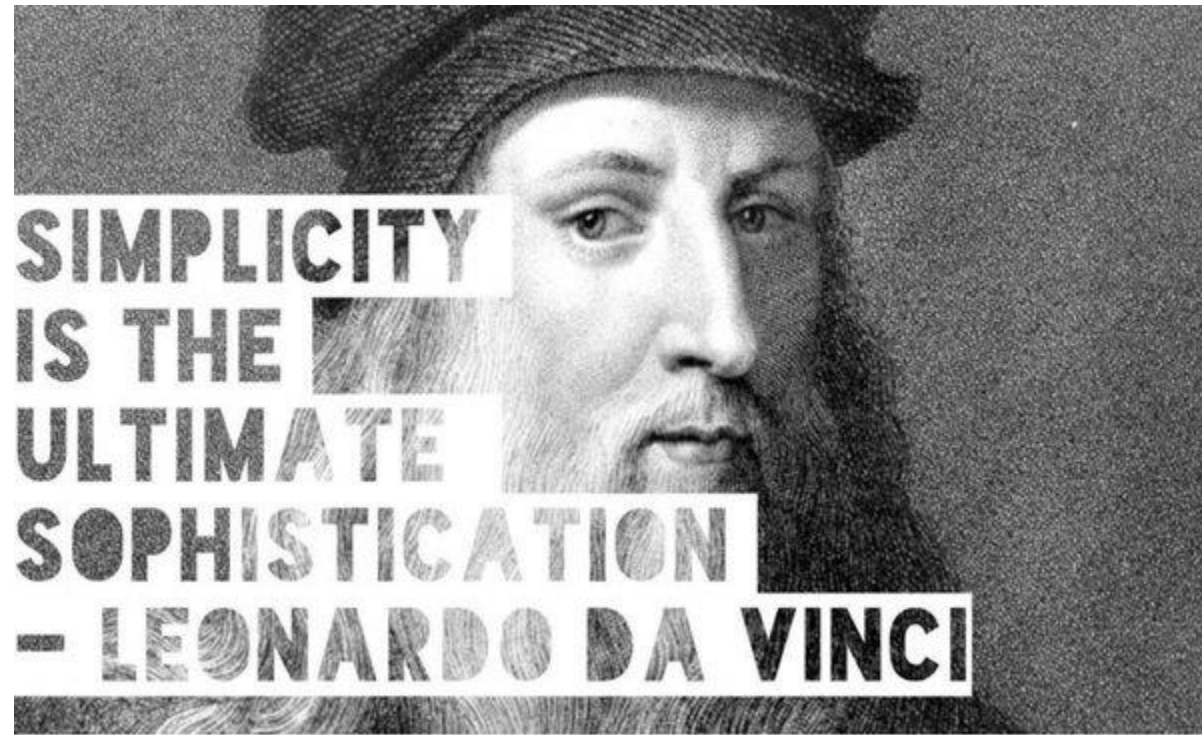
CDS Datasets



Poll

<https://goo.gl/forms/OuxdalXlphg0b6fM2>

- Link also on the indico page
- Please fill it by noon



DEMO
IODS BASICS

Demo : Sharing any data

- Unique identifier (& free) for each dataset
- Different access level :
 - public, share (for partner of CDS) or private
- Description :
 - What : Card, Licence, size, localisation, audience, etc
 - Distributions : Files (excel, image,...) & API
 - Exemples : queries
 - Links with the knowledge bases

Wikidata

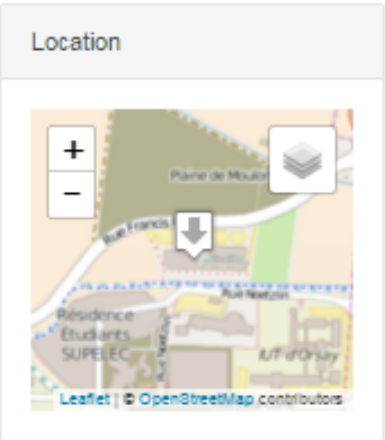
[general knowledge](#) [Semantic Web](#) [Wikipedia](#) [Wikidata](#) [Wikivoyage](#) [Wikisource](#)

Description

Wikidata is a free linked database for the structured data of its Wikimedia sister projects including Wikipedia, Wikivoyage, Wikisource and others.

The content is available under a free license, exported using standard formats, and can be interlinked to other open data sets on the Internet data base.

<https://www.wikidata.org> karima.rafes@gmail.com



Links to Wikidata.
Scientific topics (light blue)
Tag about the content (dark blue)

Distribution

SPARQL endpoint [CC0](#) [srx](#) [json](#) [oqv](#) 2018-01-31 [Details](#) [Q](#) [Go](#)

RDF Exports [CC0](#) [nt](#) [gz](#) [json](#) 2018-01-31 [Details](#) [Go](#)

Description

Data distribution

Publishers by Occupation [The content is not displayed?](#)

[science](#) [computer science](#) [scientist](#) [profession](#)

This query retrieves the number of scientists publishing in Wikidata.

5



Examples to reuse

Star rating about the capacity to reuse these data



SPARQL endpoint



Search in the examples

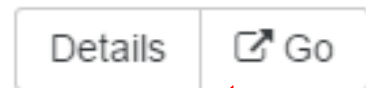
Make a new query in IODS

Open the official query form

See details



RDF Exports



License, formats and date

Access type :
WWW Web site
Torrent
File
DATA (API)

Access to raw data

Multi-criteria search

Domains ☰

- 1 aerospace engi...
- 1 branch of scie...
- 4 computer science
- 1 cosmology
- 2 ecology
- 1 economics
- 3 economy
- 1 environmental ...
- 2 geography

Tag ☰

- 1 asteroid
- 1 automobile
- 1 digital library
- 1 European Union
- 2 female
- 1 file format
- 1 France

32 results (0.07 seconds)

Aires protégées par taille

[écologie](#) [aire protégée](#)

Protected areas sorted by size, whatever the unit (km², mi², ha, ac)

Asteroids discoverers by countries

[cosmology](#) [asteroid](#)

Asteroids discoverers by countries

Automobile manufacturers by country

[economy](#) [automobile](#)

Piechart of automobile manufacturers by country.

Cemetery of Père-Lachaise of Paris - year of death

[Wikidata](#) [Wikidata](#)

Cimetière du Père-Lachaise - année de décès

Personnalités enterrées au cimetière du Père-Lachaise (hors columbarium) par année de décès

Demo : Linking the knowledge

With the linked data platform :

You facilitate the discovering of your data.

1. Describe a dataset in IODS
2. Install CDS' APPs in Wikipedia
3. Search a topic in your preferred search engine...
4. Open the Wikipedia's page

➔ Discover the available data and more !



Raman spectroscopy

(Q862228 [ttl](#) [rdf/xml](#) [json](#) [sparql](#) [Export](#) : [labels](#))

Datasets in relation

- Drug Classification** : The dataset contain Raman spectra of 4 types of chemotherapeutic agents diluted in 9 different solutions, and having different concentrations. Measures were made by the Lip(Sys)². ([source](#))

Devices in relation



RamanEvolution - Organization : Lip(Sys)²

Inria teams in relation

[Hyperspectral imaging](#)

[Mutual exclusion](#)



DEMO

IODS ADVANCED

Demo : Test and reuse

With the linked data platform :

You facilitate the reusing of your data.

1. Choose an example
2. Copy
3. Change the query
4. Save it
5. Copy/Past the code in your algorithm/program/Wiki

➔ Reuse the Linked Data is simple.

Linked Data example : table, graphs...

Asteroids discoverers by countries

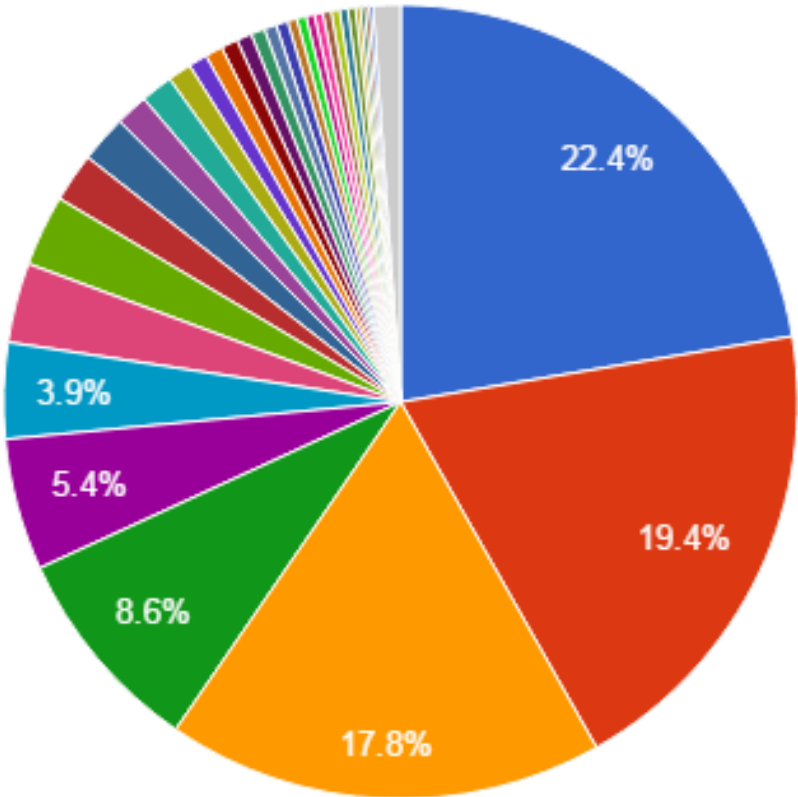
cosmology asteroid

Asteroids discoverers by countries

Links to Wikidata.
Scientific topics (light blue)
Tag about the content (dark blue)

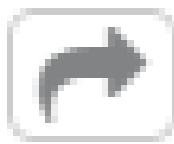
Description

Tools



- United States of America
 - Netherlands
 - Japan
 - Belgium
 - Germany
 - Soviet Union
 - Italy
 - Czech Republic
 - Australia
 - Russia
 - France
 - Croatia
 - Ukraine
 - Canada
- ▲ 1/3 ▼

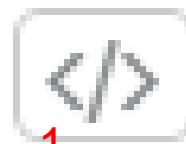
Table ou graph



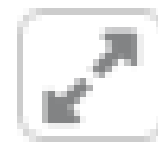
Share



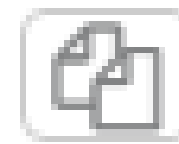
Catch



Reuse



Full
screen



Do a copy

Reuse these data in your code

Query, endpoint and code for reusing the same data

SPARQL

Javascript

HTML

Matlab

Python

R

PHP

Wiki

Choose your language

```
from SPARQLWrapper import SPARQLWrapper, JSON
```

Copy

3 Copy/past the code

```
sparql = SPARQLWrapper("https://query.wikidata.org/bigdata/namespace/wdq/sparql")
```

```
sparql.setQuery("""
```

```
    PREFIX bd: <http://www.bigdata.com/rdf#>
```

```
    PREFIX wikibase: <http://wikiba.se/ontology#>
```

```
    PREFIX wd: <http://www.wikidata.org/entity/>
```

```
    PREFIX wdt: <http://www.wikidata.org/prop/direct/>
```

```
select ?countryLabel (COUNT(?asteroid) as ?nb)
```

```
where {
```

```
    ?asteroid wdt:P31 wd:Q3863 ;
```

```
            wdt:P61 ?discoverer .
```

```
    ?discoverer wdt:P27 ?country .
```

```
    SERVICE wikibase:label {
```

```
        bd:serviceParam wikibase:language "en" .
```

```
    } .
```

```
}
```

```
GROUP BY ?countryLabel
```

```
ORDER BY DESC(?nb)"""
```

```
sparql.setReturnFormat(JSON)
```

```
results = sparql.query().convert()
```

```
for result in results["results"]["bindings"]:
```

```
    print(result["label"]["value"])
```

Howto use SPARQL with Python ?

4

May be also, read the doc ;-)



Copy and change a query

Copy of Asteroids discoverers by countries

Title

Domain

Tag

Short description

Dataset

Query

```
1 PREFIX bd: <http://www.bigdata.com/rdf#>
2 PREFIX wikibase: <http://wikiba.se/ontology#>
3 PREFIX wd: <http://www.wikidata.org/entity/>
4 PREFIX wdt: <http://www.wikidata.org/prop/direct/>
5
6 select ?countryLabel (COUNT(?asteroid) as ?nb)
7 where {
8     ?asteroid wdt:P31 wd:Q3863 ;
9     ?asteroid wdt:P31{instance of} wd:Q3863{asteroid} ;
10    ?asteroid wdt:P61 ?discoverer .
```

Copy of Asteroids discoverers by countries

Title

Domain

Tag

Short description

Dataset

Query 17 } .
18
19 GROUP BY ?countryLabel
20 ORDER BY DESC(?nb)"}
Press **Ctrl + Space** to activate auto completion.

Chart

Options

Style

Description

Select a SPARQL Endpoint

SPARQL Query

Test the query

Configure the graph/table

What next in 2018 ?

Connecting **experts** and **problems**

Connecting **experts** and **tools**

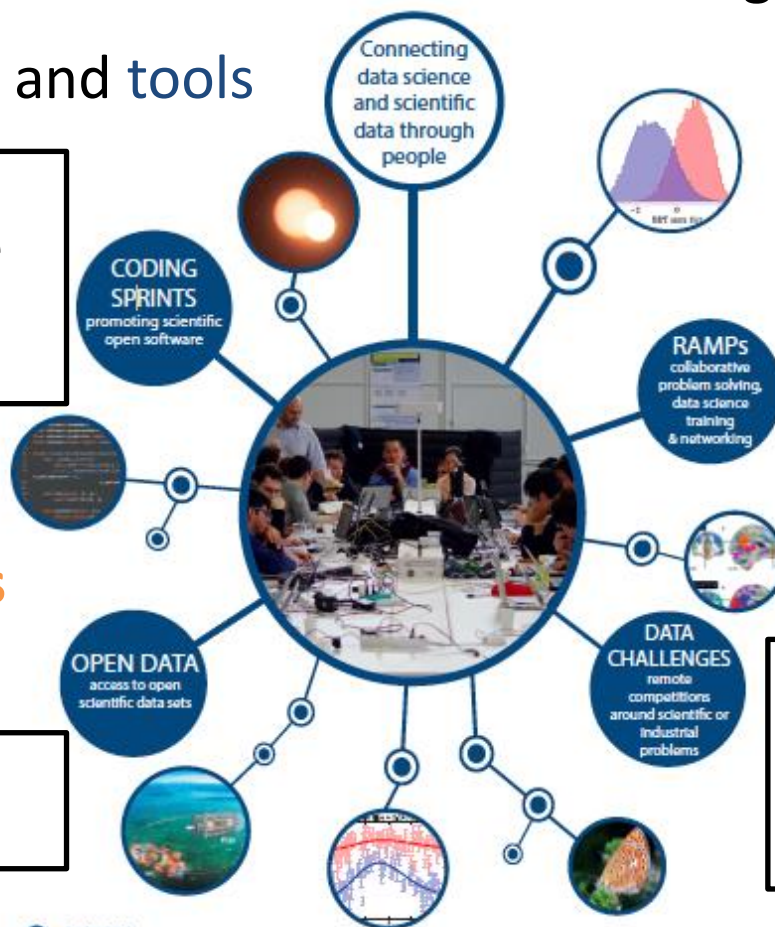
- State-of-the-art data science in easy-to-use tools
- High-quality software

- Prototyping
- Training
- Collaboration building

Connecting **experts** and **data**

- Data as a Service
- Linked (Open) Data

- Impact on science
- Visibility
- Benchmarks



@SadayCDS

Thanks, questions ?

Karima.rafes@bordercloud.com

Hangout : Karima.rafes@gmail.com

Twitter : @karima_rafes

« Long live the open science and proper »



Poll