



**HAL**  
open science

## Data dictionary cookbook for research data and software interoperability at global scale

Romain David, Laurent Bouveret, Lorraine Coché, Pedro Pizzigatti Corrêa, Rorie Edmunds, Ana Heredia, Jean-Luc Jung, Yasuhisa Kondo, Iwan Le Berre, Yvan Le Bras, et al.

### ► To cite this version:

Romain David, Laurent Bouveret, Lorraine Coché, Pedro Pizzigatti Corrêa, Rorie Edmunds, et al.. Data dictionary cookbook for research data and software interoperability at global scale. Research Data Alliance Plenary 17 (RDA P17), Apr 2021, Edinburg (virtual), United Kingdom. , Research Data Alliance Plenary 17 (RDA P17), Edinburg, remotely, 20-22 april 2021 (Session poster session), 2021, 10.5281/zenodo.4683066 . hal-03214743

**HAL Id: hal-03214743**

**<https://hal.science/hal-03214743v1>**

Submitted on 2 May 2021

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

We are now facing profound changes (biodiversity, climate, pandemic, etc.). Human impacts and their mitigation will depend on our ability to mobilize research at the global level. The sustainable development of the society will largely depend on the **sustainable development of global science and scientific research tools, outputs, and research ecosystems**. This globalization of research requires inter-operating our observation and experimentation systems in order to better understand these changes, to better simulate their effects.

The Covid-19 pandemic is now raging around the world. The reproducibility of research and results across regions in different contexts should accelerate human responses. Data sharing and the development of Synthesis Research with data aggregation at large scale is critical to enable such processes.

**The use of common knowledge, vocabularies, standards and procedures at a large scale is necessary.**

## Objectives:

This poster proposes a draft common methodology, a data dictionary cookbook, which will provide a roadmap towards the building of large scale data dictionaries.

The objective is to report on the challenges met while building data dictionaries in three global projects related to biodiversity and/or disease research:

**PARSEC, Kakila, ERINHA-Advance.**

# Data dictionary cookbook

## Generic

Define your community perimeter  
**[Scientific Questions]**

Explain and convince  
**[Data dictionary principles]**

Identify your scientific objects  
**[List of ENTITIES]**

Identify the aspects of your scientific objects you need to assess  
**[List of QUALITY]** for each Entity

Name and define necessary variables and their dimensions:  
**[Variables names: Entity\_Quality\_Dimensions]** for each Quality of each Entity

Reuse existing standards, vocabularies, concepts and definitions

Validate by the whole community  
**[List of DEFINITIONS]**

## Parsec

Community perimeter defined by:  
**[Scientific Questions]**  
 Can some satellite images be used as proxies for socio-economic indicators?

**[Data dictionary principles]**

- Choose common definitions to enable comparisons between countries
- Validate reusable definitions to ensure reproducibility of project results

About the scientific question, examples of  
**[ENTITIES]**

- Municipality
- Satellite images
- Sectors
- ...

For each entities, examples of  
**[QUALITIES]**

- Human Development Indicators of [Municipality or Satellite images or Sectors]
- Population of [Municipality or Satellite images or Sectors]
- ...

Due for each qualities, examples of  
**[Variables names: Entity\_Quality\_Dimensions]**

- Municipality\_HDV\_value
- Municipality\_population\_value
- ...

Reuse **existing vocabularies and ontologies**, with a challenge: when there are several, choosing the better terms.  
 All terms are **community approved**

## Kakila

Community perimeter defined by:  
**[Scientific Questions]**  
 Can we characterize cetacean presence in the waters of the Guadeloupe archipelago using several citizen science databases?

**[Data dictionary principles]**

- Choose common vocabulary to describe fields from heterogeneous databases
- Adopt common variables description and values between databases

About the scientific question, examples of  
**[ENTITIES]**

- Cetacean observation
- Observers
- ...

For each entities, examples of  
**[QUALITIES]**

- Individual Species (Cetacean observation)
- Individual Length (Cetacean observation)
- Experience (of Observer)
- ...

For each qualities, examples of  
**[Variables names: Entity\_Quality\_Dimensions(&unit)]**

- Cetacean\_Obs\_Length
- Observer\_Exp\_[category]
- ...

The two challenges are:

- to obtain a **consensus for variables values** between databases
- to **align with the existing data** to Darwin Core vocabulary

All terms are **community approved**

## Erinha Advance

Community perimeter defined by:  
**[Scientific Questions]**  
 Can Biosafety level 4 laboratories compare their virologic *in vivo* studies results?

**[Data dictionary principles]**

- List all *in vivo* models, viruses, protocols and their data sensitivity
- List and homogenize variables description for possible meta-analyses

About the scientific question, examples of  
**[ENTITIES]**

- Virus
- Hamster
- ...

For each entities, examples of  
**[QUALITIES]**

- Mortality rate (of the virus)
- Age (of the hamster)
- Breed (of the hamster)
- ...

For each qualities, examples of  
**[Variables names: Entity\_Quality\_Dimensions]**

- Virus\_Mortality\_Percentage
- Hamster\_Age\_Nbweeks
- Hamster\_Breed\_[BreedName]
- ...

The two challenge are:

- to Reuse **existing vocabularies and ontologies** (if several, choosing the better one).
- to obtain a consensus between variables names and definitions.

All terms must be **community approved**

Data dictionary with **[List of DEFINITIONS] approved by the whole community** contain ALL **[Scientific Questions], [List of ENTITIES], [List of QUALITIES], [Variables names: Entity\_Quality\_Dimensions\_Units]** and **variables definitions**

## 3 projects, the same cookbook

The **PARSEC Project** is building new tools for data sharing and reuse through a transnational investigation of the socioeconomic impact of protected areas.

The **Kakila database** centralizes and harmonizes marine mammal observation data for the AGOA sanctuary around the French archipelago of Guadeloupe, French Antilles.

The **ERINHA-Advance project** aims to support the operations of the ERINHA research infrastructure which is designed to generate data from transnational access research activities on highly pathogenic agents.

In these 3 global case-studies, similar challenges have arisen: **to aggregate and interoperate pre-existing heterogeneous data** at the global scale, and **to share common tools** to monitor, maintain quality, scan scale and cope with uncertainty.

## Dealing with complexity

- Developing data dictionary literacy is an essential work to imply scientists in variables definitions,
- Enabling adaptability, portability, replicability and reproducibility implicates that software and workflows must be defined as all data in the data dictionary,
- Addressing dimension issues in each context is necessary for all variables.



## As simple as possible!!

The common experience of our three projects showed that we **need to proceed step by step as simply as possible** and to ensure that each step is **understandable for the whole community**. It is necessary to *improve access and re-use of all existing semantic materials* and not trying to build a cathedral with a little spoon.

### \* References:

- ★ David, R., Mabile, L., Specht, A., Strycek, S., Thomsen, M., Yahia, M., Jonquet, C., Dollé, L., Jacob, D., Bailo, D., Bravo, E., Gachet, S., Gunderman, H., Hollebecq, J.-E., Ioannidis, V., Le Bras, Y., Lerigoleur, E., Cambon-Thomsen, A. and Research Data Alliance – SHARing Reward and Credit (SHARC) Interest Group, T.R.D., 2020. FAIRness Literacy: The Achilles' Heel of Applying FAIR Principles. Data Science Journal, 19(1), p.32. DOI: <http://doi.org/10.5334/dsj-2020-032>
- ★ Coché, L., Arnaud E., Bouveret L., David R., Foulquier E., Gandilhon N., Jeannesson E., Le Bras Y., Lerigoleur E., Lopez P., Madon B., Sananikone J., Sèbe M., Le Berre I., Jung J.-L., 2021. Kakila database: Towards a FAIR community approved database of cetacean presence in the waters of the Guadeloupe archipelago based on citizen science. Biodiversity Data Journal: Data paper. submitted Dataset: <https://doi.org/10.48502/cg6n-1103>



### Acknowledgements:

PARSEC is funded by the Belmont Forum through the National Science Foundation (NSF), The São Paulo Research Foundation (FAPESP), the French National Research Agency (ANR), and the Japan Science and Technology Agency (JST). ERINHA Advance is funded by ERINHA-Advance european program under grant agreement N°824061. Kakila database is funded by the LabEx DRIIHM French program "Investissements d'Avenir" (ANR-11-LABX-0010) and supported by the SO-DRIIHM project (ANR-19-DATA-0022). This work is partially funded by the EOSC-Life European program (grant agreement No. 824087).