

Session 4 International and open access strategy for the future

A trip to the “data management” wonderland

Cécile Pignol - Engineer

Sedimentary Platform & Archives collections

*Data management group (LTER-Fr -BED) Bancariser Ensemble les Données
(Give access to Data together)*

13:30-16:00

Fondamental in Data management - F.A.I.R. Principles

Data management (Samples & data registering strategies)

Summary :

Party I : Introduction & State of Data management situation ----- Slide 3 -> 11

Issues & benefits of data sharing

« Long tail of Data »

Process for collecting scientific data and publications

Link for Self-training courses

Party II : Fundamental in Data management ----- Slide 12 -> 28

Data life cycle

What kind of Data ?

Principes FAIR fundamentals

Party III : Praticals ----- Slide 29 -> 56

Usecase1 for Sample : Collec of cores

Usecase2 about Analytical data

Return to Principes FAIR (Specific)

Data-repository

Party IV : You have some 1st keys to make a ----- Slide 57 -> 61

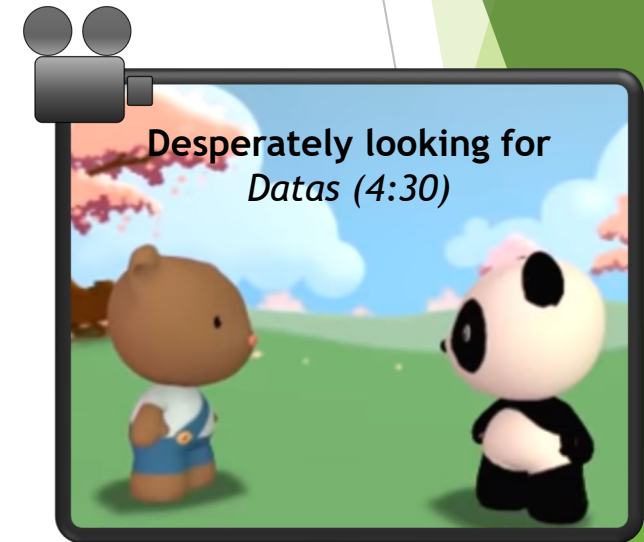
DMP (Data management plan)

Data Paper

Conclusion : CHECK-LIST----- Slide 62 -> 67

PART I- INTRODUCTION

- 1- Issues & beneficts of data sharing
- 2- Process for collecting scientific data and publications
- 3- What are the brakes in 2021 ?
- 4- Learning path / Self-training courses



BENEFICITS OF DATA SHARING



FOR RESEARCHERS

TRANSCRIPTION OF NEW REQUIREMENT & LAWS

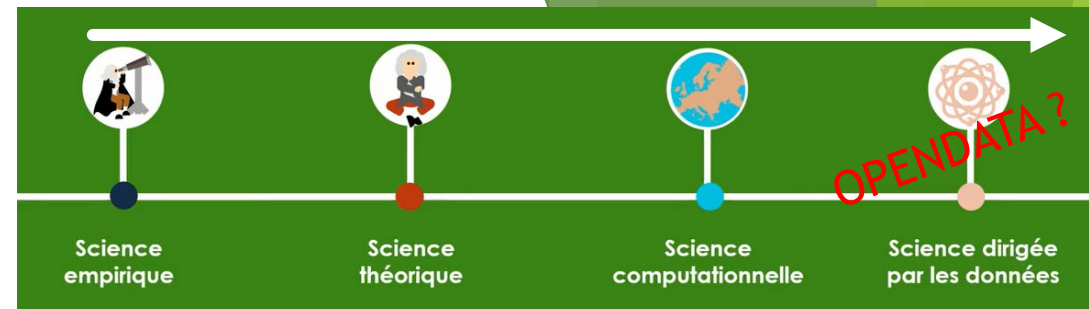
Research data from **publicly funded** must be **openly** disseminated

→ Européen **H2020 Projets (DataMmangementPlan)**,

→ In France, the **Plan national pour la science ouverte** (some of which aim to structure and open research data)
Research Agency **ANR (DMP)**

OPPORTUNITIES

- increase the **visibility & re-use of the work** of researchers



FOR SCIENCE

NEW PERSPECTIVES

Data sharing allows researchers to build on a larger data foundation and paves the way for **new methods of investigation**.

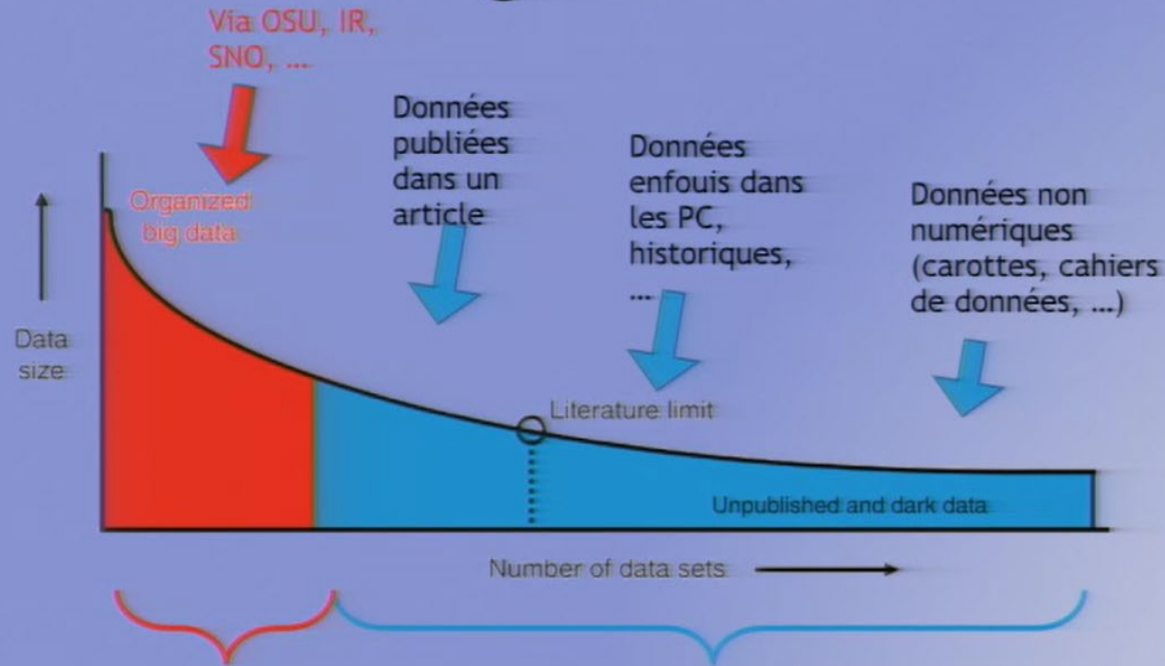
EXAMPLE: Use reference databases, approach the same dataset in different ways...

STONGER VALIDATION - QUALITY - REPRODUCIBILITY

- conditions for **validating articles** by access to **data**
- Filing of **data** is increasingly requested by the **reading committees**.
- Re-use of the data by other researchers may lead to **reconsideration of the initial results**.

« LONG TAIL OF DATA »

Données de longue traine

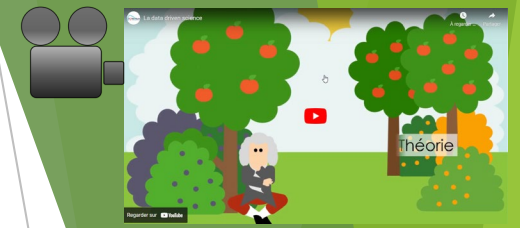


Données structurées

Les grands efforts scientifiques ne représentent qu'une petite proportion du total des données collectées par les scientifiques

Données de longue traine

De nombreux petits efforts de recherche indépendants produisant une riche variété de jeux de données de recherche spécialisés. Non structurés, souvent inaccessibles au monde extérieur (et intérieur)



Accessible structured data (Big data, Environmental observatories, Ecology, etc.)

Vs

Unstructured data, +/- available for some time, to researchers around a publication or, available on request

Non-Digital Data (Samples, Paper Lab Book,...)

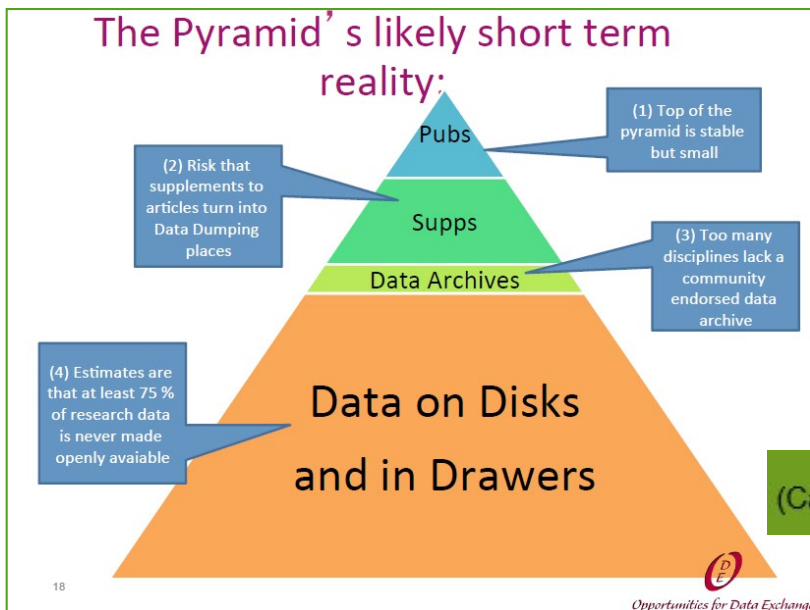
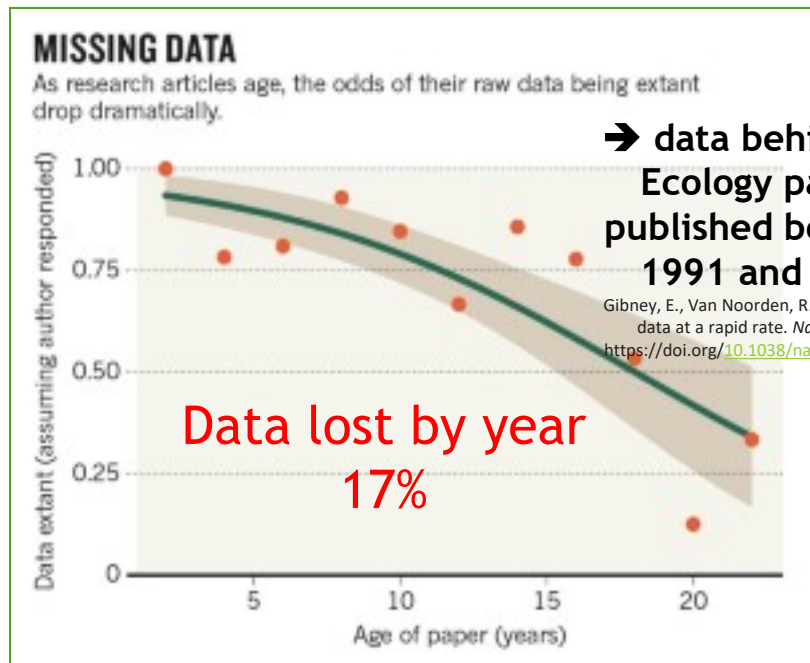
→ Talk about Smalls & Darks Data !

Vidéo 2:30 :

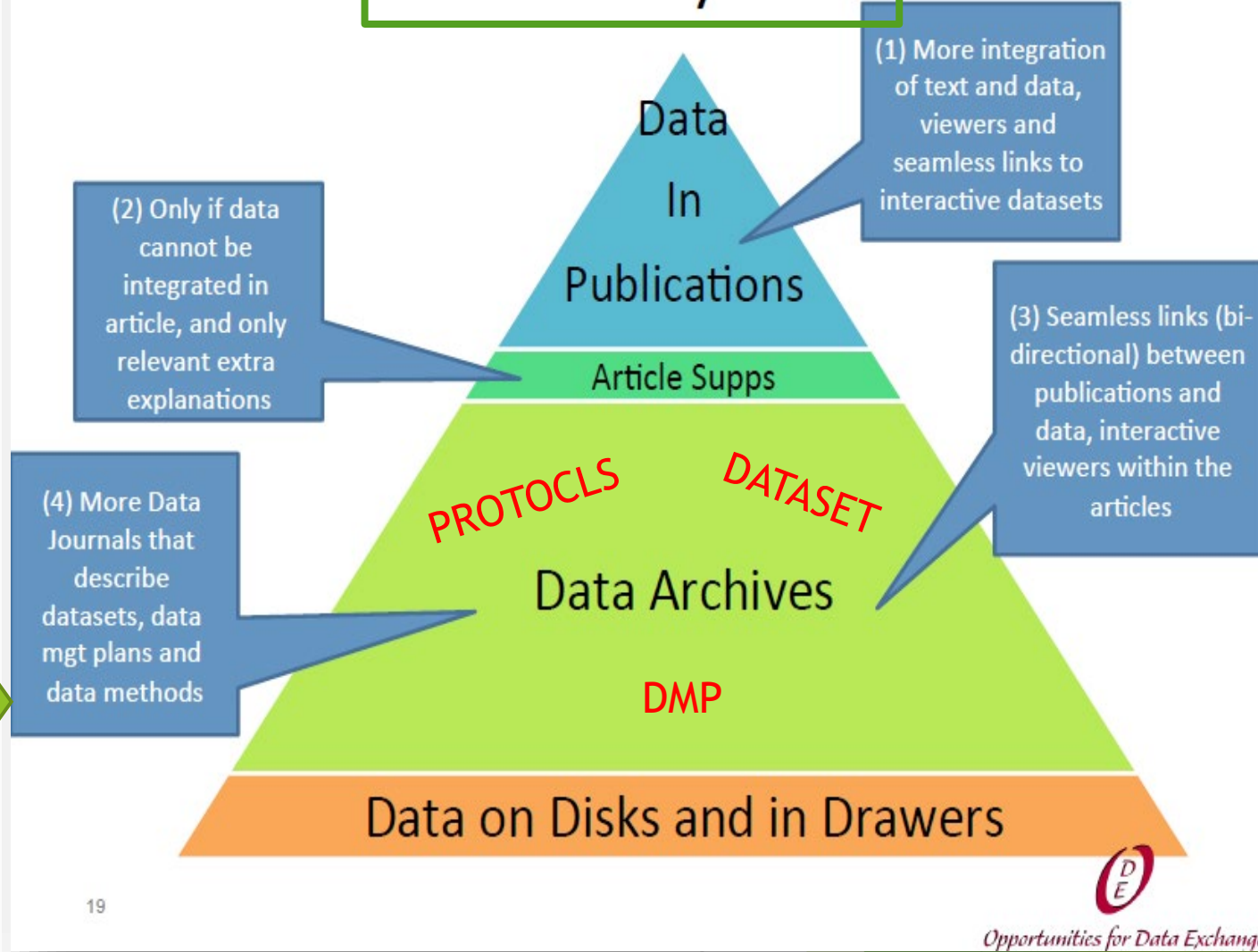
Video (in french) à 2:30 : S Galle, V Stoll, Données de longue traine, accès aux données attachées aux publications, données hors SNO, données non numériques, Atelier défi 14 de la prospective transverse INSU, Du 20 janv. 2020 au 21 janvier 2020

Ferguson, A., Nielson, J., Cragin, M. *et al.* Big data from small data: data-sharing in the 'long tail' of neuroscience. *Nat Neurosci* **17**, 1442–1447 (2014). <https://doi.org/10.1038/nn.3838>

The reality !



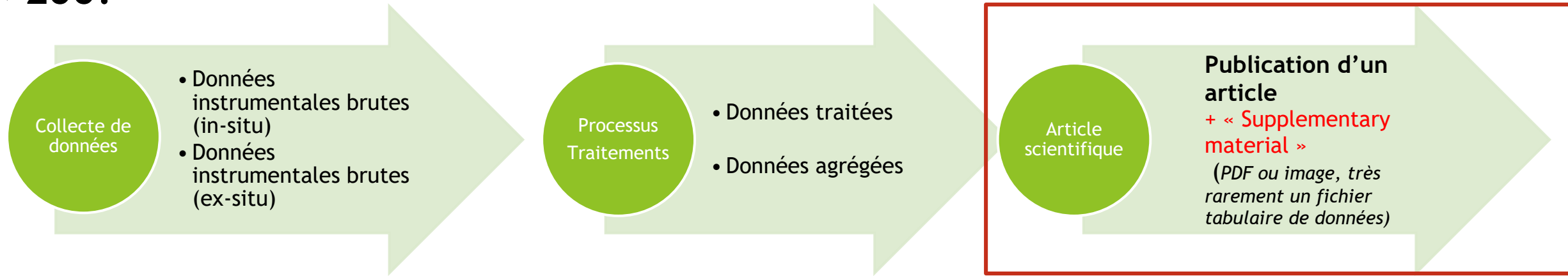
The Ideal Pyramid



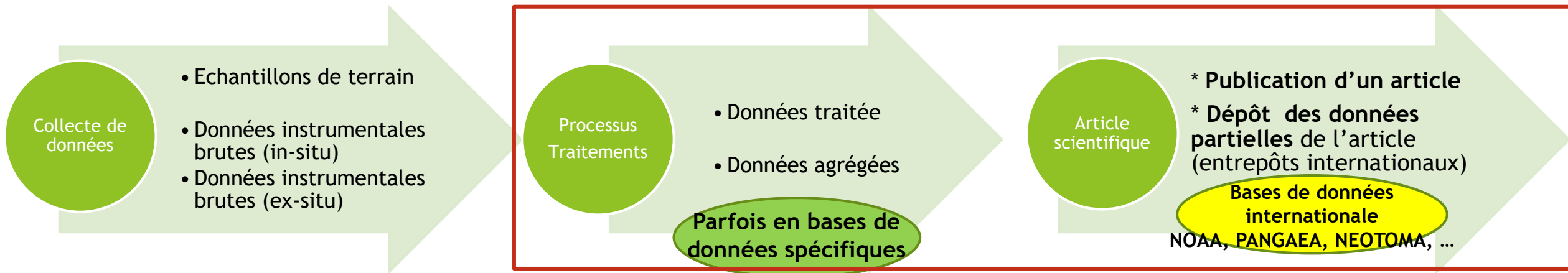
(Callaghan, 2013) d'après : (Reilly, Schallier, Schimpf, Smit, & Wilkinson, 2011)

Process for collecting scientific data and publications (1)

< 200?



200?-2020



IS THERE ANY OBSTACLES IN 2021 ?

▶ Less incentive from local institutions ?

→ **No**, country build nationals' roadmap

▶ Lack of appropriate IT infrastructure ?

(where to put your data, safe place and easily accessible)

→ **No !** but **Yes** we still need easy & ergonomic IT & Curation procedures to describe & deposit data

▶ Lack of time ? managing your data is a time-consuming activity ?

→ **True !** especially if it taken into account **too late**

→ Need to well organize / to anticipate / to create or optimize procedures ?

▶ Lack of skills ?

→ **True !** data management and sharing require various skills in IT, data curation, legal knowledge, etc.)

→ **1- Increase skills and use easy tools to save time**

→ **2- Create teams : Researcher - Data-librarian - IT Profils**



<https://data.gov.be/fr/info-faq>

<https://www.snf.ch/fr/FAiWVH4WvpKvohw9/dossier/points-de-vue-politique-de-recherche>

YES

YES



Self-training courses

► In french : Big effort for training
DoraNUM (INST CNRS), MITI (CNRS), CoopIST (CIRAD),
➔ Guide, Fiches, Vidéos...

The image shows two overlapping website screenshots. The top one is DoraNum, titled 'DONNÉES DE LA RECHERCHE APPRENTISSAGE NUMÉRIQUE', with a navigation menu including 'ACCUEIL', 'RESSOURCES', 'ACTUALITES', 'À PROPOS', and 'CONTACT'. The bottom one is CoopIST, titled 'Gérer les données de la recherche', with a navigation menu including 'Trouver l'information', 'Être auteur', 'Rédiger', 'Publier et diffuser', 'Evaluer', and 'Gérer des données'. Both sites offer resources for researchers on data management.

sur-la-gestion-des-donnees-de-la-recherche/

The logo for 'Atelier Données' features a central circular icon with a gear and data points, surrounded by text: 'Atelier Données', 'Guide de bonnes pratiques', and 'MITI'. To the right, a list of data management topics is shown in a word cloud style: 'Métadonnées', 'Stockage & Archivage', 'Licence', 'Sécurité', 'Partage', 'Accessibilité', 'Interopérabilité', 'Normalisation', 'Documentation', 'Citation', 'FAIR', 'Reproductibilité', 'Transparence', 'Traçabilité', 'Généralisation', 'Reproduction', 'Métiers', 'Recherche', 'Généralisation', 'Reproduction', 'Métiers', 'Recherche'.

Atelier Données – groupe de travail inter réseaux de la MITI – CNRS
Guide de bonnes pratiques sur la gestion des données de la recherche
Rechercher dans ce livre ...
version 1.0 - Janvier 2021

➔ Parcours interactif adapté aux domaines de l'environnement



Gestion des données de recherche en Environnement

► In english :

The screenshot shows the MANTRA website interface. It features a header with the logo 'MANTRA Research Data Management Training' and a main heading: 'MANTRA is a free online course for those who manage digital data as part of their research project.' Below this, there are navigation tabs: 'Home', 'About', 'Acknowledgements', 'RDMS MOOC', and 'Feedback'. A section titled 'Learning Units: Select one to start' lists various topics like 'Research data in context', 'File formats & transformation', 'Protecting sensitive data', 'Data management planning', 'Documentation, metadata, citation', 'FAIR sharing and access', 'Organising data', 'Storage & security', and 'Data handling tutorials'. There are also icons for 'Research Student', 'Career Researcher', 'Senior Academic', and 'Professional'.

<https://mantra.ed.ac.uk/#profiles>

The screenshot shows the GFZ Data Services website. It features a header with the logo 'GFZ Data Services' and a main heading: 'Data Readiness Roadmap'. Below this, there is a section titled 'Data Readiness Roadmap' with a sub-heading: 'The Data Readiness Roadmap is a tool to help researchers assisting with the challenges of data management. Research data management (RDM) is a process which begins before the data is collected and in most cases should conclude with data being shared publicly. This process is termed the Data Life Cycle. According to each stage of the data life cycle, handling research data may require different RDM practices, although some are common to all.' There is also a section titled 'A Word on Data Management' with a sub-heading: 'Data represent a fundamental output of the scientific process. Managing data can be a daunting process with ever-increasing data quality, multiplying with each transformation, model iteration, or quality control steps. Therefore, organizing, storing, backing, documenting and sharing of research data are essential data management skills for a modern researcher. Successfully managing data provides many benefits in daily routines, but also helps to satisfy journal, institutional and funder requirements.'

The logo for 'GFZ Data Services' features a circular icon with a gear and data points, surrounded by text: 'GFZ Data Services', 'Data management recommendations for PhD students', and 'GFZ Data Services'.

Before continuing the rest of the presentation

Opendata concept (1 janv. 2020)

[... **data should be “as open as possible and as closed as necessary”** “open” in order to foster the reusability and to accelerate research, but at the same time they should be closed” to safeguard the privacy of the subjects*...]

according to the H2020 Program Guidelines on FAIR Data

<https://direct.mit.edu/dint/article/2/1-2/47/9998/The-A-of-FAIR-As-Open-as-Possible-as-Closed-as>

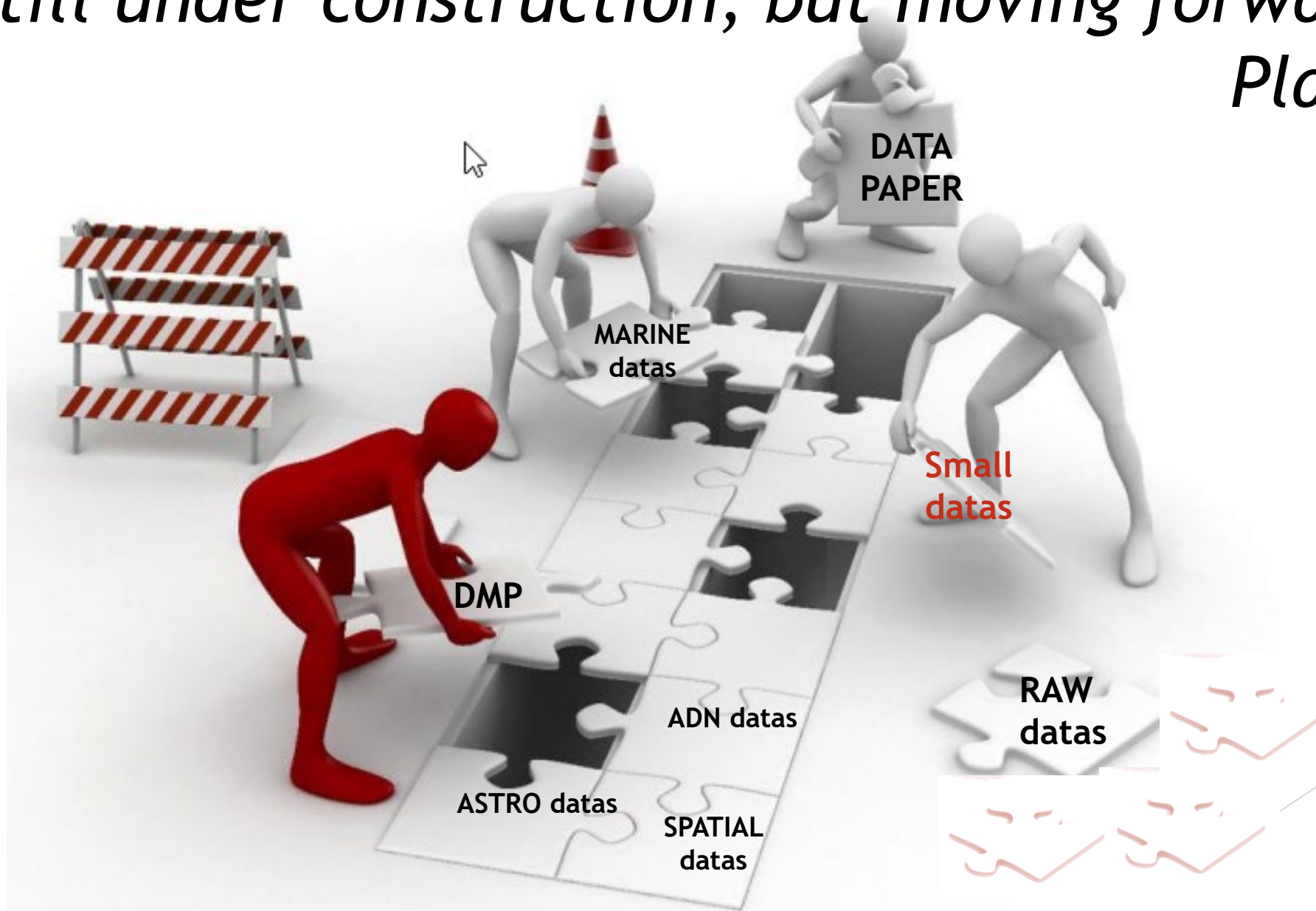
* persons, industrial contact, military information, ...

As Opendata is a way of work for the future decades ,
think that **you do what you want, within the limit of :**

- your funder injunction & laboratory or doctoral institution
- (your data manager need's because the institution tell him to do)
- the time you dedicated to this job : this (new) job isn't even an option now !

Still under construction, but moving forward !

Play together !



-PART II-

Fundamentals in Data management

- 1-Data Life cycle*
- 2-What kind of data ?*
- 3-F.A.I.R. Principles*

THE DATA LIFE CYCLE

6 STEPS

But **not necessarily** in the research timeline



Where
are
you
the
best?

- ★ Nouvelle recherche
- ★ Réexaminer les résultats/les données
- ★ Enseigner et apprendre
- Respecter la licence de réutilisation

- Définir les droits d'accès aux données :
- Niveau d'ouverture (de l'ouverture interne à l'ouverture publique)
 - Conditions d'utilisation des données : choix d'une licence de réutilisation

Plan de Gestion des Données
Localiser les données existantes
Collecter/acquérir des données ★
Créer les métadonnées

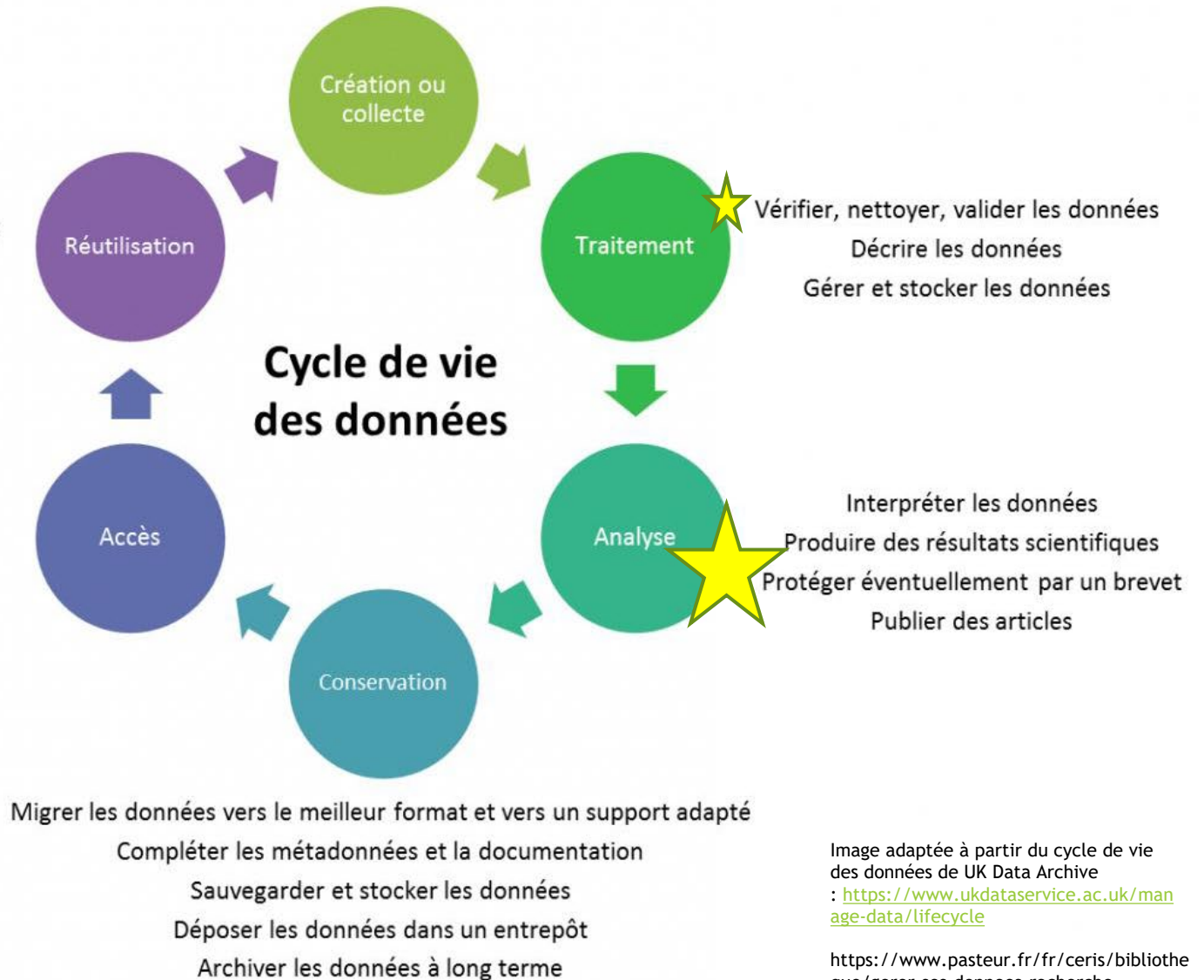


Image adaptée à partir du cycle de vie des données de UK Data Archive
: <https://www.ukdataservice.ac.uk/manage-data/lifecycle>

<https://www.pasteur.fr/fr/ceris/bibliotheque/gerer-ses-donnees-recherche>

What kind of data is involved ? (in the environment)

- ▶ identify and list the **types of data sets** that will be collected or generated (in the project/thesis, internship) :

Ex : Samples, experimental data, observation data, data from social surveys, simulation, software, images, space data *etc.*)

What Kind of informations do we need ?

- specify the **study object** (feature of interest), **geographical origin**, **time period**
- identify the **purpose** for collecting
- specify the **generation** of these different types of data (method, protocols, instruments, ...)

SOUS QUELLE FORME SE PRÉSENTENT CES DONNÉES ?

Il peut s'agir de :

+ OBJETS

+ DONNÉES TEXTUELLES

+ DONNÉES NUMÉRIQUES

+ DONNÉES AUDIOVISUELLES

+ MODÈLES , CODES INFORMATIQUES

+ DONNÉES SPÉCIFIQUES LIÉES À UNE DISCIPLINE

+ DONNÉES SPÉCIFIQUES PRODUITES PAR CERTAINS INSTRUMENTS

<https://doranum.fr/plan-gestion-donnees-dmp/origine-description-donnees-recherche/>

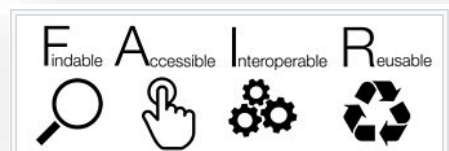
<https://coop-ist.cirad.fr/gerer-des-donnees/rediger-un-pgd/6-garantir-la-comprehension-et-l-accessibilite-des-donnees>

F.A.I.R. Principles

“it evokes a proactive and altruistic behaviour of the data producer, which seeks to make them easier to find and use by all, while facilitating downstream sourcing (possibly automatic) by the data user”

comportement proactif et altruiste du producteur de données, qui cherche à les rendre plus facilement trouvables et utilisables par tous, tout en facilitant en aval le sourçage (éventuellement automatique) par l'utilisateur des données

➔Technics & tools
enabling FAIRisation of
data



<https://www.go-fair.org/fair-principles/>



<https://fairaware.dans.knaw.nl/>



Australian FAIR
Tools diagnostic

<https://ardc.edu.au/resources/working-with-data/fair-data/fair-self-assessment-tool/>



Européen FAIR
Tools diagnostic

<https://www.f-ujl.net/index.php>

F.A.I.R. Principles : Who can applic & help for this ?

IT skills are everywhere in FAIR !

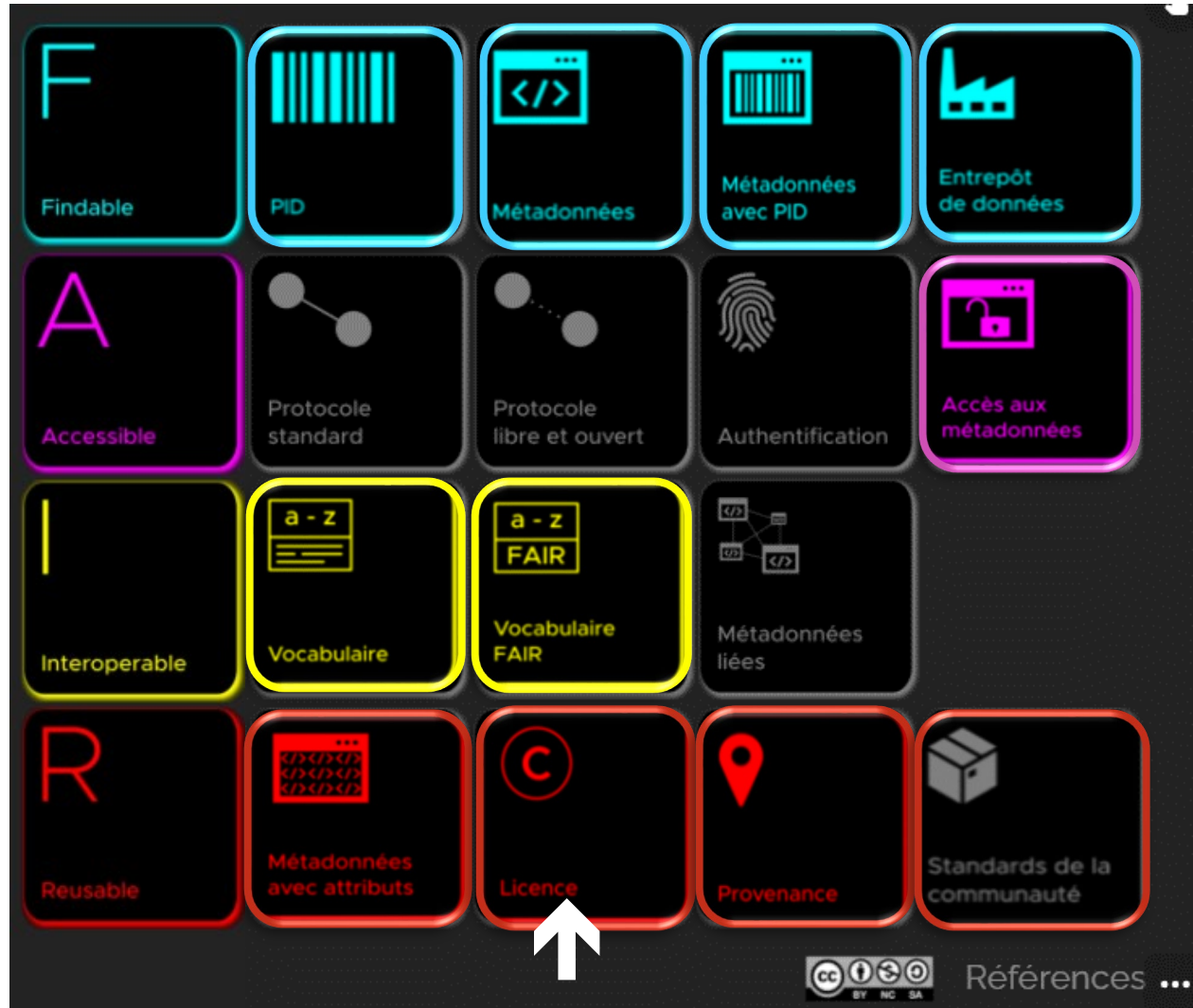
But one IT guy don't have all the time and skills

➔ **Make a team !**

Data librarian ➔

R=Researcher ? ➔

From you, for you !



← Data manager
← Data librarian

Data jurist

F.A.I.R Principles : 2 levels of management



level
« Discovery »

► Level 1 = Findable + Accessible

→ Don't expect a high quality of description

high level
of quality

► Level 2 = Level 1 + Interoperable & Reusable

beyond « Discovery »

→ Expect a high quality of description, precision for protocols, sensor, measurements, ...

Nb : you can't do a Data-paper only with this level !



Some basic F.A.I.R. principles on practices

To be FAIR, the data must be finely described using metadata.

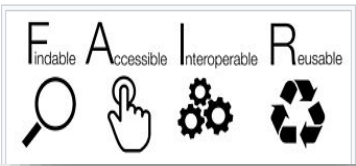
- Metadata automates sorting and prioritization tasks when **searching** for data.
- They also allow the **reuser** to better understand the **context of the data**, the conditions under which they were created or collected, their characteristics, etc.

The **more** information **you give** about your data, the **easier** they will be to **find** and **understand**

- LEVEL 1** → *Discovery Metadata (generalist or interdisciplinary) (for level 1)*
- LEVEL 2** → *Thematic Metadata (for level 2)*



The basic principles on practices



To be FAIR, the data must be described using a **controlled vocabulary** allowing interoperability.

By describing a data set according to a model of representation and a controlled vocabulary, you will provide resources understandable by both humans and machines

To be FAIR, metadata must **remain accessible in catalogues** even if the data is no longer accessible.

- Acces by (harvested) catalogues
- Over time, data may disappear. Metadata can be very useful in this case, as it will provide **valuable information about the missing dataset** and allow other researchers to take over and continue the associated research.



Focus on METADATA STANDARDS : WHY AND WHICH ?

Within a scientific framework, **metadata helps to describe the resources**, research data and productions made (article, repository, photo, measurement, software, web page, etc.)

- Facilitating the **search** for scientific research data and publications in catalogues (Level 1).
- Help to **understand the datasets** (collection, structure and context), for their validation (peers) and/or reuse (***Level 2***)



Les standards de métadonnées

Les métadonnées sont des **données décrivant d'autres données**. Elles doivent être **interopérables** afin de permettre à d'autres systèmes de les exploiter à la condition qu'elles **respectent certains standards**. Cette adéquation vous permet, ainsi qu'à d'autres utilisateurs, de faciliter la recherche **avec des critères précis et uniformes** entre des objets équivalents.

COMMENCER →

DORANum

<https://doranum.fr/metadonnees-standards-formats/standard-metadonnees/>



METADATA STANDARDS : WHICH ?



→ Catalogue of standards metadata & vocabulary

Level 1 : exemples interdisciplinary or generalist

- ▶ Dublin Core : interdisciplinary, description of digital resources.
- ▶ DataCite Metadata Schema : fairly generalist and can be useful to you regardless of your context (*Use for DOI PID*)

Level 2 : exemples thematics

- ▶ Iso19.115/139 : Used for geo-localized environmental data framework or geographic data
- ▶ EML (Ecological Metadata Language) : Used in Ecology
- ▶ Data Documentation Initiative (site officiel) : Used in Social science for social surveys

→ Some of them are « cross mapp » : ex Author = Creator

IGSN
Geological samples description)
PACTs
Standard for Paleoclimate Data
.../...



Controlled vocabularies, thesauri, ontologies = interoperability & quality increase



- ▶ **Vocabularies (for « Human reading ») : controlled vocabularies hierachic therms + **defintitions****

Level 1 : generalist // Level 2 : Specific & Thematic

- ▶ **Vocabularies *F.A.I.R* (for « marchine reading ») : vocabularies hierachic therms with defintitions + having resolvable globally unique and persistent identifiers (URI)**

Géographie thématique: Pollution
Pollution (en), Contaminación (es)
<http://data.loterre.fr/vocabs/BGI/GT/2440>

Géographie thématique: Pollution chimique
Chemical pollution (en), Contaminación química (es)
<http://data.loterre.fr/vocabs/BGI/GT/2442>

Géographie thématique: Pollution radioactive
Radioactive pollution (en), Contaminación radioactiva (es)
<http://data.loterre.fr/vocabs/BGI/GT/2447>



Term's catalogues (multi-thesaurii)

1-Catalog's vocabularies :



2- Search in multi-catalogs :

- **LOTERRRE** : <https://www.loterre.fr>
FR (Fr/Eng/Sp/It, Ger)

F.A.I.R.

- **Quantity & Units** : *Units of Measure, Quantity Kinds, Dimensions*

(Eng) → *the more complete*

- **NERC / NOC BODC Vocabularies** :

GB (Eng) → *Marine's Community (EU-SeaDataNet)*

- **CZA Critical Zone** :

US (Eng) → *Thematics [Observation Data Model v2 \(ODM2\)](#), [CUAHSI CVs](#), [IEDA EarthChem](#)*

- **PANGAEA « List » Vocabularies** :

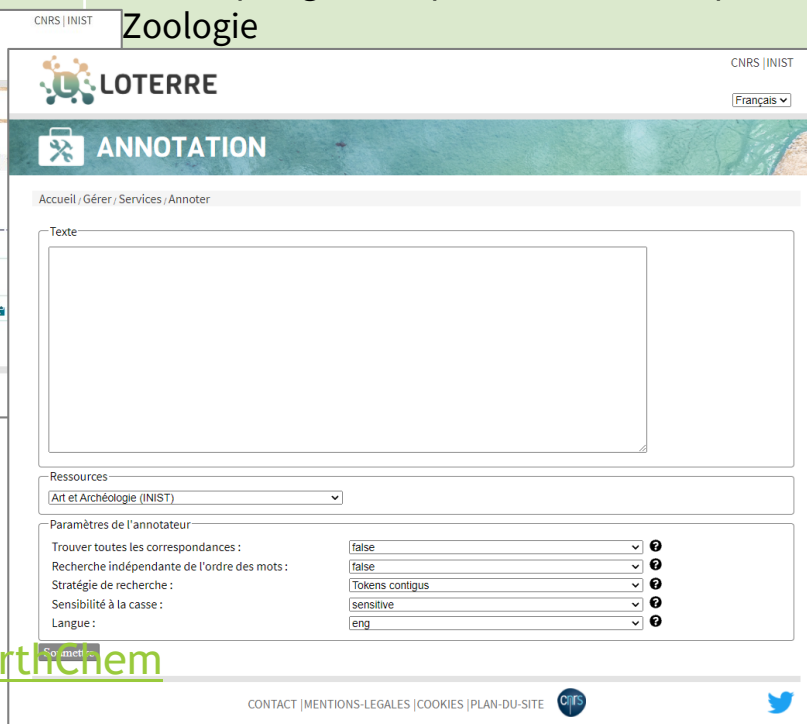
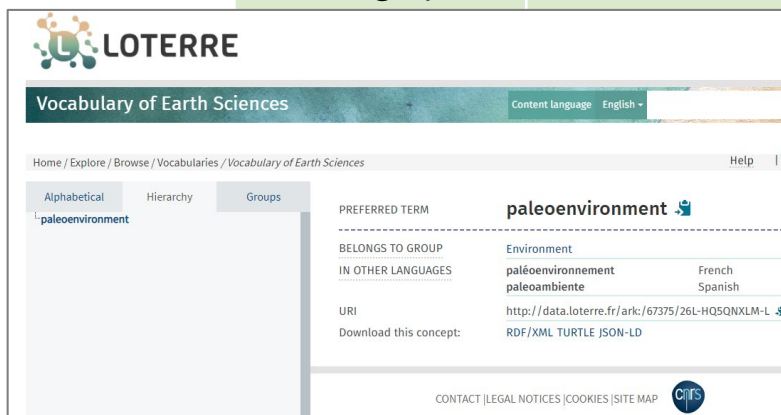


Parameter list

- [Complete list of parameters used in PANGAEA](#)



Sciences de la terre et de l'univers	Chimie	Sciences du vivant et environnement (non-exhaustif)
Aéronomie, Météorologie Climatologie Astronomie Géologie Géophysique Glaciologie Océanographie	Chimie analytique Chimie générale Chimie théorique Chimie physique Chimie minérale (inorganique) Chimie organique	Agronomie, écologie, environnement Biologie du développement végétal, botanique Évolution, écologie, biologie des populations Biotechnologies, sciences environnementales, agronomie, foresterie Biologie, santé Biologie moléculaire et structurale, biochimie Génétique, génomique, bio-informatique Zoologie





Term's catalogues (multi-thesaurii)

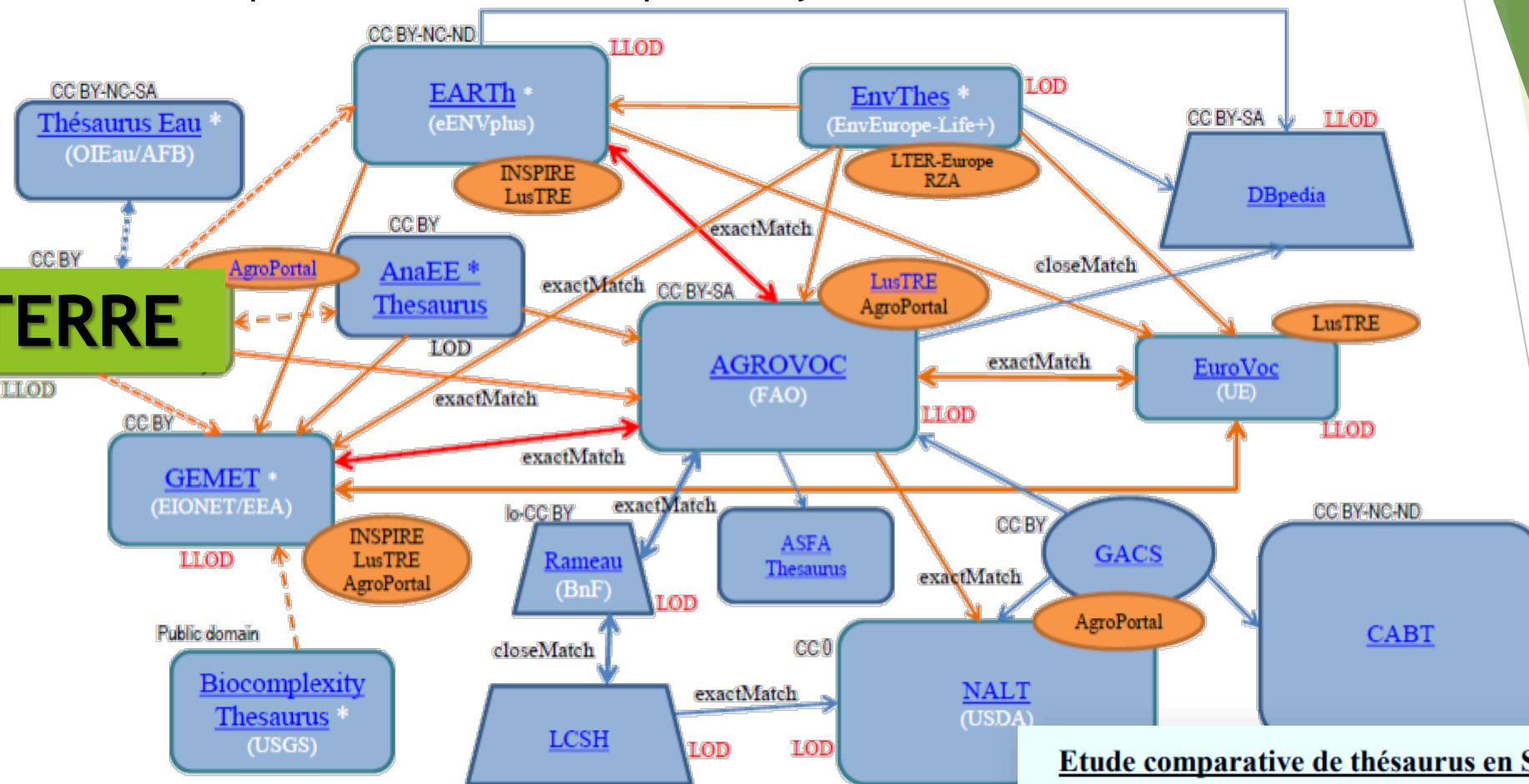
in 2018 exemple of thesaurii interoperability

In 2020

thesaurus
OZCAR-Theia

LOTERRÉ

Chimie
In 2020



Etude comparative de thésaurus en Sciences de l'Environnement
Bonnes pratiques de conception et FAIRisation de thésaurus

Institut de l'information scientifique et technique (Inist-CNRS) 2018- Dominique Vachez



HAL - Inist
Publications scientifiques de l'Inist

Dominique Vachez
dominique.vachez@inist.fr
INIST-CNRS, 2 rue Jean Zay CS 10310 F-54519 Vandoeuvre-lès-Nancy, France



Exemple of interest of Metadata for catalogs

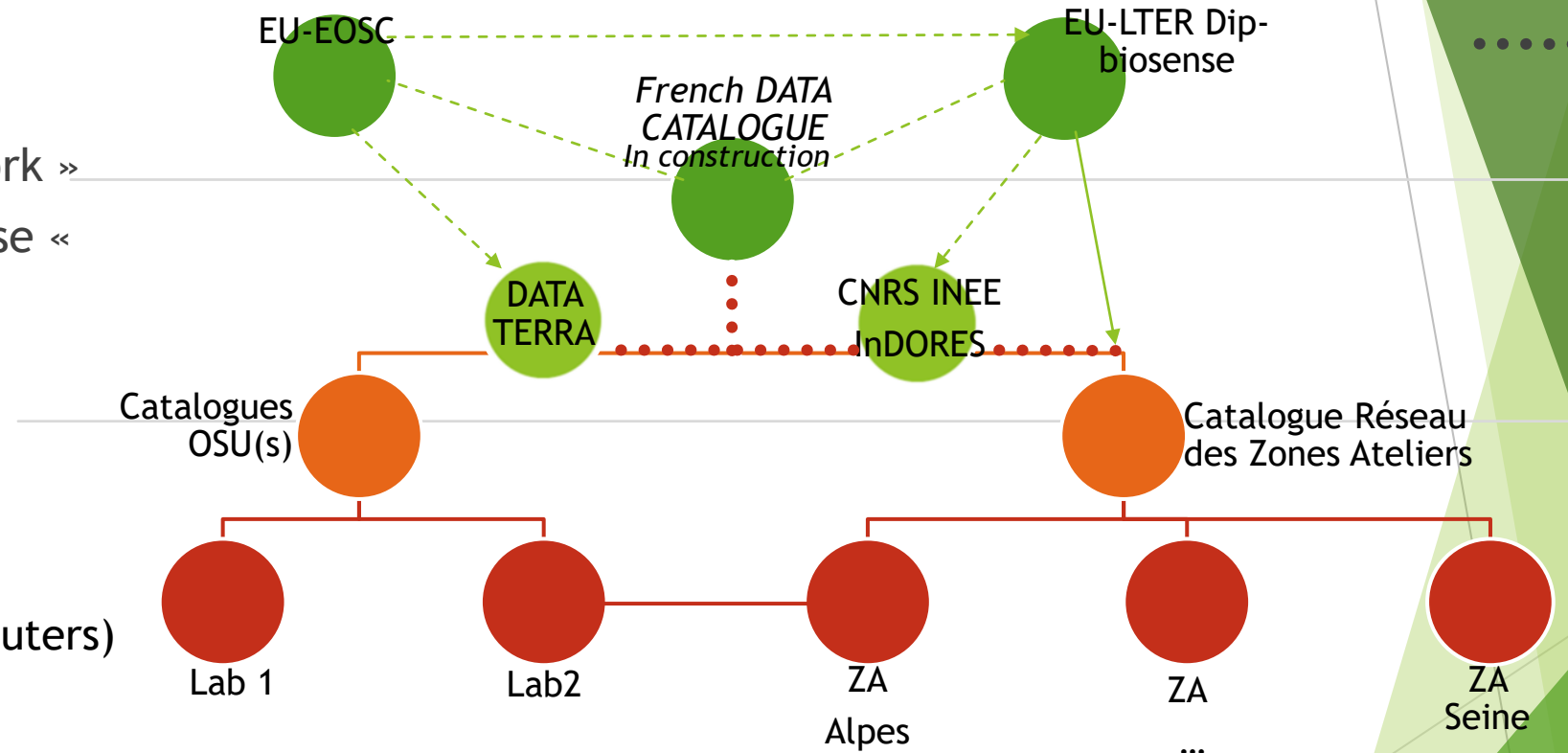
Harvesting : *Interest of catalogues in data visibility* (oct 2021)

Catalogs & Repository

- ▶ Metadata « Geonetwork »
- ▶ Repository « Dataverse »
- ▶ Others ...

Meta-search tools

- ▶ Data Cite search
- ▶ Data Citation Index (Reuters)
- ▶ Google Dataset Search
- ▶ Data Search (Elsevier)



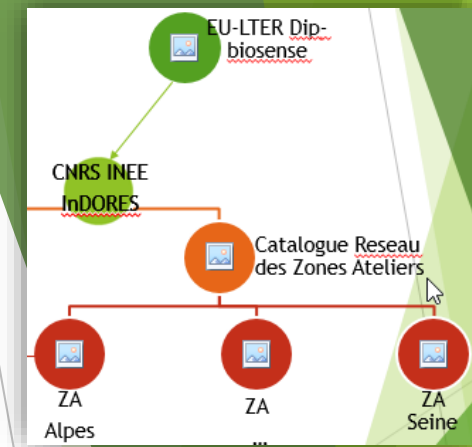
..... In progress

harvested by catalogues, integrators, European data infrastructures... more and more numerous (EOSC, OpenAIRE, DataONE, etc.)



Exemple of « metadata sheets trip »

(metadata standard - Iso19115 & Geonetwork catalogue)



1-metadata sheets with low metadata



Pollen analysis - Grozon (F-39), L'étang (Marais)

Téléchargements et liens



041165667

<http://www.sudoc.fr/041165667>

(suspicion of IT harvesting problems ...)

2-metadata sheets with more metadata



Section de carottage MUZ12-01A (IEFRA094Y) - Lac de la MUZELLE, Carottage (04/2012) Zone Atelier Alpes ZAA (LTER-FRANCE ROZA)

(do not have the link to the data)

3-metadata sheets with complete high metadata quality



Lac de la Muzelle - Carottage Muz12-I – Plant cover and erosion dynamics : impacts of pastoral activities and climatic changes over the last 1600 years - 2017 – (IGSN:IEFRA00A4)

Parcourir par **Thèmes INSPIRE** / Thématique

- Habitats et biotopes: 150
- Hydrographie: 83
- Usage des sols: 45
- Sols: 15
- Parcelles cadastrales: 9
- Caractéristiques géographiques...: 5
- Systèmes de maillage géograp...: 3
- Répartition de la population — ...: 2
- Réseaux de transport: 1
- Occupation des terres: 94
- Répartition des espèces: 76
- Altitude: 33
- Ressources minérales: 11
- Référentiels de coordonnées: 8
- Bâtiments: 4
- Zones à risque naturel: 3
- Zones de gestion, de restriction ...: 2
- Services d'utilité publique et ser...: 2
- Ortho-imagerie: 84
- Installations de suivi environne...: 64
- Interest of catalogues in data vi...: 17
- Installations agricoles et aquaco...: 10
- Dénominations géographiques: 5
- Géologie: 4
- Caractéristiques géographiques...: 3
- Unités administratives: 2
- Régions maritimes: 1
- Jeu de données: 745
- Carte: 14
- Jeux de données non géograph...: 10
- Collection de données: 10
- Dataset: 4
- Carte statique: 4
- Collection de session: 1
- Type d'entité: 1
- Modèle: 1
- Raster: 1



Exemple of « metadata sheets trip »

thanks to metadata standard (Iso19115 - Geonetwork catalogue)

ex of keyword : study area

→ Interest of indexation catalogues in data visibility

Local Catalogue

International thematic Repository

International Catalog (Sheet harvesting from à Local Catalog)

Local Catalogue

International thematic Repository

The screenshot shows a Google Dataset Search interface. At the top, the search bar contains 'muzelle' and the results section shows '36 ensembles de données trouvés'. The first result is 'Lac de la Muzelle - Bathymétrie_grid data - 2016' from 'catalogue.parcnational.fr'. A blue button 'Découvrir sur catalogue.parcnational.fr' is visible. Below it, the data is described as 'Ensemble de données mis à jour Sep 10, 2020' and 'Ensemble de données fourni par EDYTEM - UMR5204 - Environnements DYnamiques et TErritoires de la Montagne'. A map shows the geographical area around Lake Muzelle in the Massif des Écrins. The description states: 'Les lacs de montagne sont souvent situés dans des zones naturelles protégées, une caractéristique dans le lac Muzelle dans le Massif des Écrins.'



TOOLS FOR METADATA CREATING

Many tools / Many ways to do

(template, formulars, package R, database,...)

→ Find the one(s) that best suits you,
according to your profile, your(s) standard(s) needs, ...

Exemples for our « Paleo-community » :

- 1) (eng) International metadata Geosample : **IGSN**
- 2) (eng) interdisciplinary community Tool, base on standards DublinCore, DataCite, Iso19.115 (et EML) : **Geoflow**
- 3) (fr) OSU (OTELo) recommended template for fields to be filled Observ. Terre Environnement Lorraine
- 4) (fr) a webtool provided by DataCite
- 5) etc ...

} Initiate by world
Marine consortium

} See Slide N° 55
Funder FAO, IRD, INARe,
CNRS LTER-fr ...

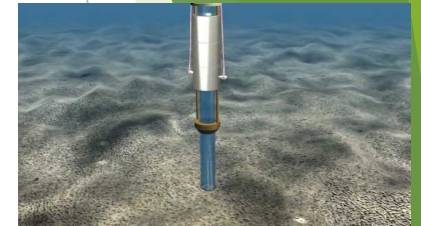
} <https://doranum.fr/metadonnees-standards-formats/outils-creation-metadonnees/>



-PART III- Practical from case studies & Basics FAIR Principles

LEVEL 2

- ▶ Usecase 1 for Sample collection : ex of cores
- ▶ Usecase 2 about Analytical data
- ▶ Which « Licences » ?
- ▶ Why How deposit your research data in a data-repository ?
- ▶ Use case of Workflow deposit data



	Total defects	A	B	C	D	E
A4836	131	37	21	28		45
A2524	86	20	24	21	1	20
A3719	75	17	13	18		27
A4452	73	5	33	17		18
A4088	72	14	16	12	2	28
A2103	68	14	13	14	1	26
A2156	68	16	13	19	2	18
A3681	66	12	16	9	1	28
A1368	50	11	15	12		12
A2610	39	5	7	12		15
Total	728	151	171	162	7	237

Remember F.A.I.R Principles : 2 levels of management



level
Discovery »

► Level 1 = Findable + Accessible

→ Don't expect a high quality of description

Nb : you can't do a Data-paper only with this level !



high level
of quality

► Level 2 = Level 1 + Interoperable & Reusable

beyond « Discovery »

→ Expect a high quality of description, precision for protocols, sensor, measurements, ...



Before usecases, let's continue to see some basics F.A.I.R. principles



To be FAIR, the data must indicate **their origin (Provenance)**
data are reusable = information to contextualize them
as the authors, their institution, the date of creation of the data, sources if aggregated data,

Link the relationship with other data already published, etc... (add "Read me" for additional information if neccessary)



To be FAIR, the data are **richly described** with a plurality of **specific and relevant attributes**

The more you know the context in which data was created (technical, setting, limit of detection,...), the more you can benefit from it. It can be technical metadata (contextual information on the data, even those which may seem useless to you)

addition of a text file of the type "DataDictionary"



To be FAIR, the **data** must be uniquely identifiable and sustainable **using a PID.**

DOI - IGSN - ORCID ...

To be FAIR, **the metadata** must **contain the PID of the dataset** described





Case of SAMPLES : *in Environnemental Sciences*

Most of the time, no science without samples !

UseCase 1

Samples are DATA

and all « sample's data » become
metadata of your future
Analytical measurement

→ Important to well describe them a the
begining of the project (+ fielwork report)

→ They are part of your « FAIR
Provenance » for your futur analytical datas

**Only one standard for (geological)
samples = IGSN**

Nb : Find your own Allocation Agent in your
country to declare your samples

The image shows a screenshot of the IGSN website. At the top, there is a navigation bar with links: HOME, ABOUT US, MEMBERSHIP, IGSN (logo), ALLOCATING AGENTS, RESOURCES, and IGSN 2040. Below the navigation bar is a large image of a hand holding a sample in a laboratory setting. A circular inset in the center of the image shows a close-up of a sample with a white label that reads 'IGSN: IEZHP0001'. Below the label is a QR code and a table of metadata:

IGSN:	IEZHP000
Sample Name:	M15-82
Other Name(s):	
Sample Type:	Individ
Parent IGSN:	Not F

On the right side of the image, there are two blue-bordered icons. The top one is a barcode icon with the text 'PID' below it. The bottom one is a barcode icon with the text 'Métadonnées avec PID' below it.

Link to
metadata &
data

<https://www.igs.org>



SAMPLES : Sample collection in Environmental Sciences, Data become Metadata !

- ▶ Common Metadata standard to describe a general geological sample
- ▶ A unique Identifier (for each sample and children) link in your articles, datapaper or IT Data infrastructure

exemple

Advanced Settings [Clear](#) Collector: SABATIER P

Total sample counts are listed, but only public sample metadata can be accessed. Please

46 Core(s) 22 Core Section(s) 132 Core Section Half(s)

Editors (Recommends IGSN)

- ▶ [Copernicus Publications](#)
- ▶ [AGU Publications](#), ...

- ➔ Assigne PID
- ➔ Good common metadata but it is « Poor » technical metadata
- ➔ Sufficient for discovery level

IGSN: IEFRA00BB

IGSN: IEFRA00BB
 Sample Name: THU10-I
 Other Name(s):
 Sample Type: Site
 Parent IGSN: Not Provided

Description

Material:	Sediment
Classification:	Not Provided
Field Name:	Not Provided
Description:	Not Provided
Age (min):	20000 years
Age (max):	Not Provided
Collection Method:	Coring>PistonCorer
Collection Method Description:	Not Provided
Size:	1900
Geological Age:	Holocene
Geological Unit:	Not Provided
Comment:	Not Provided
Purpose:	Not Provided

Geolocation

Latitude (WGS84):	45.530217
Longitude (WGS84):	6.056567
Northing (m) (UTM NAD83):	Not Provided
Easting (m) (UTM NAD83):	Not Provided
Zone:	Not Provided
Vertical Datum:	Not Provided
Elevation Start:	874 meters
Elevation End:	Not Provided
Nav Type:	Not Provided
Physiographic Feature:	lake
Name Of Physiographic Feature:	THUILE
Location Description:	Not Provided
Locality:	Not Provided
Locality Description:	Not Provided
Country:	Not Provided
State/Province:	Not Provided
County:	Not Provided
City:	Not Provided

Collection

Field Program/Cruise:	PALEOMAG26/04/10
Platform Type:	Small craft
Platform Name:	Not Provided
Platform Description:	Not Provided
Launch Type:	Not Provided
Launch Platform Name:	Not Provided
Launch ID:	Not Provided
Collector/Chief Scientist:	CROUZET C
Collector/Chief Scientist Detail:	curateur-edytem@univ-savoie.fr
Collection Start Date:	2010-04-25
Collection End Date:	Not Provided



Pit the PID in your général MD standard

UseCase 1



UseCase 1

Citation:

Bajard, Manon; Sabatier, Pierre; David, Fernand; Develle, Anne-Lise; Reyss, Jean-Louis; Fanget, Bernard; Malet, Emmanuel; Arnaud, Daniel; Augustin, Laurent; Cruzet, Christian; Poulenard, Jérôme; Fabien (2015): Analyses of lake La Thuile sediment core. PANGAEA, <https://doi.org/10.1594/PANGAEA.855423>

Supplement to: Bajard, M et al. (2015): Erosion record in Lake La Thuile sediments (Prealps): A multi-proxy approach to mountain landscape dynamics throughout the Holocene. *The Holocene*, **26(3)**, 350-360. <https://doi.org/10.1177/0959683615609750>

Always quote citation above when using data! You can download the citation in several formats:

[RIS Citation](#) [BrisTeX Citation](#) [Copy Citation](#) [Facebook](#) [Twitter](#) [Show More](#)

IGSN: IEFRA00BA

SESAR

IGSN: IEFRA00BA
 Sample Name: THU10_P1
 Other Name(s):
 Sample Type: Core
 Parent IGSN: Not Provided

Description

Material: Sediment
 Classification: Not Provided
 Field Name: Not Provided
 Description: Not Provided
 Age (min): Not Provided
 Age (max): 150 years
 Collection Method: Coring->GravityCorer->Pilot
 Collection Method Description: Not Provided
 Size: 1900
 Geological Age: Holocene
 Geological Unit: Not Provided
 Comment: Not Provided
 Purpose: Not Provided

Geolocation

Latitude (WGS84): 45.530000
 Longitude (WGS84): 6.056700
 Northing (m): Not Provided
 Easting (m): Not Provided
 Zone: Not Provided

Program/Cruise: PALEOMAG26/04/10
 Platform Type: Small craft
 Platform Name: Not Provided
 Platform Description: Not Provided
 Launch Type: Not Provided
 Launch Platform Name: Not Provided
 Launch ID: Not Provided
 Collector/Chief Scientist: CROUZET C
 Collector/Chief Scientist Detail: curateur-edyttem@univ-savoie.fr
 Collection Start Date: 2010-04-26

<https://app.geosamples.org/sample/igsn/IEFRA00BA>

Comment: IGSN of cores: THU10-P1: IEFRA00BA; THU10-I: IEFRA00BB



... provides an 18 m long sedimentary sequence spanning the entire Lateglacial/Holocene period. The high resolution multi-proxy (sedimentological, palynological, ... reveals the Holocene dynamics of erosion in the catchment in response to landscape modifications. The mountain belt is at relevant altitude to study past human activities and the ... large valleys to capture a local sedimentary signal. From 12,000 to 10,000 cal. BP (10 to 8 ka cal. BC), the onset of hardwood species triggered a drop in erosion following the ... 10,000 to 4500 cal. BP (8 to 2.5 ka cal. BC), the forest became denser and favored slope stabilization while erosion processes were very weak. A first erosive phase was initiated at ca. 4500 cal. BP ... presence in the catchment. Then, the forest declined at approximately 3000 cal. BP, suggesting the first human influence on the landscape. Two other erosive phases are related to anthropic activities: ... BP (550 cal. BC) during the Roman period and after 1600 cal. BP (350 cal. AD) with a substantial accentuation in the Middle Ages. In contrast, the lower erosion produced during the Little Ice Age, when climate ... generally considered to result in an increased erosion signal in this region, suggests that anthropic activities dominated the erosive processes and completely masked the natural effects of climate on erosion in the ...

Latitude: 45.530000 * Longitude: 6.056700
Date/Time Start: 2010-04-25T00:00:00 * Date/Time End: 2013-05-22T12:25:30

THU10-Mastercore * Latitude: 45.530000 * Longitude: 6.056700 * Date/Time: 2010-04-25T00:00:00 * Elevation: 874.0 m * Method/Device: Piston corer (PC) * Comment: IGSN of cores: THU10-P1: IEFRA00BA; THU10-I: IEFRA00BB

- Chemical composition of sediment core THU10-Mastercore. <https://doi.org/10.1594/PANGAEA.855423>
- Declination of sediment core THU10-Mastercore. <https://doi.org/10.1594/PANGAEA.855423>
- Grain size composition of sediment core THU10-Mastercore. <https://doi.org/10.1594/PANGAEA.855423>
- Loss on ignition in sediment core THU10-Mastercore. <https://doi.org/10.1594/PANGAEA.855423>
- Radiocarbon age determination of sediment core THU10-Mastercore. <https://doi.org/10.1594/PANGAEA.855423>
- Short-lived radionuclides in sediment core THU10-Mastercore. <https://doi.org/10.1594/PANGAEA.855423>

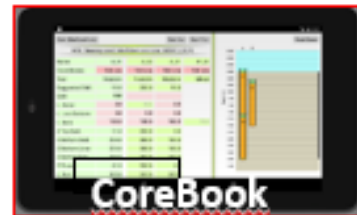
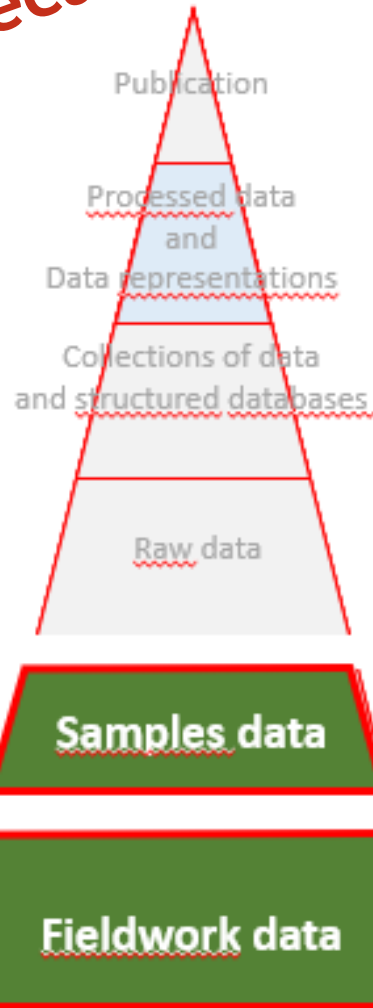
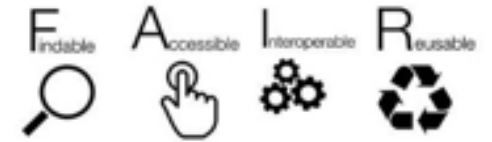


Exemple of usecase 1

How to turn kilos of mud into megabytes of data?

10 years of efforts in curating lake sediment cores and their associated results

→ Banking and displaying fieldwork data: the National Cyber-core repository (Cyber-carothèque nationale)



Export geographical point & fielddata for GIS SYSTEM

Arnaud F., Pignol C. et al.,
OZCAR-TERENO-2021



Exemple of usecase 1

Provenance : Collect your sample collection

In France :

French Cyber-Core repository

(marine, lacustrine, river, peatland, ice cores)

► Directly on the field, collect data by [APK Android](#)

OR

► after the field work, Import CCN template (or for your core legacy)



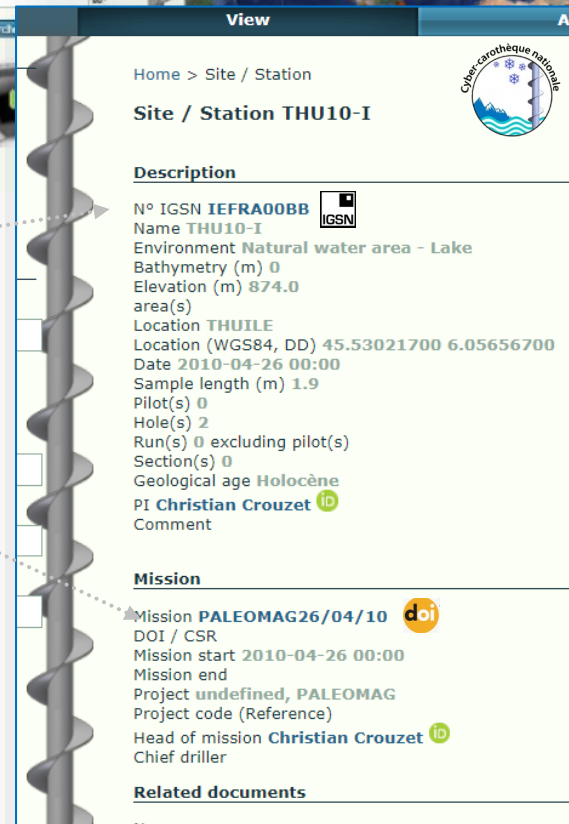
Download file here

In the end of 2021 = version v2 →


+ IGSN automatic assignation

+ Campaign PDF report (DOI)

+ Pivot format for link to Core Storage Db (« Collec-Science » ?)



Other Databases for samples :

- Other IGSN AAgent IGSN like IFREMER, [SESAR](#)
- NOAA Index Marine Lacustrine [IMLGS](#)
- CSDCO Continental Scientific Drilling Holes / LacCore
- IODP Core Repository 

Example of usecase 1

How to turn kilos of mud into megabytes of data?

10 years of efforts in curating lake sediment cores and their associated results

→ Banking and displaying fieldwork data: the National Cyber-core repository (Cyber-carothèque nationale)

Mission information

Mission PALAS19-GUYNEMER

Description

Fichier source Corebook

XML source file

Code
Début de mission 27-11-2019 06:00
Fin de mission 27-11-2019 18:00

Description

Projet PALAS

Contexte scientifique

Connaissances préliminaires

Commentaire

Confidentialité Publique

Remerciements

DOI

URL

CSR

Link to project: permit to display all surveys/cores within the same project

Automatic survey report with DOI (under development)

Participants

Chef foreur **Arnaud Fabien** ^{ID}

Chef de mission **Arnaud Fabien** ^{ID}

Acknowledgement of crew members with determined role and ORCID identification

PI(s) des prélèvements (liste non exhaustive)

Arnaud Fabien ^{ID}

Financier

Aucun(e)

Documents liés

GUY19-I-Pil-05-20191223-141634.jpg

GUY19-I-Pil-05-20191223-141716.jpg

GUY19-I-Pil-05-20191223-141736.jpg

GUY19-I.20190619-195534.Coring.Schema.png

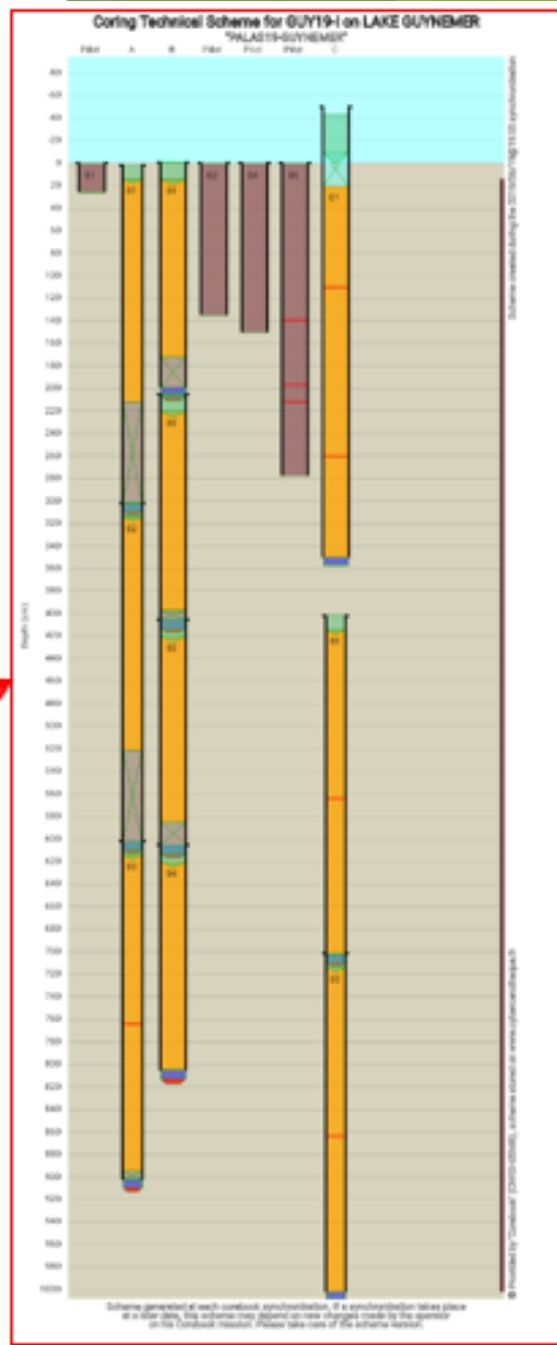
GUY19-I-B-01-20191201-064036.mp4

EPGUY19-01-20191204-113555.jpg

EPGUY19-01-20191204-113630.jpg

Any additional document (pictures, sounds, videos) taken with CoreBook App.

When applicable: link to the sketches of parallel corings



How to turn kilos of mud into megabytes of data?

10 years of efforts in curating lake sediment cores and their associated results

→ Banking and displaying fieldwork data: the National Cyber-core repository ([Cyber-carothèque nationale](#))

Sample (core) information

Description

N° IGSN [TOAE0000000078](#)

Site / Station [GUY19-I](#)

Nom [GUY19-I-A-03](#)

Matériel [Sédiment](#)

Type de prélèvement [Carotte \(recouvrement\)](#)

Environnement [Étendue d'eau \(douce\) naturelle - Lac](#)

Plateforme

Profondeur bathymétrique (m) [94](#)

Altitude (m) [131.0](#)

zone(s) [French Southern and Antarctic Lands, France](#)

Lieu [LAKE GUYNEMER](#)

Point GPS (WGS84, DD) [-49.11344500 69.02049300](#)

Date [28-11-2019 16:01](#)

Nom du carottier [U-NIEDERREITER 63-3m \(EDY\)](#)

Nom de la configuration [2-NIED63-3m PO/TC \(EDY\)](#)

Diamètre intérieur (mm) [63](#)

Longueur du prélèvement (m) [2.93](#)

Longueur matière récupérée (m) [2.81](#)

Époque géologique

PI [Fabien Arnaud](#)

Sondeur/Foreur en chef

Commentaire



IGSN unique identifier

Codified usual name :
Site-Year-Station-Hole-Section

Where is the core stored?

Link to coring tool information

Children subsamples

Mission

Mission [PALAS19-GUYNEMER](#)

DOI / CSR

Début de mission [27-11-2019 06:00](#)

Fin de mission [27-11-2019 18:00](#)

Projet [PALAS](#)

Code projet (Reference)

Chef de mission [Fabien Arnaud](#)

Sondeur/Foreur en chef [Fabien Arnaud](#)

Repository

Repository [Environnements, DYnamiques et TErritoires de la Montagne \(EDYTEM\)](#)

Documents liés

Aucun(e)

Sous-prélèvements associés

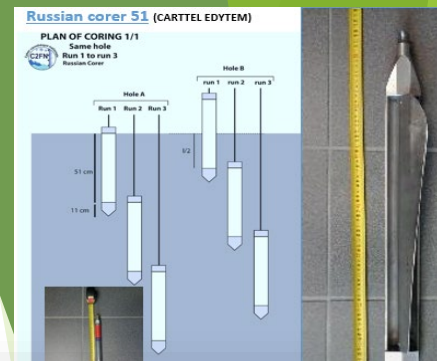
Export csv

Nom	Parent	IGSN	Profondeur supé
GUY19-I-A-03-A Section	GUY19-I-A-03 Carotte (recouvrement)	TOAE00000000582	
GUY19-I-A-03-B Section	GUY19-I-A-03 Carotte (recouvrement)	TOAE00000000583	

How to turn kilos of mud into megabytes of data?

10 years of efforts in curating lake sediment cores and their associated results

→ Banking and displaying fieldwork data: the National Cyber-core repository (Cyber-carothèque nationale)



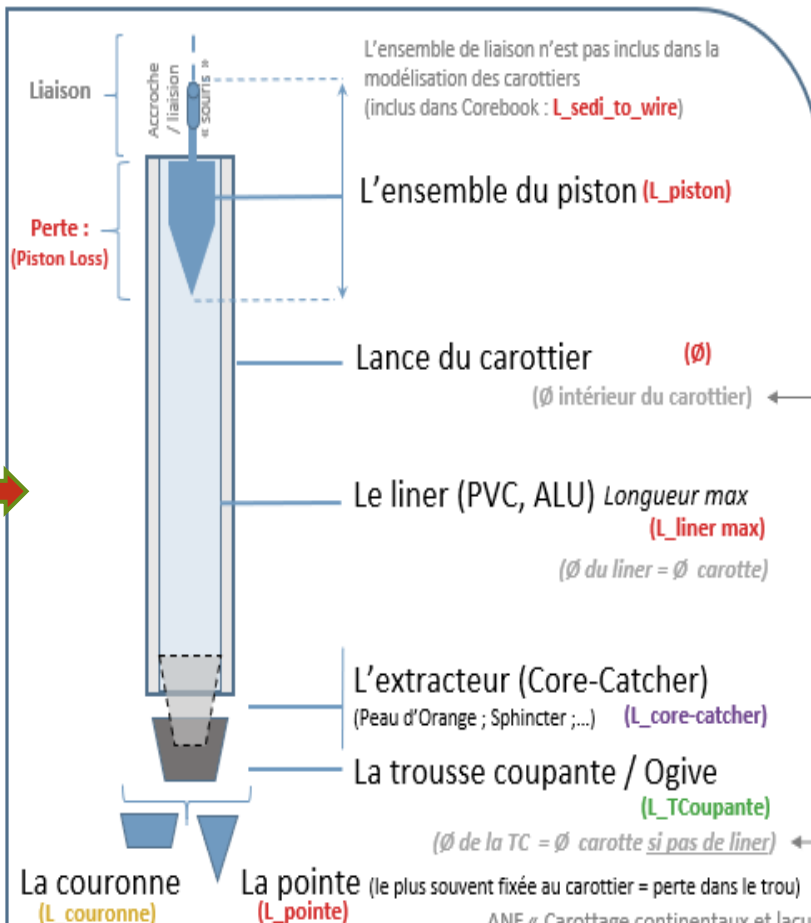
Sample (core) information

Description

N° IGSN **TOAE000000078**
 Site / Station **GUY19-I**
 Nom **GUY19-I-A-03**
 Matériel **Sédiment**
 Type de prélèvement **Carotte (recouvrement)**
 Environnement **Étendue d'eau (douce) naturelle - Lac Plateforme**
 Profondeur bathymétrique (m) **94**
 Altitude (m) **131.0**
 zone(s) **French Southern and Antarctic Lands, France**
 Lieu **LAKE GUYNEMER**
 Point GPS (WGS84, DD) **-49.11344500 69.02049300**
 Date **28-11-2019 16:01**
 Nom du carottier **U-NIEDERREITER 63-3m (EDY)**
 Nom de la configuration **2-NIED63-3m PO/TC (EDY)**
 Diamètre intérieur (mm) **63**
 Longueur du prélèvement (m) **2.93**
 Longueur matière récupérée (m) **2.81**
 Époque géologique
 PI **Fabien Arnaud**
 Sondeur/Foreur en chef
 Commentaire

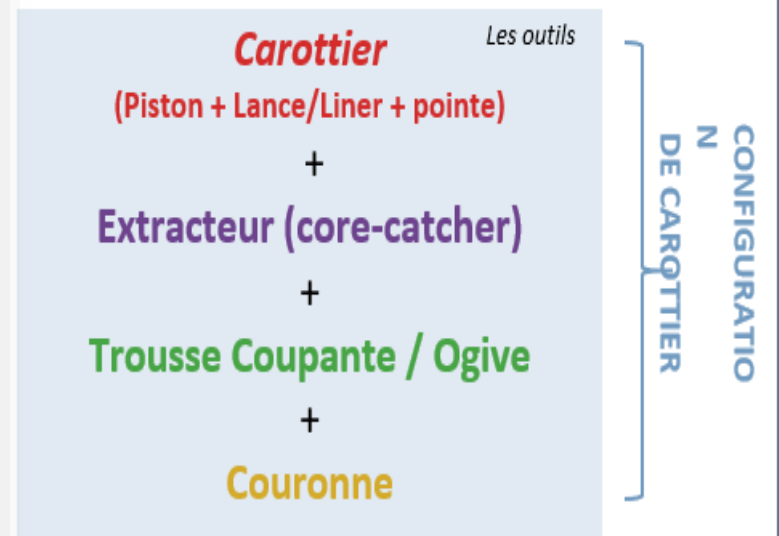
Link to coring tool information

Variables d'un carottier, extracteur et trousse coupante



Qu'est ce qu'une « CONFIGURATION » de carottier ?

?? Selon les carottiers, multitudes de variation de 'Configuration' autour de :



?? Les 'configs' entre dans le référentiel envoyé au Corebook



UseCase 2

Provenance for Analytical Data

→ Well structure data in tables

→ **MAKE REPETABLE “TEMPLATE”**

What
reference
metadata
standard
for this
type of
data ?

- ▶ **Sheet Tab « Read-me »** : General & Context information contexte (from MD std) + thematic control vocabulary
(*Geoflow template slide 55*)
- ▶ **Sheet Tab « Data dictionary »** :
 - variable(s) names** (+ definitions + control vocabulary)
 - unit of measure** (+ control vocabulary),
 - quality** : define unmeasured values (9999, -9999), n/a, value below the measurement ...), uncertainty, standard calibration, suppressed value,



Keep **PROTOCOLS** close to your data & share them !

→ well describe your lab or study protocols

1- txt procols, sheet of laboratory notebook,

Electronic labnotebook (ELN)

2-publication of a protocols => Link DOI to your article, Data deposit or Data paper



nature protocols

protocolexchange

“Five keys to writing a reproducible lab protocol”, Nature, 2021



Provenance : Analytical Data

Well structure data in tables

UseCase 2

- ▶ **Sheet Tab** « Rawdata » (if necessary, otherwise link to the permanent repository which store row data)
- ▶ **Sheet Tab** « Clean Data » : only variables values, Good name for header of colonnes (vocabularies?)
- ▶ *any information other than define in the Read-me*

➔ Do **VERSIONNING** data as your work progresses
(with date consistent name : AAAA-MM-DD_XXXXX.xls)

- **Sheet Tab** « Work in progresse Interprétation session date 1 »
- **Sheet Tab** « Work in progresse Interprétation session date 2 »
- **Sheet Tab** « Work in progresse Interprétation session date 3 »

Keep what
you do at
the end of
time session

ROZA : standardised data & metadata template

Exemple of a 2014 « READ ME » of a sheet from a projet « ROZA »

UseCase 2



Usual name

IGSN

Usage license

NOM CAROTTE	POU 14P 1w	SEBAR ID SN	IGSN:EDYFER0025	LICENCE UTILISATION (Creative Commons)	
METADONNEES / REFERENCES					
Version 3.0.0 - Révisé: 2016-01-01					
WHAT KIND OF ANALYSE	XRF core scanning				
WHERE / LAKE	Fouie				
WHO / NAME RE SEAR CHER	Pierre Sabatier				
WHEN	25/01/2016				
WHY	Application				
DOI	PALAS				
Tube length (mm)	890				
Top sediment (mm)	38				
Bottom sediment (mm)	889				
Sediment length (mm)	851				
Tube Diameter	53mm				
INSTRUMENT & SETTINGS					
Source	Avaatech Core Scanner				
Vol tage (kV)	10	30	50		
Intensity (mA)	1.2	0.75			
LiveTime (s)	30	30			
Step (mm)	1	1			
Measured Area (Core score/DownCore)	15*1	15*1			
Filter	No filter	Pd-Thin	Cu		
ELEMENTS concerned :					
	Al_Area	Cu_Area	Ag_Area		
	Si_Area	Zn_Area	Sr_Area		
	P_Area	Ga_Area			
	S_Area	Br_Area			
	Q_Area	Rb_Area			
	K_Area	Gr_Area			
	Ca_Area	Y_Area			
	Tl_Area	Zr_Area			
	Cr_Area	Nb_Area			
	Mn_Area	Mo_Area			
	Fe_Area	Pb_Area			
	Rh_Area	Bi_Area			
Standards (Yes/No, Local Number, Certified, standard's Name)	YES	134	MESS-3, SARM-4		
Operator	Anne-Lise deville (orcid:0000028244106)				
Processing softwares	Winaxi, Winaxbatch (Avaatech)				

Sample characteristics (from Cyber-core repository)

Measurements settings

From Lab Note book

Measurement variables + definition (here: list of elements)

Operator's acknowledgement (with ORCID)



Ongoing effort

When possible, use thesaurus/controlled vocabularies (CHEBI; IUPAC; QUTD; ...) to provide unambiguous variables and units (interoperability)

→ Referenced in metadata catalogs



Provenance : Analytical Data

Well structure data in tables

UseCase 2

A Data Dictionary

	A	B	C
1	Data sheet names	Data description	Name of the file where data are located
2	Water silicone	These data are giving the concentration in the water of the different PCB, concentrations calculated with two methods : Smedes and Yates. These data are named like this because silicone is used as a medium for the analysis of PCB contained in water with a chromatograph.	PCB_Muzelle_2014 or PCB_Muzelle_2015
3	Air	These data are giving the deposit flow between air and water of PCB for the both phase of the atmosphere : dry air and wet air.	PCB_Muzelle_2014 or PCB_Muzelle_2015
4	Sediment trap	These data are giving the sedimentation flow of PCB at the surface and the bottom of the lake .	PCB_Muzelle_2014 or PCB_Muzelle_2015
5	LOD	These data are giving the limit of detection for each PCB in this study.	PCB_Muzelle_2014 or PCB_Muzelle_2015

7	Remark to users	
8	Remark	All the given data in the files are first scale data. Some "second scale" data are available but the fluxes of transfers between compartments and mass balance in the lake. They are calculated by have access to it by contacting Christine Piot, Edytem, Savoie-Mont Blanc University.
11	Signification of abbreviation and specific values in both data (attribute) dictionary and datafiles	
12	Cw	Concentration in water
13	<LOD	Measures below the limit of detection. In chromatography the LOD is equal to three times the value of the background noise.
14	Profile	Ratio between concentration of compound and the sum of concentrations of all compounds in percent
15	PCB	Polychlorinated biphenyl compounds
16	NA	Not Available

	A	B	C	D
1	Attribute / variable	Title	Définition	Comments
2	PCB 28	Polychlorinated biphenyl 28	N°cas : 7012-37-5	7012-37-5@https://commonchemistry.cas.org/detail?cas_rn=7012-37-5&search=7012-37-5
3	PCB 52	Polychlorinated biphenyl 52	N°cas : 35693-99-3	35693-99-3@https://commonchemistry.cas.org/detail?cas_rn=35693-99-3&search=35693-99-3
4	PCB 101	Polychlorinated biphenyl 101	N°cas : 37680-73-2	37680-73-2@https://commonchemistry.cas.org/detail?cas_rn=37680-73-2&search=37680-73-2
5	PCB 118	Polychlorinated biphenyl 118	N°cas : 31508-00-6	31508-00-6@https://commonchemistry.cas.org/detail?cas_rn=31508-00-6&search=31508-00-6
6	PCB 153	Polychlorinated biphenyl 153	N°cas : 35065-27-1	35065-27-1@https://commonchemistry.cas.org/detail?cas_rn=35065-27-1&search=35065-27-1
7	PCB 138	Polychlorinated biphenyl 138	N°cas : 35065-28-2	35065-28-2@https://commonchemistry.cas.org/detail?cas_rn=35065-28-2&search=35065-28-2
8	PCB 180	Polychlorinated biphenyl 180	N°cas : 35065-29-3	35065-29-3@https://commonchemistry.cas.org/detail?cas_rn=35065-29-3&search=35065-29-3
9	Sum PCB	Sum of PCB		
10				



LICENCES : explicitly define the rights of reusers regarding your data

-Should they cite you?
-Can they modify the data?
-Make commercial use of it?


The choice depend on the **rights you can grant** depending on the legal **nature of the data**

For example, they may already be protected by copyright (droit auteur) or be subject to exceptions to their open access (collaboration with private, military...).



LICENCE OUVERTE
OPEN LICENCE

La Licence Etalab a été conçue par le Gouvernement français pour faciliter la mise en place de l'Open Data. Elle équivaut à la licence CC-BY.



Les licences CC permettent de définir plusieurs restrictions, comme l'interdiction d'usage commercial ou de modification.

ODbL

L'Open Database Licence (ODbL) est une licence spécifique permettant d'exploiter publiquement des bases de données.

Careful with  "Personal Data »

The EU-General Data Protection Regulation – **GDPR (RGPD)**

**ODC-BY,
CC-BY 2.0**



In France : Institutions get an « Déléguee à la protection des données (DPO)

CC0 / CC-By / CC-By-SA ...



- ➔ A metadata catalogue **must inform** the licenses of the dataset
- ➔ A data repository **must provide** the chosen license for its reuse

Ca commence a bien F.A.I.R !



5 minutes ?



Why deposit your research data in a data-repository ?

Depositing your data in a warehouse (data and metadata) ensures their **preservation, visibility and access**, thus facilitating their **sharing and reuse**

Several benefits :

- ✓ compliance with the **Opendata** recommendations of funders and institutions
- ✓ **data visibility** and easy access for catalogues or search engines discovery
- ✓ **reuse by third parties** and citation of the dataset facilitated (PID identifier)
- ✓ management of data sharing arrangements through **licensing**
- ✓ **data interoperability** (use of metadata standards and vocabularies)
- ✓ data retention in a **secure environment**

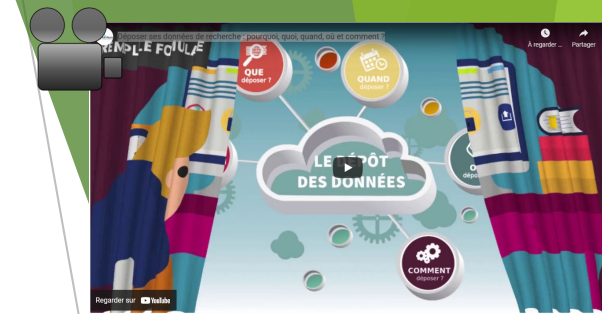
→ **improved research reproducibility, integrity and scientific validation**

→ **reuse in new studies and innovations**

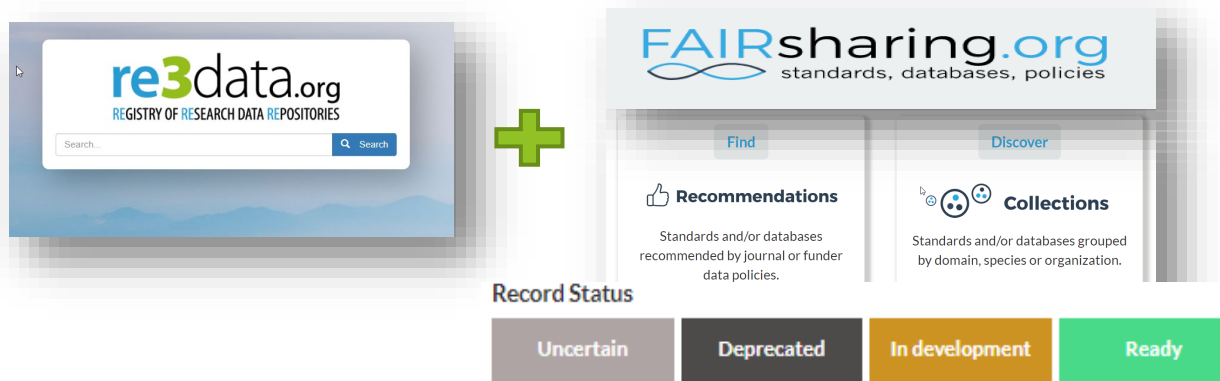




How to find a data-repository ?



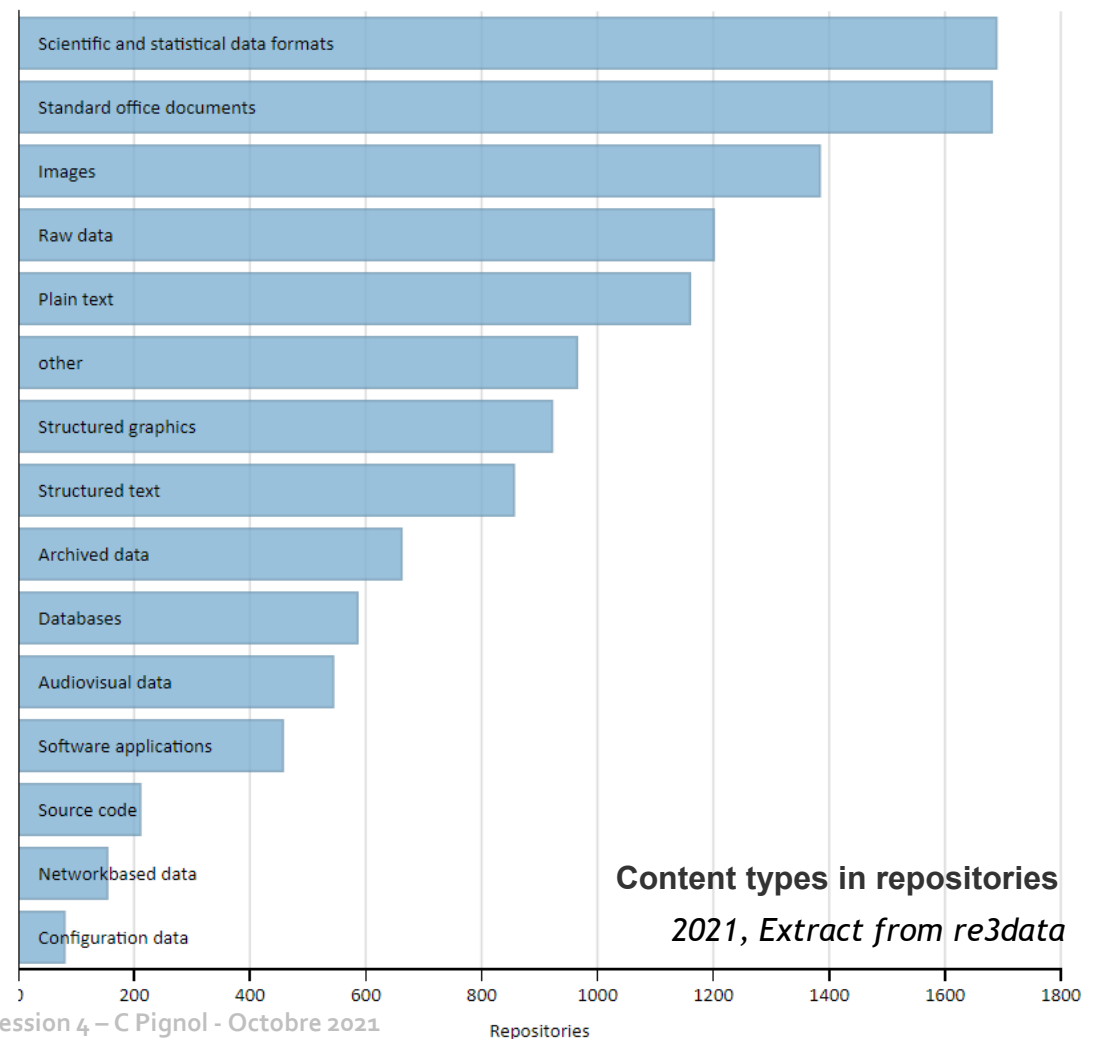
1) Searching repository by « annuaire » (directories) or according editor recommendations



→ Search for storage « Subject(s) » = Oceanography Geology and Palaeontology Geochemistry

→ *Nature* Data Repository Guidance (+90 entrepôts recommandés) Scientific Data repository

SIST https://sist20.sciencesconf.org/data/pages/SIST20_entrepot_de_donnees_Desconnets.pdf <https://doranum.fr/depot-entrepots/criteres-choix-entrepot/>





Where : how to choose a data repository ?

Si vous ne savez pas lequel choisir, vous pouvez vous appuyer sur des critères présentés ci-dessous :

<https://doranum.fr/depot-entrepots/criteres-choix-entrepot/>

2) Aspects to be checked directly on the website :

- ✓ **Easy deposit** : easy to use ? user-friendly, ergonomic? Some warehouses also offer support, tutorials
- ✓ **The quality of the description (metadata)** : The quality of the description of the datasets is important for easy retrieval (based on standards ?)
- ✓ **Usage statistics** : Does the warehouse offer usage, consultation and data download statistics ?

3) check the warehouse “policies”:

Exemple, consulting the [Policy](#) (*Zenodo free & open*)

- ✓ **Server hosting location** : in Europe ? US ? Data sovereignty ? Lockdown ?
- ✓ **Long-term preservation of data** : **how many years** ?
- ✓ **Cost of deposit**
Dryad, Figshare are paid/ Pangaea, Zenodoo are free.
The cost can also be related to the volume of files
- ✓ **Statut public / academic or private ?**
- ✓ **Supported data types**
Dryad, only accepts data related to a publication.
Zenodo, accepts any type of data.
- ✓ **Accepted file formats**
recommend specific file formats.
(Non-proprietary formats, list recommended formats, preferred format)
- ✓ **Link(s) with articles ?**



A Comparative Review of Various generalist Data Repositories (1)

Extraction of 4 repository

From 2018
<https://dataverse.org/bl/og/comparative-review-various-data-repositories>

	Dryad	figshare	Mendeley Data	Zenodo
Data citation automatically generated	Yes	Yes	Yes	Yes
Ability to embargo files	Yes	Yes	Yes	Yes
Analyzing tabular data (aside from geospatial mapping)	No	No	Yes	No
Filetypes supported (list)	"All data types and formats within reason"	All file types	All file types	All file types
Handling large data 	Yes - Limit of 20GB per file for "package", charge for more storage	Yes - Default limit of 5GB per file, can support up to 5TB	Planned Q2 18	Yes - limit of 50GB per file. Upload page mentions you can contact them about larger datasets.
Mapping of Geospatial files	No	Yes	Yes	No
Previewing and/or analyzing uploaded non-tabular data	No	Yes	Yes (preview images, tabular data)	Yes (previewing, not analyzing)
Previewing tabular files	No	Yes	Yes	Yes
Provenance	No	No	Yes	No
Users are able to control dataset file hierarchy + directory structure	No	No (can re-order files, but can't create hierarchy or directories)	Planned Jan 18	No
Robustness of metadata	Yes	Yes	Yes (depositors can describe files today;	Yes
Support for controlled vocabulary terms with URIs	?	No	No	Yes 
Tracking citations with altmetrics	No	Yes	No	No
Ability to access older versions of files	Yes - on data package page	Yes	Yes	Yes
Dataset-level versioning	Yes	Yes	Yes	No
Faceted search	Yes	Yes	Planned Q1 18	Yes



A Comparative Review of Various generalist Data Repositories (1)

Extraction of 4 repository

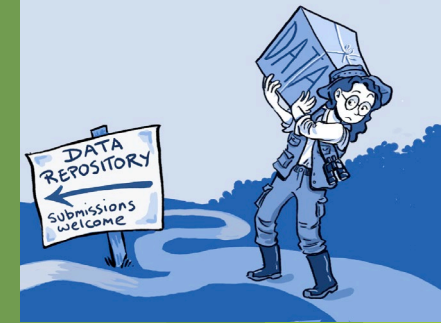
From 2018
<https://dataverse.org/blog/comparative-review-various-data-repositories>

	Dryad	figshare	Mendeley Data	Zenodo
Permissions	Yes (ability to apply embargo)	Yes (ability to apply embargo or make data "confidential")	Yes, with embargo	Yes (Includes embargoed, restricted, and closed access)
Terms of use, copyright	No - all uploaded data is CC0	Yes	Yes	Yes
Open Source	Yes	No	No	Yes
How is maintenance or development of repository funded?	Grants; Service charges	Unknown (Funded by Digital Science while retaining autonomy)	Private investments (Elsevier; Mendeley Ltd)	Institutional (OpenAIRE / European Commission); donations (via CERN & ...)
Free To Use	No	Yes	Yes	Yes
Institutional Fees?	No	Yes - "Figshare for Institutions" is priced based on "research intensity of the institution"	No	No
Journal Fees?	No	Yes - "Figshare for Publishers"	No	No
Paid services?	Yes	Yes - dedicated support team	No	No
Can depositors pay for extra storage?	Yes - "For data packages in excess of the 20GB size limit, submitters will be charged \$50 for each additional 10GB, or part thereof. (Packages between 20 and 30GB = \$50, between 30 and 40GB = \$100, and so on)."	No	No	No
Preservation technology?	Yes - DataONE	Yes - Chronopolis	Yes - Amazon S3 servers in Ireland and dark archive storage with DANS	Yes - All data files are stored in CERN Data Centres, primarily Geneva, with replicas in Budapest.
Most popular subject tags	Adaptation / Population Genetics - Empirical / Speciation	Unknown	Unknown	Taxonomy / Biodiversity / Animalia
Subject vocabulary	Folksonomy	Australian and New Zealand Standard Research Classification (ANZSRC), 2008	Elsevier's OmniScience taxonomy	Folksonomy



Exemple of data-repository

non-exhaustive



- **Institutional (in France)** : [Dataverse Cirad](#), [Datapartage](#) (INRAe), [DataSuds](#) (IRD),

[DataTERRA](#) (4 Pôles de données ODATIS, AERIS, THEIA, Form@Terre),

In nov 2021 [Dataverse InDORES](#) (INEE-CNRS) , *in mars 2022* [RechercheDataGouv](#)

- **Institutional (in Europe)** : [NERC Data Centres](#) (GB), ...

- **International multidisciplinary** : [Zenodo](#), [Dryad](#), [Figshare](#)

- **Own by editors** : [GigaDB](#) (Oxford Univ. Press), [Dataverse Ubiquity Press](#),

- **Thematic and disciplinary** : **Repository or Database !**

[GenBank](#) (séquences génétiques), [TRY](#) (caractères botaniques), [GBIF](#) (biodiversité), [NOAA](#), [HydroShare \(CUAHSI\)](#),
[EarthChem](#), [Pangaea](#) (sc de la terre&environnement)

- **For « Protocols »** : [Protocols.io](#) (protocols), [Nature ProtocolsExchange](#), [Plos Protocols](#) , etc. (*txt, video*)

- **For «Code and software »** : [Software Heritage \(or via HAL\)](#)



Repository Data-repository in connection with other devices/ infrastructure

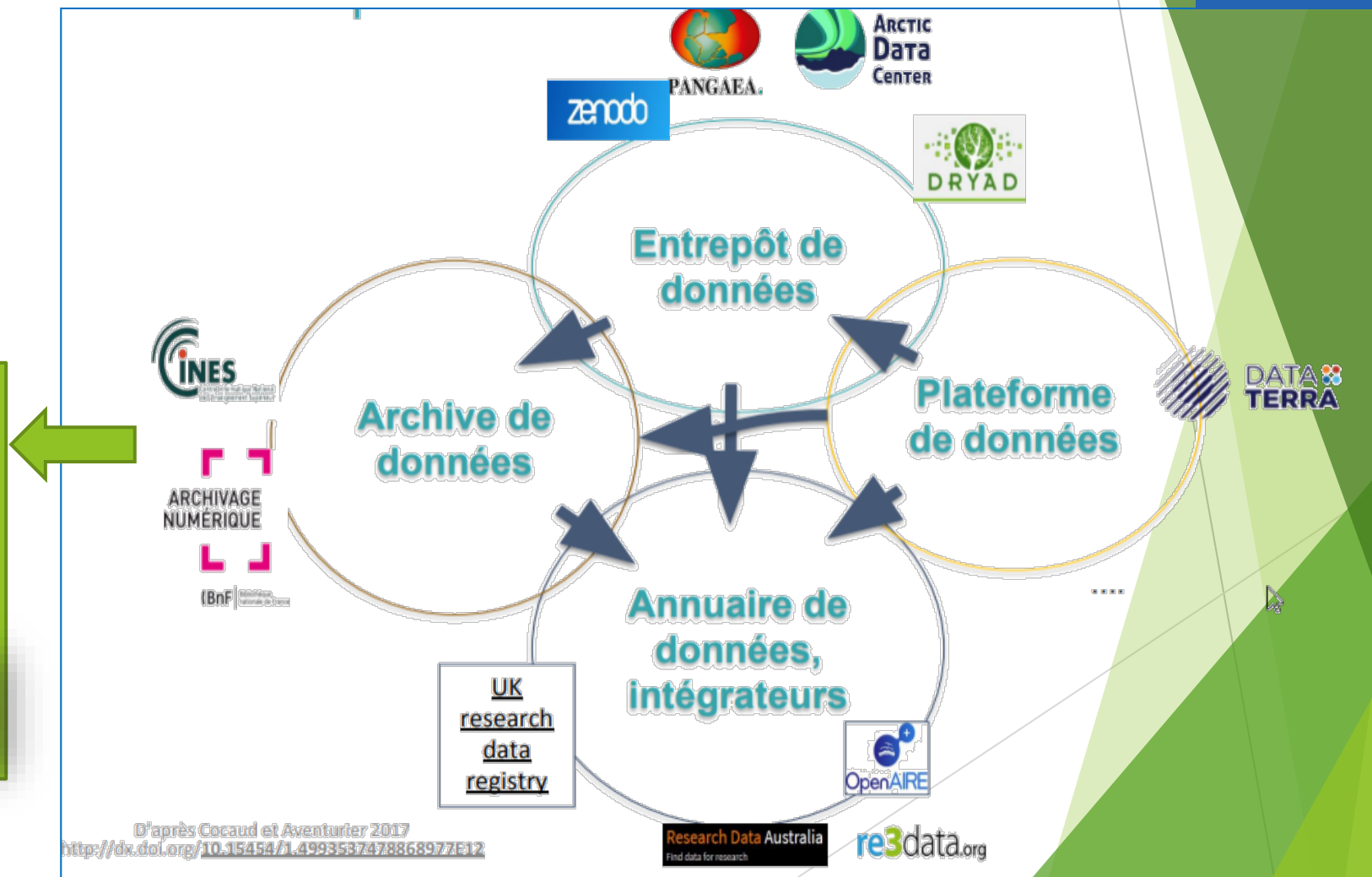
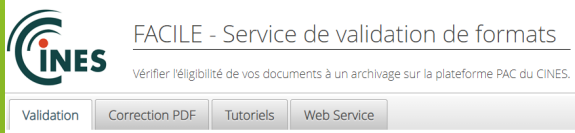


Stockage
≠
Archive

FOCUS

data archive
longevity of format file

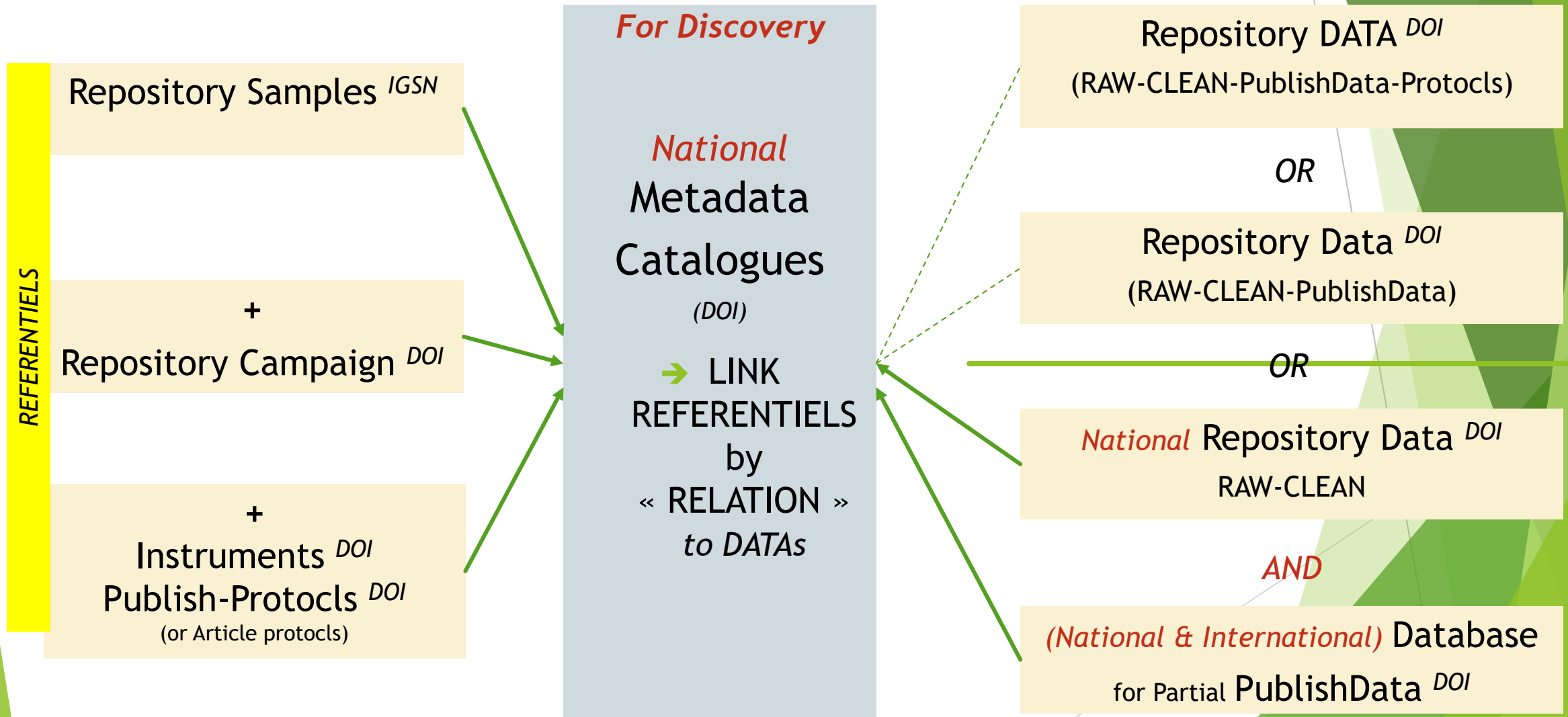
<https://www.cines.fr/archivage/des-expertises/les-formats-de-fichier/>



D'après Cocaud et Aventuret 2017
<http://dx.doi.org/10.15454/1.49935374738868977912>



How repository-data network could be structured next ? *It's depend when you do it in your Datacycle !*





An example of tools that can help *Geoflow* a metadata management workflow for everyone !

→ **1a,b,c** - You prepare your data (at different moments of the collect)

→ **2** - You use *Geoflow* to expose Metadata & Data

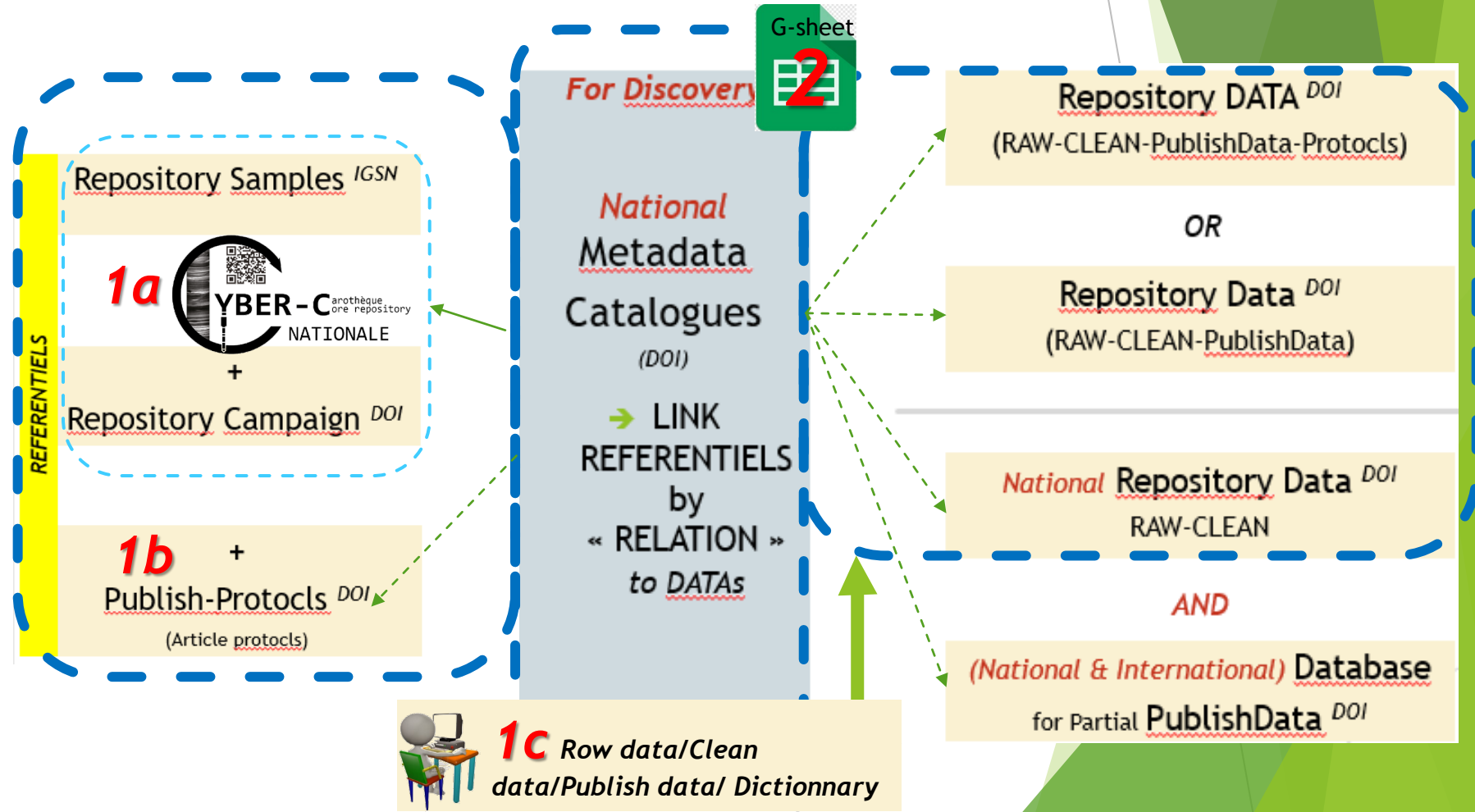
→ GEOFLOW

Multi metadata pivot for publishing ***in one time*** to catalogMD and Data Repository

geoflow : workflow R pour gérer les données spatiales

Julien Barde, Emmanuel Blondel, Wilfried Heintz

<https://github.com/eblondel/geoflow>

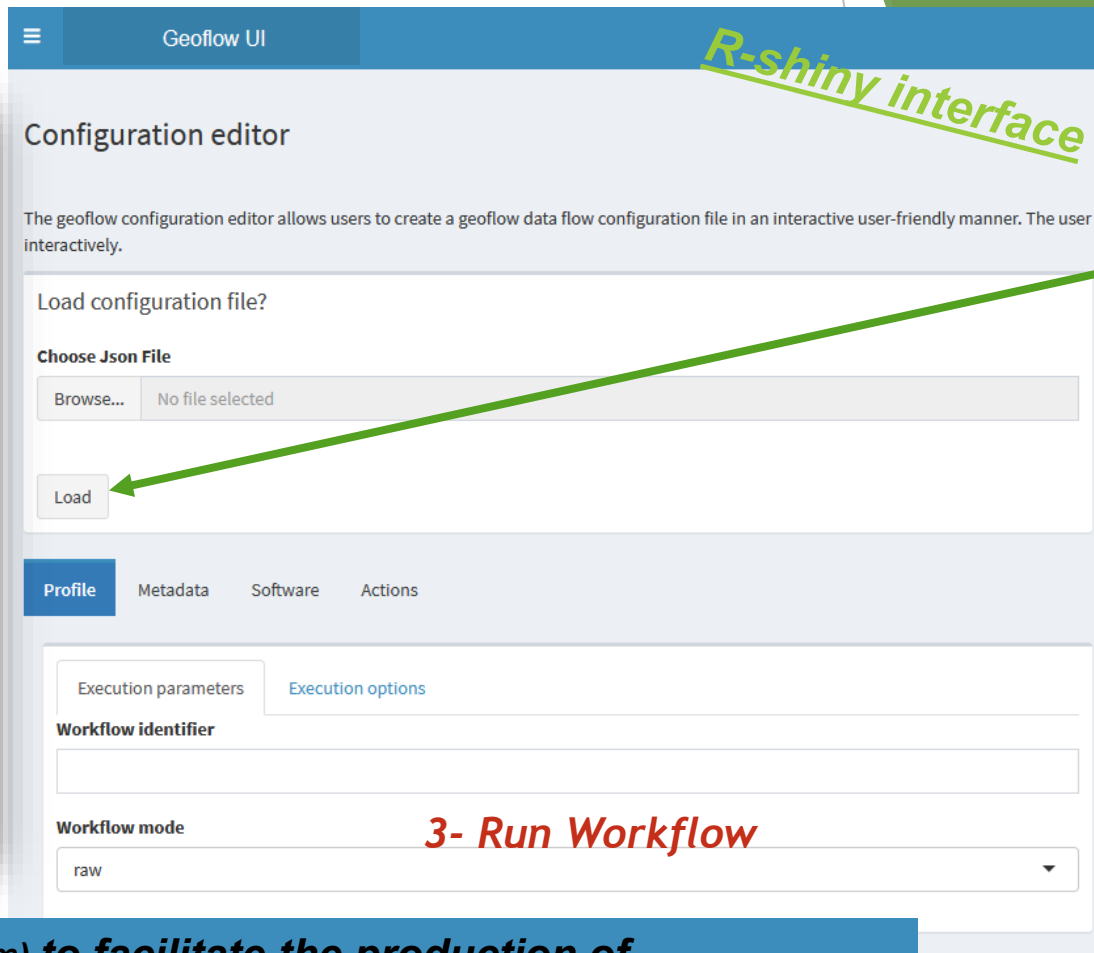
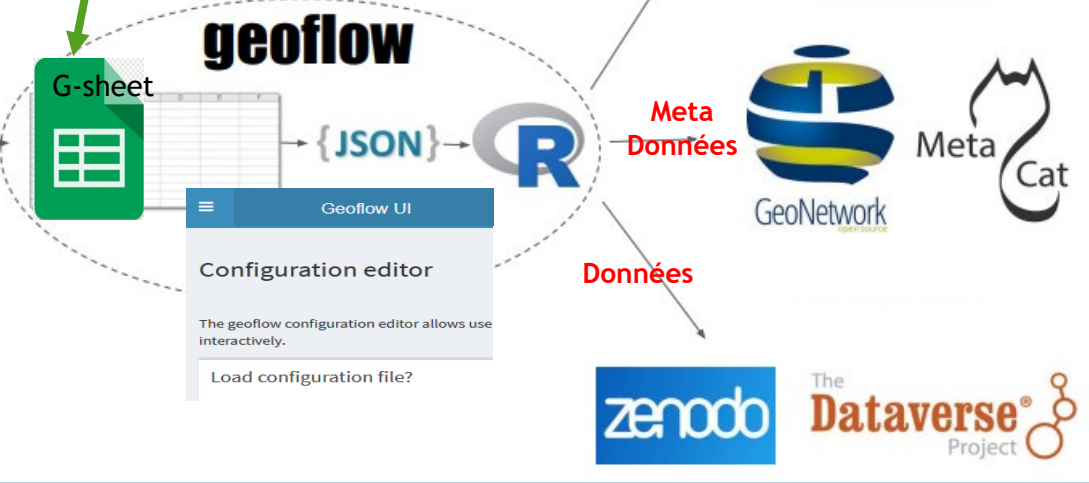




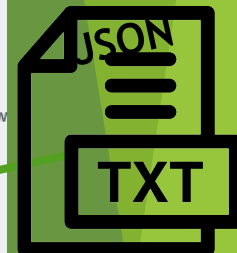
An example of tools that can help you *Geoflow* a metadata management workflow for everyone !

Workflow

1- Simple list of common metadata



R-shiny interface



2- Text Parameters to link Metadata & data to catalogues & repository

3- Run Workflow

geo

R-coded metadata management tool, (supported by the R Consortium) to facilitate the production of standardized metadata (DublinCore, Datacite, ISO 19115, Ecological Metadata Language) through a very simplified interface (*R-shiny interface-test* (INRAE), enabling collaborative metadata editing (Googlesheet,...))



Exemple of the geoflow guidline :

Part of metadata in the geoflow workflow

	A	B	C	D	E
1	Colonne	Type de données	Définition	Besoin	Nomenclature clés - valeurs
5	Subject	Liste de mots-clé décrivant le jeu de données	Sujets / matières associés au(x) jeu(x) de données. Un sujet se résume par un ou plusieurs mots-clés. Un sujet peut être composé de plusieurs mots-clés. On sépare les mots-clés d'un même sujet par une virgule. Plusieurs sujets peuvent décrire un jeu de données. On sépare les sujets par des underscores. Un mot clé peut être associé à un thésaurus. Dans ce cas on associe l'URL au mot clé avec @.	Recommandé	Pas de clé par défaut. Il faut créer soit même ses couples de "clé:valeur". Pour ajouter un sujet : utiliser comme clé le type de mot clé ou bien le thesaurus utilisé. On a ainsi des couples "clé:valeur" de cette forme : clé --> libre valeur --> [mot clé] ou [mot clé@url] clé-libre:mot-clé1,mot-clé2@url
6	Creator	Noms des parties ayant un rôle dans la gestion des jeux de données.	Noms des parties ayant un rôle dans la gestion des jeux de données. Chaque partie citée ici doit apparaître dans un autre tableur : le tableur contact. Une personne peut être citée dans plusieurs catégories. Catégories : - owner pour les personnes auteurs possédant les jeux de données - publisher pour les personnes responsables de la publication des jeux de données - metadata pour les personnes responsables de la gestion des métadonnées - pointOfContact pour les autres contributeurs, les personnes à contacter pour plus d'information. La catégorie owner est celle fortement recommandée.	Recommandé	owner:mail personne 1,mail personne 2_ publisher:mail personne 1_ metadata:mail personne 3_ pointOfContact:mail personne 2,mail personne4
7	Date	Date associée au jeu de données, généralement sa création.	Dates associées au jeu de données.	Recommandé	creation:YYYY-MM-DD ou YYYY
8	Type	Type de données (vecteur, texte, tabulaire,...)	Type de données (vecteur, texte, tabulaire,...)	Recommandé. Si ce champ n'est pas renseigné "dataset" sera choisit par défaut.	generic:dataset
9	Language	Langue des données et fiches de métadonnées.	Langue des données et fiches de métadonnées. La meilleure pratique est d'utiliser l'ISO3 Ppour les deux cas les plus courant (anglais et français) indiquer en abréviation les trois première lettre de la langue en anglais (respectivement : eng et fre).	Recommandé. Si ce champ n'est pas renseigné "eng" est choisit par défaut.	Pas de clé. Mettre directement le code langue. Ex : eng
10	SpatialCoverage	Couverture spatiale des données et système de coordonnées (SRID) correspondant.	Couverture spatiale des données et système de coordonnées (SRID) correspondant. Attention ! On sépare le SRID du WTK avec un point virgule ; lien pour dessiner l'étendue géospatiale : https://arthur-e.github.io/Wicket/sandbox-qmaps3.html	Optionnel	SRID=code;WKT
11	TemporalCoverage	Couverture temporelle des données. Peut être un moment donné ou une période.	Couverture temporelle des données. Peut être un moment donné ou une période.	Optionnel	Pour un moment donné : YYYY-MM-DD ou YYYY 2021-08-01 ou 2021 Pour une période : YYYYYYYY ou YYYY-MM-DDTHH:MinMin:SSZ/YYYY-MM-DDT HH:MinMin:SSZ 2010/2015 ou

-PART IV-

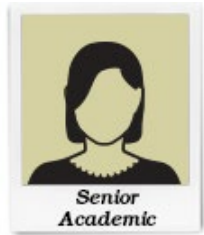
End of the trip !

You have some keys to make a :

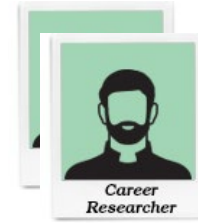
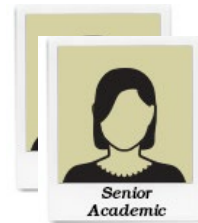


► Data Management Plan (DMP) for Project,

*Platform,
Database,
...*



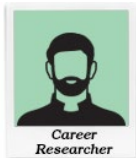
assisted by



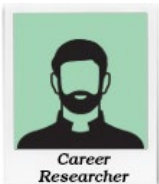
*datamanger
librarian*



► Data Paper



Researcher make homogeneous collection of data on a purpose



Researcher don't have time (or needs) to exploit all measurement data

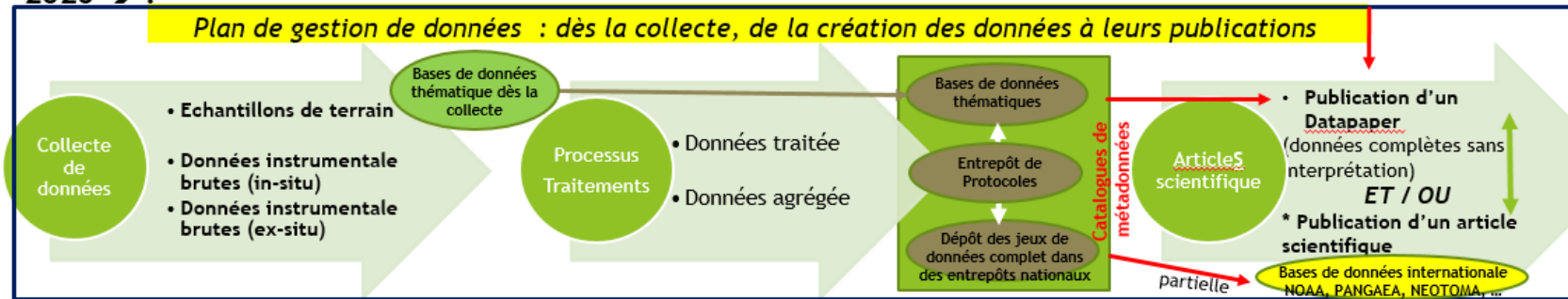
Data Management Plan (DMP) for Project

→ It is a **DOCUMENT** (*deliverable follows a frame*) that **FEEDS** THROUGHOUT the PROJECT/THESE

→ **LINKS** the overall strategy and results for **ALL DATAS** (*raw -> published*)



2020 → ?



→ So what's the plan ?

Start at the beginning, not at the middle or the end !

<https://doranum.fr/plan-gestion-donnees-dmp/>

<https://coop-ist.cirad.fr/gerer-des-donnees/se-familiariser-avec-les-pgd/3-exemple-de-trame-d-un-plan-de-gestion-de-donnees>



Data Management Plan (DMP) for Project

Different Templates for the same result

<https://coop-ist.cirad.fr/gerer-des-donnees/se-familiariser-avec-les-pgd/6-des-modeles-de-plans-de-gestion-de-donnees-en-anglais>



- ▶ In french / english
 - modèle Horizon 2020 english
 - modèle français/anglais de l' ANR
 - modèle INRAe, IRD,...



Or
Template in cloud
collaborative

Collaborative work

PI project + 1 DMP Coordinator
+ All the participant of the project ...

▶ In English

- DMP Horizon 2020 [Guidelines on FAIR Data Management in Horizon 2020: Version 3.0.](#) European Commission, 26 July 2016, 12 p. (version anglaise)
- Digital Curation Centre's DMPonline tool [DMPonline](#) (2018)
- University of California's DMPTool
- ESRC Data management plan: guidance for peer reviewers
- NERC Data Management planning





Data Paper ? *pair-review publication (Rang A) made aware for Raw, Cold, unpublished data*

In DataPaper = NEVER data analysis !

- ✓ No hypothesis
- ✓ No interpretation
- ✓ No derived conclusions



► Data Paper « Data compilation » type

Synthesis / compilation of data scattered according to criteria

Work of grouping coherent (generally published) datas (on the same theme, of the same variables, on the same or similar site, etc...)

➔ But some of them have an intention in the collection

► Data Paper « Observatory » type

Collection protocols, data organization

1 year of photography in a research site
Only 2 years of monitoring a site
Annual data from a platform (with embargo)

► Data Paper « data collection » type

Provision of data collected but even used in a project, PhD, ...

High metadata level



avoid falling into « the long tail of » Darkdata

<https://doranum.fr/data-paper-data-journal/>

<https://coop-ist.cirad.fr/gerer-des-donnees/rediger-un-data-paper/1-qu-est-ce-qu-un-data-paper>



Data Paper (2)

→ *Nature* Data Repository Guidance [Scientific Data repository](#)
→ A list of peer-review Data papers : [PDF](#) (CIRAD, 2018)

NEVER ! data analysis
✓ No hypothesis
✓ No interpretation
✓ No derived conclusions

- ▶ *Data paper* inform the scientific community of the **availability of these datasets** and their **potential** for future use

Describes scientific data :

→ circumstances and methods of its collection.

→ Quality : technical and statistical analyses for validating data

- ▶ *Data paper* shows the **originality and scope** of the dataset it describes

Journals publishing data papers are particularly interested in

→ the scope of the data submitted

→ The potential for reuse by other scientists (major argument)

High metadata level



- describes a set of scientific data
- using precise informations (metadataS generalist, interdisciplinary & thematics, quality), scope, potential ...
- data storage in a permanent warehouse (DOI)



Check-list « What can I do ? »

At the beginning (or whenever possible) of your PhD or Project



Check-list (1)

ORCID



Clean your personnel id network
→ only one ORCID or one IDHal

1- Think beforehand about all this actions :

- good hierarchical folder/file tree (raw, processed, a aggregated). Think you could need to share it at the end
- properly name your folders/files with consistent name/date, version,...
- search the **metadata standards of your thematics** (sample, geospatial, climate, chemistry, etc.) and **vocabularies** (from generic to specific) to describe your variables



DON'T WAIT if :

- ERC, ANR(>2020)= DMP mandatory
→ Project manager take the template of a DMP and think about a plan !
(Take a good coordinator & collaborators to help you along the project)
- Anticipation of a Datapaper ("observatory" or "collection" type)
→ prepare files as soon as they are collected in connection with the acquisition platform & metadata standard

-Senior Researcher
-Data Scientist
Help The young !



Workshop

with your expert community to decide your specific Guideline & Template data

(Before see thematic's Groups in the Research Data Alliance)

Thematics interest groups



RESEARCH DATA ALLIANCE

<https://coop-ist.cirad.fr/gerer-des-donnees/rediger-un-pgd/6-garantir-la-comprehension-et-l-accessibilite-des-donnees>

<https://doranum.fr/stockage-archivage/comment-nommer-fichiers/>

<https://dataservices.gfz-potsdam.de/portal/phd.html>



Check-list (2)

- 2. SAMPLE :** Properly name and describe your “Parent” and “Child” Samples (sub-samples)
 - « Sampling plan” (and update if changed) => DMP
 - Deposit in sample repository = get an IGSN identifier + search by catalogues
For Cores => [French Cyber-core repository](#) , For other samples [CNRS AA IGSN](#), [SESAR](#),
- 3. PROTOCOLS :** well describe your measurement protocols
(txt procols, sheet of laboratory notebook, Electronic labnotebook (ELN), publication of a protocols (article or [Protocols repository](#))
- 4. ANALYTICAL DATA :** well structure your data in tables → **MAKE TEMPLATES MODEL REPETABLE**
 - *Sheet Tab « Read-me »* : with a **data dictionary** (variable name + unit of measure + vocabularies + definitions), define unmeasured values (9999, -9999, n/a, value below the measurement threshold ,...), suppressed value,
 - *Sheet Tab « rawdata »* (if necessary, otherwise link to the permanent repository)
 - *Sheet Tab « data »* : **only variables values, any information other than definine in the Read-me**
- 5. Versioning** your data as your work progresses (with date : AAAA-MM-DD_XXXXX.xls)



Check-list (3) : 2 cases

► Scientific Article Publication :

- ✓ I submit an article with the link to [datas](#) warehouse (samples, campaigns, protocols, data)

IN THE POSTPRINT PAPER (some editors ask links in the paper abstract)

- ✓ Data deposit with complete [data dictionary and metadata](#) (standard & vocabulary)

► Datapaper Publication ? Why ?



- My research is an observatory (eg. without Database)
- Part of my research will lead the creation of an cross-checking of data
- I don't have time & needs & expertise to use all the dataset of an instrumental acquisition
- I don't have time to write an article

Don't aliment the long tail of data anymore !

➔ I make a Datapaper

► Supply the infrastructure at your disposal (*IT Infrastructures available in France*)

- ❑ **Metadata catalogues** : OSUs Cat.⁽¹⁾, [Réseau des Zones Ateliers \(LTER\)](#)⁽¹⁾, [Pôles DataTERRA](#) (ODATIS, OZCAR-THEIA, FORM@Terre, AERIS)⁽¹⁾, Pole national de la Biodiv. ([PNDB](#))⁽²⁾, CNRS-INEE ([IndoRES-MD Catalogue](#))⁽¹⁾
- ❑ **Thematic Database (OSU, International,...)**
- ❑ **National Data repository** : Dataverses IRD, INRAe, OSU (Lorraine), InDORES-data (INEE) *Nov2021*, Entrepôt national fédéré de la Recherche (*mars 2022*)

Conclusion of this trip ?

Faut pas trop s'en F.A.I.R ?





A lot of bibliography, websites' tools throughout the présentation and,

1. Gil, Y. *et al.* Toward the Geoscience Paper of the Future: Best practices for documenting and sharing research from data to software to provenance. *Earth and Space Science* **3**, 388-415 (2016).
2. Khider, D. *et al.* PaCTS 1.0: A Crowdsourced Reporting Standard for Paleoclimate Data. *Paleoceanography and Paleoclimatology* **34**, 1570-1596 (2019).
3. Lehnert, K., Elger, K. & Ulbricht, D. IGSN - More than Unique Identification of Samples. (2020) doi:10.5281/zenodo.3776495.
4. McNutt, M. *et al.* Liberating field science samples and data. *Science* **351**, 1024-1026 (2016).
5. Morrill, C. *et al.* The Paleoenvironmental Standard Terms (PaST) Thesaurus: Standardizing Heterogeneous Variables in Paleoscience. *Paleoceanography and Paleoclimatology* **36**, e2020PA004193 (2021).
6. Peng, G. *et al.* *International Community Guidelines for Sharing and Reusing Quality Information of Individual Earth Science Datasets.* <https://osf.io/xsu4p/> (2021) doi:10.31219/osf.io/xsu4p.
7. Monya Baker . Five keys to writing a reproducible lab protocol Effective sharing of experimental methods is crucial to ensuring that others can repeat results. An abundance of tools is available to help. Monya Baker. (*2021) <https://www.nature.com/articles/d41586-021-02428-3>