

**Principes FAIR et
cycle de vie des données**

15 février 2024

→ Institut de l'information scientifique et technique
DVDR | Service Formation-DoRANum

cnrs

VALORISER LES
DONNÉES DE LA
RECHERCHE

Ce webinaire a pour objectif de vous expliquer à quoi correspondent les principes FAIR. Ces derniers s'appliquent tout au long du cycle de vie des données de recherche. Celui-ci sera détaillé en expliquant à chaque étape les bonnes pratiques à mettre en place pour respecter les principes FAIR.

Sommaire

01 Définition des données de recherche

02 Principes FAIR

03 Cycle de vie des données de recherche

Étape 1 : planification (PDG / DMP)

Étape 2 : création, collecte et description des données de recherche

Étape 3 : traitement et analyse, stockage sécurisé durant le projet

Étape 4 : partage, diffusion des données, dépôt dans un entrepôt

Étape 5 : archivage pérenne

Étape 6 : réutilisation et valorisation des données

04 À retenir

01

Définition des données de recherche

Définition des données de recherche

« Les données de la recherche sont définies comme des **enregistrements factuels** (chiffres, textes, images et sons), qui sont utilisés comme **sources principales pour la recherche** scientifique et sont généralement reconnus par la communauté scientifique comme **nécessaires pour valider les résultats** de recherche. »

([Définition l'OCDE](#))

Si vous voulez en savoir plus : [What is data?](#)

Plusieurs définitions des données de la recherche existent.

« Les données de la recherche sont définies comme des **enregistrements factuels** (chiffres, textes, images et sons), qui sont utilisés comme **sources principales pour la recherche scientifique** et sont généralement reconnus par la communauté scientifique comme **nécessaires pour valider les résultats de recherche**. »

(OCDE, *Organisation de Coopération et de Développement Économiques. Principes et lignes directrices de l'OCDE pour l'accès aux données de la recherche financée sur fonds publics*. 2007. <https://doi.org/10.1787/9789264034020-en-fr>)

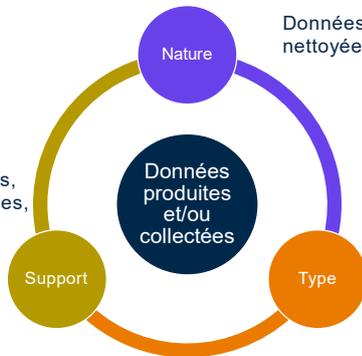
On peut également citer la définition issue de l'ARDC - Australian Research Data Commons (traduction Inist-CNRS) :

« Le terme de données de la recherche désigne les données sous forme de faits, d'observations, d'images, de résultats de programmes informatiques, d'enregistrements, de mesures ou d'expériences sur lesquels un argument, une théorie, un test ou une hypothèse, ou un autre produit de la recherche est basé [...]. Les données peuvent être numériques, descriptives, visuelles ou tactiles. Elles peuvent être brutes, nettoyées ou traitées, et peuvent être conservées dans tout format ou support [...]. » (ARDC, *Australian Research Data Commons. What Is Research Data*. 2019)

Définition des données de recherche



Documents électroniques, programmes informatiques, carnets de laboratoire, supports papier, logiciels...



Données brutes, dérivées, formatées, nettoyées, primaires, secondaires, traitées...

Archives, audio, vidéos, bases de données, codes sources, données géospatiales, images, photographies, langages de programmation, matérielles et physiques, modèles, visualisations, 3D, numériques, textuelles, numérisations, scans, qualitatives, quantitatives, statistiques...

Selon leur contexte de création (capture ou production), leur exploitation, leur analyse et les traitements qu'elles subissent, les données de recherche peuvent être :

- de différentes **natures** : données brutes, dérivées, formatées, nettoyées, primaires, secondaires, traitées....
- de tous **types** : audio, vidéos, images, bases de données, codes sources, données géospatiales, photographies, langages de programmation, matérielles et physiques, modèles, visualisations, 3D, numériques, textuelles, numérisations, scans, qualitatives, quantitatives, statistiques...
- contenues dans divers **supports** : documents électroniques, programmes informatiques, carnets de laboratoire, supports papier, logiciels...

Sources :

- Rivet Alain, Bachèlerie Marie-Laure, Denis-Meyere Auriane, Tisserand Delphine. *Traçabilité des activités de recherche et gestion des connaissances. Guide pratique de mise en place.* 2018.
<https://qualite-en-recherche.cnrs.fr/gt/tracabilite-des-activites-de-recherche/>
- UNIL, Université de Lausanne. *Les données de recherche.*
<https://www.unil.ch/openscience/fr/home/menuinst/open-research-data/les-donnees-de-recherche.html>

Ressource :

DoRANum. *L'origine et la description des données de la recherche.* 15 mars 2022.

https://doranum.fr/plan-gestion-donnees-dmp/origine-description-donnees-recherche_10_13143_e9zh-w908/

Exemples de données de recherche



Photo d'une peinture pariétale représentant une scène de chasse

Code	Code provenance	Code description	Pays d'origine	Nom de la zone dans laquelle elle se trouve	Autre nom	Coordonnées géographiques (lat/long)	Altitude (m)	P (mm)	substrat	Poids de 1000 graines (g)	% graines viables (méthode + triage)	Présence	01-4800	01-4800	01-4800	01-4800
7	0078		Espagne	Andalucía		36°20'N 0°50'W	631		Lunuel	653,6	4,6	x	x	x	x	x
8	0079		Espagne	Cataluña		41°04'N 2°06'W	710		Carbanel	719,2	5,7	x	x	x	x	x
9	0079		Espagne	Madrid		41°30'N 4°20'W	875		Carbanel	588,2	12,1	x	x	x	x	x
10	0080		Espagne	Cataluña		40°29'N 4°00'W	900		Carbanel	772,2	2,9	x	x	x	x	x
11	0081		Grèce	Sabouni / Oraklioti		40°14'N 23°34'W	400		granite	689,7	0,4	x	x	x	x	x
12	0082		Grèce	Melochi		39°50'N 23°18'W				513,9	0	x	x	x	x	x
13	0083		Italie	Fregene (ST)		42°25'N 11°17'W	250		gypse	714,3	12,7	x	x	x	x	x
14	0084		Italie	Majano		42°42'N 10°17'W	530		gypse	1026,4	9,1	x	x	x	x	x
15	0088		Liban	Beit Moucar		34°29'N 35°31'W	1300-1400		gribois	1041,9	6,2	x	x	x	x	x
16	0088		Liban	Komel		33°44'N 35°13'W	1400		gribois	788,6	9,1	x	x	x	x	x
17	0088		Liban	Qadieh		33°39'N 35°02'W	1000		gribois	784,3	7,8	x	x	x	x	x
18	0087		Israëlle	Blakness		33°17'N 35°02'W	900			812,4	4,8	x	x	x	x	x
19	0088		Liban	Hamir		33°06'N 35°02'W	1000			7	0	x	x	x	x	x
20	0088		Liban	Zahle		33°06'N 35°02'W	1000			641	2,7	x	x	x	x	x
21	0088		Liban	Hamir		33°06'N 35°02'W	1000			775,2	6,1	x	x	x	x	x
22	0088		Liban	Hamir		33°06'N 35°02'W	1000			488,4	11,2	x	x	x	x	x
23	0088		Liban	Hamir		33°06'N 35°02'W	1000			527,7	6,1	x	x	x	x	x
24	0088		Liban	Hamir		33°06'N 35°02'W	1000			720,9	12,2	x	x	x	x	x
25	0021		Turquie	Alayun		41°11'N 41°15'W	225			574,7	4,4	x	x	x	x	x
26	0044		France	Mayens (G)		47°10'N 0°36'W	600		sablon	918,8	6,7	x	x	x	x	x
27	0044		France	Mayens (G)		47°10'N 0°36'W	600		sablon	710	0,2	x	x	x	x	x
28	0044		France	Mayens (G)		47°10'N 0°36'W	600		sablon	710	0,2	x	x	x	x	x
29	0044		France	Mayens (G)		47°10'N 0°36'W	600		sablon	509,9	0	x	x	x	x	x
30	0044		France	Mayens (G)		47°10'N 0°36'W	600		sablon	517,5	7,1	x	x	x	x	x
31	0020		Turquie	Ballisara		39°40'N 27°30'W	370			603	11,8	x	x	x	x	x
32	0020		Turquie	Ballisara		39°40'N 27°30'W	370			448,4	10,6	x	x	x	x	x
33	0020		Turquie	Ballisara		39°40'N 27°30'W	370			950,3	14	x	x	x	x	x
34	0020		Turquie	Ballisara		39°40'N 27°30'W	370			466,2	18,2	x	x	x	x	x
35	0020		Turquie	Ballisara		39°40'N 27°30'W	370			919,7	6,1	x	x	x	x	x
36	0020		Turquie	Ballisara		39°40'N 27°30'W	370			712,5	6,6	x	x	x	x	x
37	0020		Turquie	Ballisara		39°40'N 27°30'W	370			561,5	5,2	x	x	x	x	x
38	0020		Turquie	Ballisara		39°40'N 27°30'W	370			665,6	10,6	x	x	x	x	x
39	0020		Turquie	Ballisara		39°40'N 27°30'W	370			662,3	6,7	x	x	x	x	x
40	0020		Turquie	Ballisara		39°40'N 27°30'W	370			819,7	16	x	x	x	x	x
41	0020		Turquie	Ballisara		39°40'N 27°30'W	370			680,3	16,2	x	x	x	x	x
42	0020		Turquie	Ballisara		39°40'N 27°30'W	370			513,3	23,9	x	x	x	x	x

Données sur les traits phénotypiques de *Pinus pinea* mesurés dans les jardins communs du réseau INRAE de génétique forestière pour la recherche et l'expérimentation (GEN4X)



Voici des exemples concrets de données de recherche dans différentes disciplines :

- Cette photo d'une peinture pariétale représentant une scène de chasse dans le désert du Wadi Ram, en Jordanie est un exemple de donnée pour un chercheur en archéozoologie qui étudie le rôle du monde animal au sein des sociétés humaines sans écriture.

Source de l'image : <https://pixabay.com/fr/photos/jordan-wadi-d%C3%A9sert-sable-paysage-4158477/>

- Données sur les traits phénotypiques de *Pinus pinea* mesurés dans les jardins communs du réseau INRAE de génétique forestière pour la recherche et l'expérimentation (GEN4X). Ce jeu de données est composé d'observations telles que la hauteur totale, la circonférence, le nombre de cônes, la survie et la fourchaison, réalisées jusqu'en 2012 dans cinq tests de provenance de *Pinus pinea* du réseau français de jardins communs GEN4X (réseau de génétique forestière pour la recherche et l'expérimentation).

Fady Bruno, Vauthier Denis, Michotey Celia. Phenotypic trait dataset of *Pinus pinea* from five common gardens of the INRAE GEN4X network. 2021. <https://doi.org/10.15454/MQKQ10>

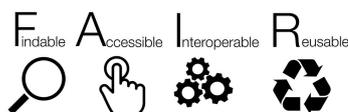
02

Principes FAIR

Principes FAIR

4 principes à respecter pour garantir une utilisation optimale des données de recherche et des métadonnées associées, à la fois **par les hommes et par les machines**

- **F** (Findable) = Facile à trouver
- **A** (Accessible) = Accessible
- **I** (Interoperable) = Interopérable
- **R** (Reusable) = Réutilisable



Principes admis par les différentes communautés scientifiques au niveau international, ainsi que par les financeurs (ex : Commission européenne, ANR, etc.)

Applicables tout au long du cycle de vie des données

Tout au long de cette présentation des petites bulles vous permettront de repérer à quel(s) principe(s) FAIR correspondent principalement les bonnes pratiques décrites.

Ressource :

DoRANum. *Les principes FAIR*. 4 décembre 2019. https://doranum.fr/enjeux-benefices/principes-fair_10_13143_z7s6-ed26/

Principe F

F

Le principe F facilite la découverte des données par les humains et les systèmes informatiques :

- Catalogues ou entrepôts
- Identifiant pérennes
- Métadonnées riches



Le principe F signifie que vos données doivent être faciles à trouver par des humains et par des systèmes informatiques.

Bonnes pratiques :

- Déposer vos jeux de données dans un entrepôt de données
- Leur attribuer un identifiant pérenne
- Les décrire précisément à l'aide de métadonnées

Source :

- *INRAe, Institut national de recherche pour l'agriculture, l'alimentation et l'environnement. Produire des données FAIR. <https://science-ouverte.inrae.fr/les-donnees-et-le-numerique-scientifiques/produire-des-donnees-fair>*

Principe A

Le principe A encourage à stocker durablement les données et les métadonnées et à faciliter leur accès et/ou leur téléchargement, en spécifiant les conditions d'accès (accès ouvert ou restreint)

- Protocoles standardisés et ouverts
- Authentification et autorisation si besoin
- Accès permanent aux métadonnées



Accessible

Le principe A encourage à stocker durablement les données et les métadonnées et à faciliter leur accès et/ou leur téléchargement, en spécifiant les conditions d'accès (accès ouvert ou restreint).

Bonnes pratiques :

- Protocoles de communication standards : les données doivent être rendues accessibles soit via une URL, soit via un téléchargement de fichier(s). L'idée est de ne pas avoir besoin de recourir à des outils spécialisés ou propriétaires pour pouvoir accéder aux données, exemples : Skype, Microsoft Exchange...
- Préférer le dépôt dans des entrepôts certifiés qui proposent un accès ouvert.
- Définir qui a le droit d'accéder aux (méta)données, par quelle procédure (ouverture totale ou demande d'accès par authentification ou demande d'autorisation)
- Les métadonnées restent accessibles même si les données ne le sont pas ou plus.

Sources :

- INRAe, Institut national de recherche pour l'agriculture, l'alimentation et l'environnement. Produire des données FAIR. <https://science-ouverte.inrae.fr/les-donnees-et-le-numerique-scientifiques/produire-des-donnees-fair>
- Université Paris-Saclay. Produire des données FAIR. <https://www.universite-paris-saclay.fr/recherche/science-ouverte/produire-des-donnees-fair>

Principe I

Le principe I signifie que les données sont téléchargeables, utilisables, intelligibles et combinables avec d'autres données, par des humains et par des systèmes informatiques

- Formats ouverts
- Lien vers d'autres (méta)données
- Vocabulaires contrôlés (standards)

Comment ouvre-t-on
un fichier .xzq ?



Interopérable

Le principe I signifie que les données sont téléchargeables, utilisables, intelligibles et combinables avec d'autres données, par des humains et des systèmes informatiques.

Bonnes pratiques :

- Utiliser des formats ouverts et indépendants.
- Utiliser des vocabulaires standards, des thésaurus, des ontologies... Le but étant de « parler le même langage ».
- Contextualiser les données en indiquant les liens vers d'autres données et métadonnées ou vers des publications.

Source :

- INRAe, Institut national de recherche pour l'agriculture, l'alimentation et l'environnement. Produire des données FAIR. <https://science-ouverte.inrae.fr/les-donnees-et-le-numerique-scientifiques/produire-des-donnees-fair>

Principe R

Le principe R encourage à rendre vos données réutilisables pour de futures recherches ou d'autres finalités (enseignement, innovation, reproduction et transparence de la science) :

- Description riche
- Provenance
- Standards communautaires
- Licences



Le principe R encourage à rendre vos données réutilisables pour de futures recherches ou d'autres finalités (enseignement, innovation, reproduction et transparence de la science). Les métadonnées et les données doivent être bien décrites afin de pouvoir être reproduites et/ou combinées dans différents contextes.

Bonnes pratiques :

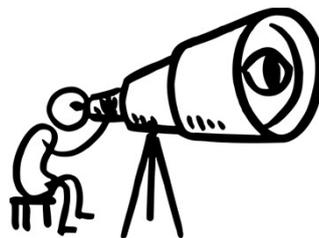
- Les données sont très précisément décrites avec le plus d'informations possibles (métadonnées).
- Leur provenance est renseignée.
- Elles répondent à des standards communautaires pertinents pour le domaine.
- Il est important de leur attribuer une licence d'utilisation.

Source :

- *INRAe, Institut national de recherche pour l'agriculture, l'alimentation et l'environnement. Produire des données FAIR. <https://science-ouverte.inrae.fr/les-donnees-et-le-numerique-scientifiques/produire-des-donnees-fair>*

Qu'apportent les Principes FAIR ? Pour le chercheur

- Mise en place de bonnes pratiques
- Gain de temps dans la gestion d'un projet de recherche
- Meilleure visibilité et citation
- Favorise les collaborations



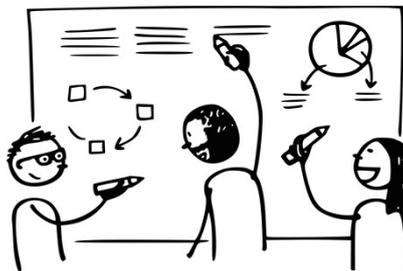
Afin de mieux saisir à quoi servent les principes FAIR, voici une liste de ce qu'ils apportent concrètement au chercheur ainsi qu'à la communauté scientifique.

Pour le chercheur

- Bonnes pratiques à adopter.
- C'est aussi un gain de temps dans la gestion du projet de recherche.
- Le chercheur gagne en visibilité.
- Ses données sont plus accessibles.
- L'intégrité scientifique est facilitée.
- Cela favorise les collaborations.

Qu'apportent les Principes FAIR ? Pour la communauté scientifique

- Accès facile à des données publiques et réutilisables
- Gain de temps et d'argent : ne pas recréer des données déjà existantes
- Meilleure reproductibilité de la recherche



Pour la communauté scientifique

- L'accès à des données publiques et réutilisables est facilité, ainsi que l'accès à des corpus utiles pour d'autres domaines.
- L'Interopérabilité des données est meilleure.
- Il y a également un gain de temps et d'argent puisque des données déjà existantes ne sont pas recréées.
- Cela permet une meilleure reproductibilité de la recherche.

03

Cycle de vie des données de recherche

Cycle de vie des données de recherche

C'est l'ensemble des étapes

- de gestion
- de conservation
- de diffusion des données de recherche associées aux activités de recherche



« Le cycle de vie des données de la recherche est l'ensemble des étapes de gestion, conservation, diffusion et réutilisation des données scientifiques, associées aux activités de recherche. »

Deboin Marie-Claude. Découvrir de nouveaux métiers liés aux données de la recherche. CIRAD. 5 p. 5 octobre 2018. <https://doi.org/10.18167/coopist/0061>

Nous allons détailler les 6 étapes dans notre présentation.

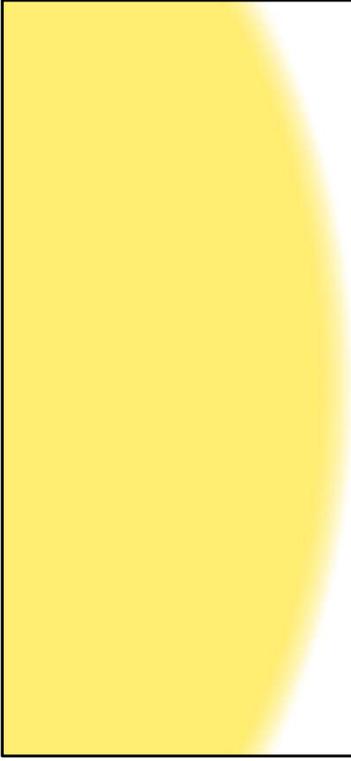
Source de l'image :

D'après Research data lifecycle – UK Data Service.

Ressource :

DoRANum. Le cycle de vie des données de recherche. 4 février 2021.

https://doranum.fr/enjeux-benefices/le-cycle-de-vie-des-donnees-de-recherche_10_13143_gzj2-j593/

A large yellow shape on the left side of the slide, resembling a quarter-circle or a stylized 'C' shape, with a gradient effect.

Étape 1 : Planification (PGD/DMP)

Place dans le cycle de vie des données



Le PGD qu'est-ce que c'est ?



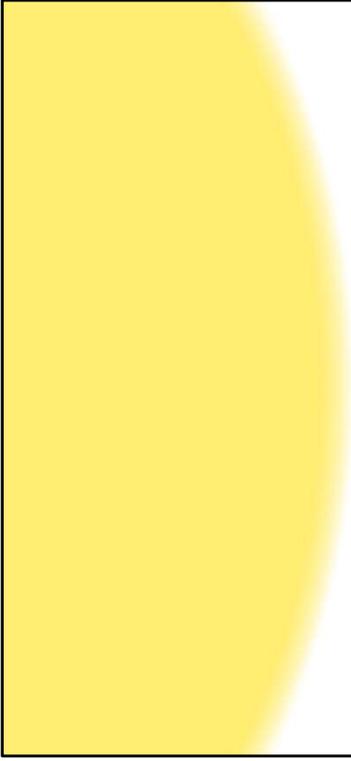
PGD (Plan de Gestion des Données) = DMP (Data Management Plan)

- Explique comment sont gérées les données, depuis leur création ou collecte, jusqu'à leur partage et leur archivage
- Aide à organiser et à anticiper toutes les étapes du cycle de vie des données
- Permet de se poser les bonnes questions pour rendre les données FAIR
- Évolutif : rédaction commençant dès le début du projet, mais mise à jour tout au long du projet

- Le DMP ou PGD est un document de quelques pages, qui permet de définir et de mettre en place les actions indispensables pour gérer les données qui vont être collectées ou produites au cours d'un projet.
- Il explique comment sont gérées les données depuis leur création ou collecte jusqu'à leur partage et leur archivage.
- Le DMP est évolutif : sa rédaction démarre dès le début d'un projet, mais il va être mis à jour régulièrement en fonction de l'évolution du projet. Vous n'avez pas à répondre à toutes les questions du DMP dans sa phase initiale mais il est bon de réfléchir à tous les points abordés dans le DMP. Ce document est là pour faciliter la gestion des données de recherche de votre projet.
- Surtout, le DMP aide à organiser et anticiper toutes les étapes du cycle de vie de la donnée.

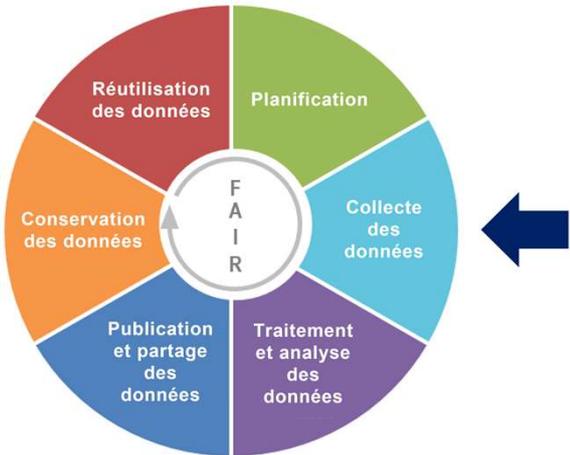
Ressource :

- DoRANum. Cours introductif sur le Plan de Gestion de Données (PGD). 27 avril 2022. https://doranum.fr/plan-gestion-donnees-dmp/cours-introductif-sur-le-plan-de-gestion-de-donnees_10_13143_t3j4-vn03/

A large yellow shape on the left side of the slide, resembling a quarter-circle or a soft-edged rectangle, with a gradient from light to dark yellow.

Étape 2 : Création, collecte et description des données de recherche

Place dans le cycle de vie des données



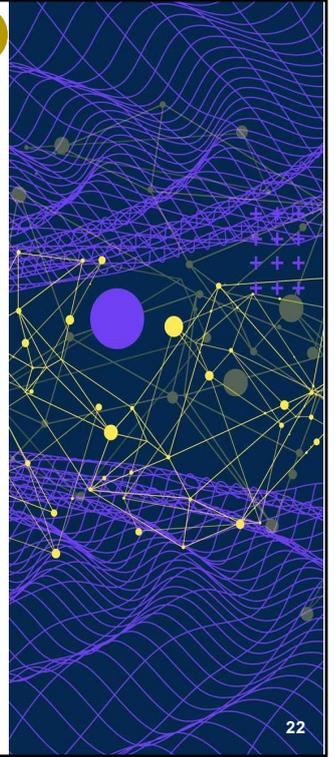
Description des données de recherche



Produites ou recueillies



Préexistantes



Selon le projet, les données de la recherche peuvent être :

- **Produites ou recueillies** : ce sont des données créées spécifiquement durant le projet. Elles sont élaborées, générées lors d'activités de recherche (collectes sur le terrain, observations, des mesures...).
- **Préexistantes** : ce sont des données déjà existantes (provenant de corpus, archives...) qui sont utilisées pour le projet. Les données utilisées peuvent avoir été recueillies initialement dans un autre contexte que celui de la recherche mais elles sont utilisées comme données de recherche dans le cadre du projet.

Selon que les données sont produites ou préexistantes, leur description ne contiendra pas les mêmes informations.

Préparation et documentation des données

Dans le PGD, il faudra indiquer de manière précise **quelles méthodes** sont utilisées pour recueillir ou produire les données :

Dans le cas de **données préexistantes** :

- leur provenance (corpus, archives...)
- sur quels critères elles ont été sélectionnées
- les conditions de réutilisations préexistantes de ces données

Dans le cas de **données produites ou recueillies** : (observations, mesures, etc.) :

- le contexte de création
- les méthodes utilisées
- les protocoles suivis ou établis
- les contrôles qualité mis en place

Il est important de **bien documenter les données** de manière à ce qu'elles soient faciles à trouver et réutilisables.

Ces informations devront apparaître dans les métadonnées et/ou dans un fichier Readme. Elles sont indispensables pour que d'autres chercheurs puissent reproduire les résultats (reproductibilité et répliquabilité).

La **reproductibilité** est la capacité, par une équipe différente, de reproduire une expérience, sans se fier au dispositif expérimental et aux codes logiciels développés par l'équipe d'origine.

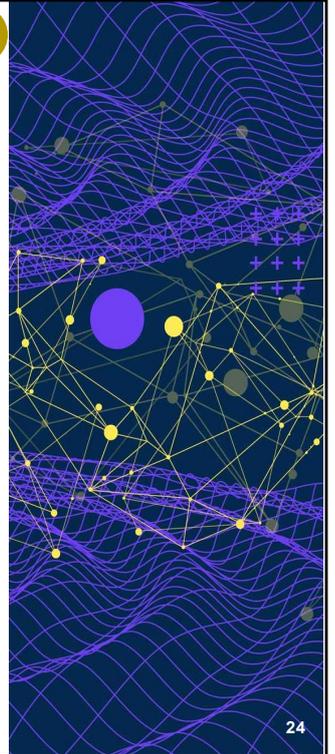
La **répliquabilité** est la capacité, par une équipe différente, de reproduire une expérience en ré-utilisant le même dispositif expérimental décrit (y compris les codes logiciels).

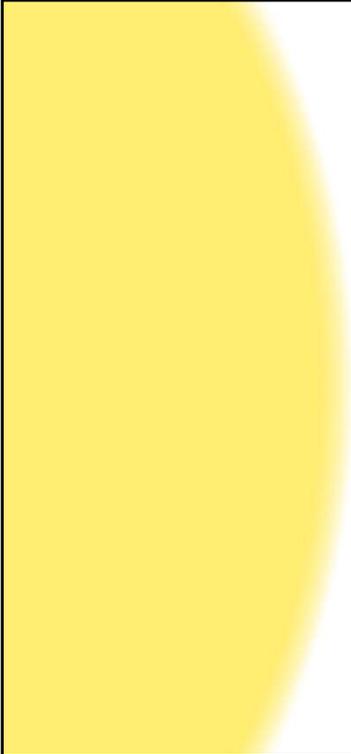
Préparation et documentation des données

Il est impératif de **bien préparer et documenter ses données** afin d'optimiser le stockage, le partage, l'archivage et la réutilisation



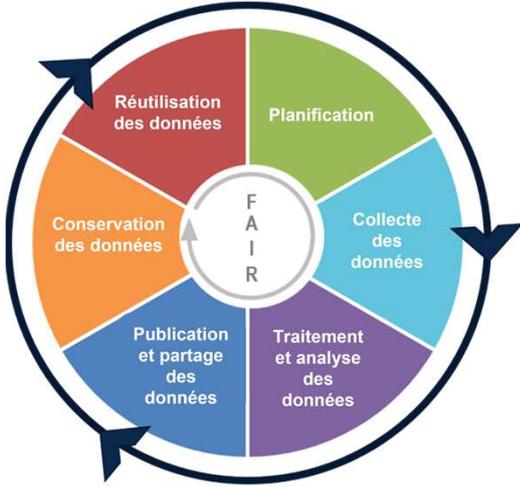
Attention aux **données sensibles, personnelles ou confidentielles** : prendre les précautions nécessaires afin de respecter les règles juridiques et éthiques en vigueur





Aspects juridiques et éthiques

Place dans le cycle de vie des données



Droits et obligations du chercheur

Cette étape est cruciale car elle détermine la latitude dont le chercheur disposera ensuite pour publier, diffuser et communiquer ses données et les résultats de ses recherches

Dès le début du projet, au moment de la collecte et de la production des données, le chercheur doit être vigilant concernant ses droits et obligations

Exemples:

- Dans le cas d'une interview ou de prises de son ou de vue, il doit recueillir le consentement écrit des personnes concernées
- Dans le cas de consultation de données d'archives, quels sont les droits afférents
- Dans le cas de collecte d'objets archéologiques, quels sont les droits liés au pays de collecte.

Ressource :

- Delplanque Catherine, Lamrini Nawale, Leclère Fabrice, Maurel Lionel, et al. *Guide : Règlement Général pour la Protection des Données*. 2019.
<https://www.u-plum.fr/guide-reglement-general-pour-la-protection-des-donnees/>

Propriété intellectuelle des données

Accompagnement par un juriste recommandé pour déterminer qui a le droit d'accéder aux données

Règle générale

- Attribution de la propriété intellectuelle des données à l'établissement de tutelle des producteurs de données



Si partenariat

- Nécessité d'établir au préalable un **accord de consortium**
- Cas de collaborations entre secteurs public et privé (déploiement industriel, commercialisation...)
- Cas de collaborations internationales

Penser à préciser la propriété des données et les responsabilités dans le PGD

Attention, **le PGD n'a aucune valeur juridique.**

Il ne s'agit pas du même droit que les publications (droit d'auteur principalement).

Les données relèvent d'un régime lié au droit des bases de données. Dans ce cas, **le droit de propriété appartient légalement au « producteur » de la base de données**, compris au sens de la personne qui réalise l'investissement financier et matériel nécessaire à la constitution de la base. Il s'agira donc en général de **l'établissement de tutelle des chercheurs qui sera considéré comme le titulaire effectif du droit de propriété.**

Mais si ce droit existe formellement, il ne peut plus être opposé aux droits des ré-utilisateurs des données (principe d'ouverture des données). En effet, la loi pour une République numérique a explicitement « neutralisé » le droit des bases de données des administrations pour faire primer le principe de libre réutilisation. Il en résulte que **les données produites par les chercheurs sont bien comprises dans le principe d'ouverture par défaut.**

Ressources :

- DoRANum. *Questions juridiques liées aux données de recherche : interview de Lionel Maurel.* <https://doranum.fr/aspects-juridiques-ethiques/questions-juridiques-liees-aux-donnees-de-la-recherche-10-13143-xjgm-hb78/>
- Ginouvès Véronique, Gras Isabelle, et al. *La diffusion numérique des données en SHS – Guide des bonnes pratiques éthiques et juridiques.* <https://hal-amu.archives-ouvertes.fr/page/guide-de-bonnes-pratiques>
- Bordignon Frédérique, Boistel Romain, Du Pasquier Delphine. *Qui a les droits, quelles obligations ?* https://espacechercheurs.enpc.fr/sites/default/files/logigramme_a_plat_2.pdf



Ressources :

- *Becard Nicolas, Castets-Renard Céline, Chassang Gauthier, Dantant Martin, Freyt-Caffin Laurence, Gandon Nathalie, Martin Caroline, Martelletti Andrea, Mendoza-Caminade Alexandra, Morcrette Nathalie, Neirac Claire. Ouverture des données de recherche. Guide d'analyse juridique en France. Décembre 2017.*
<http://dx.doi.org/10.15454/1.481273124091092E12>
- *Bordignon Frédérique, Boistel Romain, Du Pasquier Delphine. Qui a les droits, quelles obligations ?*
https://espacechercheurs.enpc.fr/sites/default/files/logigramme_a_plat_2.pdf
- *Ginouvès Véronique, Gras Isabelle, et al. La diffusion numérique des données en SHS – Guide des bonnes pratiques éthiques et juridiques.* <https://hal-amu.archives-ouvertes.fr/page/guide-de-bonnes-pratiques>

Communicabilité des données



La communicabilité des données peut être conditionnée par

- la nature ou le type des données
- l'origine des données
- leur (ré)utilisation

Elle peut être **empêchée temporairement ou définitivement**

Toute restriction doit être expliquée dans le PGD

Communication obligatoire pour certaines disciplines	Communication sous conditions	Communication interdite par principe
<ul style="list-style-type: none"> • Données géographiques • Données environnementales... 	<ul style="list-style-type: none"> • Données protégées par le droit d'auteur ou par contrat • Données personnelles • Statistiques... 	<ul style="list-style-type: none"> • Secrets professionnels • Secrets défense • Sécurité de l'établissement...

Données de différentes **natures** : brutes, dérivées, formatées, nettoyées, primaires, secondaires, traitées...

Données de différents **types** : archives, audios, vidéos, bases de données, codes sources, données géospatiales, images, photographies, langages de programmation, matérielles et physiques, modèles, visualisations, 3D, numériques, textuelles, numérisations, scans, qualitatives, quantitatives, statistiques...

Éthique des données de recherche



- Dans le cas de données devant respecter des **règles d'éthique** particulières, se référer aux normes, chartes, déclarations, codes, politiques établis dans son domaine...
- Recourir à un **comité d'éthique**, si besoin

L'**éthique** nous invite à réfléchir aux valeurs qui motivent nos actes et à leurs conséquences et fait appel à notre sens moral et à celui de notre responsabilité.

La **déontologie** réunit les devoirs et obligations imposés à une profession, une fonction ou une responsabilité.

L'**intégrité scientifique** concerne la « bonne » conduite des pratiques de recherche. Elle désigne l'ensemble des règles et valeurs qui doivent régir l'activité de recherche pour en garantir le caractère honnête et scientifiquement rigoureux. Le respect de ces règles est une condition indispensable du maintien du lien de confiance accordée par la société aux acteurs de la recherche.

Sources :

- CNRS, Centre national de la recherche scientifique. *Responsabilité de recherche*. 8 janvier 2024. <https://www.cnrs.fr/fr/le-cnrs/responsabilites/responsabilite-de-recherche>
- ANR. *L'intégrité scientifique*. <https://anr.fr/fr/lanr/engagements/lintegrite-scientifique/>

Des normes, chartes, déclarations, codes et politiques en éthique et en intégrité scientifique encadrent les pratiques pour l'ensemble de nombreux acteurs internationaux de la recherche.

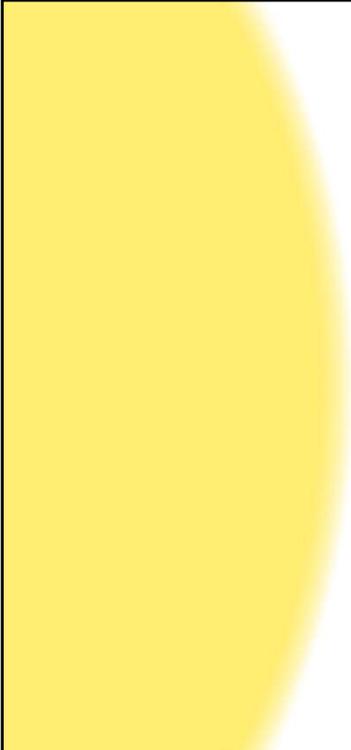
Ressources :

- ANR, Agence nationale de la recherche. *Politique en matière d'éthique et d'intégrité scientifique*. Août 2014. <https://anr.fr/fileadmin/documents/2014/Politique-ethique-integrite-scientifique-aout-2014.pdf>
- Ginouvès Véronique, Gras Isabelle et al. *La diffusion numérique des données en SHS – Guide des bonnes pratiques éthiques et juridiques* : <https://hal-amu.archives-ouvertes.fr/page/guide-de-bonnes-pratiques>
- *Questions éthique et droit en SHS* : <https://ethiquedroit.hypotheses.org/>

En résumé

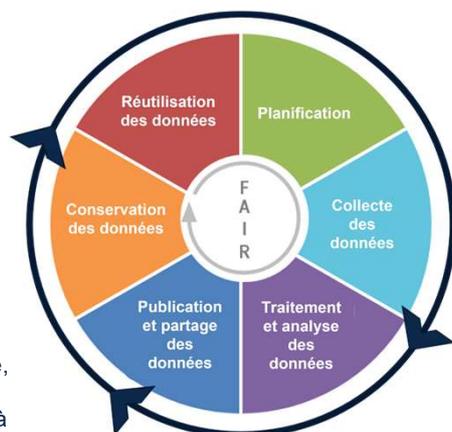
Pour être sûr que les données sont accessibles et réutilisables, veillez à respecter

Les aspects juridiques	Les aspects éthiques
<ul style="list-style-type: none">• Propriété intellectuelle• Accords de consortium• Obligations de diffusion• Communicabilité des données• ...	<ul style="list-style-type: none">• Comité d'éthique• Normes• Chartes

A large, semi-transparent yellow shape that is roughly semi-circular or fan-shaped, positioned on the left side of the slide.

Métadonnées

Place dans le cycle de vie des données



Il est recommandé de renseigner les métadonnées au fur et à mesure de l'avancée du projet

Au moment du partage puis de l'archivage pérenne, des métadonnées spécifiques seront à renseigner

Définition – utilisation des métadonnées

- Les métadonnées permettent de décrire plus précisément les données
- Elles sont à renseigner au fur et à mesure pour chaque jeu de données
- La boîte de conserve = jeu de données / l'étiquette = métadonnées



Sans métadonnées



Avec métadonnées

Renseigner les métadonnées au fur et à mesure pour chaque jeu de données, pour toutes les étapes du cycle de vie des données.

Ressource :

DoRANum. Cours introductif sur les métadonnées. 27 avril 2022.

https://doranum.fr/metadonnees-standards-formats/cours-introductif-sur-les-metadonnees_10_13143_vwce-g965/

Métadonnées embarquées et enrichies

Compléter les métadonnées embarquées par des métadonnées enrichies

Métadonnée embarquée :
produite automatiquement par les appareils (de prise de vue ou de son, de mesure...)



Ex : données GPS, type d'appareil, date, calibrage technique, etc.

Métadonnée enrichie :
ajoutée par l'auteur



Ex : mots-clés, sujet, auteur, laboratoire ou organisme, nom du projet, licence, etc.

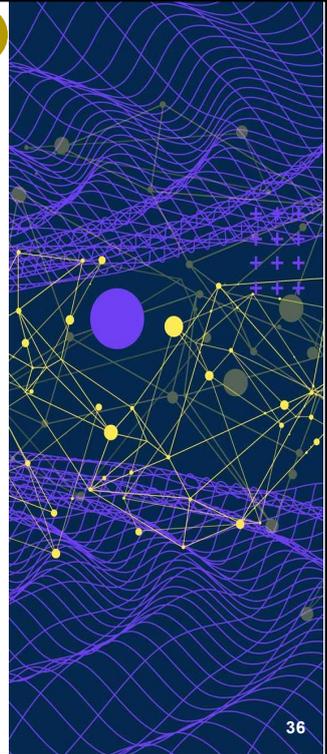


Schéma et standard de métadonnées

Utiliser un standard de votre discipline ou adapté à vos besoins et, s'il n'en existe pas, créer un schéma de métadonnées

Schéma

C'est l'organisation des métadonnées selon un plan pensé et créé spécifiquement pour les besoins d'un projet
 ⇒ Unique et personnalisé

Standard

Un standard est un schéma qui a été adopté comme modèle par un ensemble d'utilisateurs : il est reconnu, normalisé et utilisé à grande échelle
 ⇒ Modèle



Un **schéma de métadonnées** est une **construction organisée d'informations**. C'est une **liste structurée et composée d'éléments descriptifs reliés entre eux**.

L'utilisation d'un standard de métadonnées, notamment disciplinaire, est un élément clé pour atteindre un haut degré de respect des principes FAIR :

- **Facile à trouver** : une donnée n'est souvent trouvable que par les éléments de métadonnées indexés dans le moteur de recherche consulté.
- **Accessible** : notamment grâce aux métadonnées.
- **Interopérable** : grâce au standard commun qui facilite les traitements informatiques.
- **Réutilisable** : grâce à la provenance décrite dans les métadonnées, la licence attribuée et au recours à un standard disciplinaire.

Attention : opter pour un schéma de métadonnées plutôt qu'un standard rend vos données beaucoup moins interopérables et moins FAIR. Seules les métadonnées renseignées selon un standard seront interopérables. Les métadonnées spécifiques le seront moins voire pas du tout.

N'hésitez pas à vous faire aider par un documentaliste ou un informaticien.

Ressource :

- DoRANum. Les schémas de métadonnées. 25 juin 2018. https://doranum.fr/metadonnees-standards-formats/schemas-metadonnees_10_13143_j8bc-0z74/

Exemples de standards de métadonnées

Dublin Core

- Standard interdisciplinaire : description des ressources numériques

DataCite Metadata Schema

- Standard lié à l'attribution d'identifiants pérennes DOI

DDI (Data Documentation Initiative)

- Domaine des sciences sociales, comportementales et économiques

CSMD-CCLRC Core Scientific Metadata Model

- Domaines des sciences structurales (chimie, science des matériaux, sciences de la terre, biochimie)

Exemples de standards de Métadonnées :

- **Dublin Core (interdisciplinaire)** : description des ressources numériques. <http://dublincore.org/>
- **DataCite Metadata Schema** : métadonnées enregistrées dans le DataCite Metadata Store lors de la création d'un DOI pour un jeu de données. Permet l'identification précise et cohérente des données à des fins de citation et de réutilisation. <https://schema.datacite.org/>
- **DDI (Data Documentation Initiative)** : domaine des sciences sociales, comportementales et économiques. <http://www.ddialliance.org/>
- **CSMD-CCLRC Core Scientific Metadata Model** : domaines des sciences structurales (chimie, science des matériaux, sciences de la terre, biochimie) <http://icatproject-contrib.github.io/CSMD/>

Exemples de standards de métadonnées

DwC (Darwin Core)

- Domaine de la biodiversité

EML (Ecological Metadata Language)

- Domaine de l'écologie

MIDAS-Heritage

- Domaine de l'architecture

Exemples de standards de Métadonnées :

- **DwC (Darwin Core)** : domaine de la biodiversité. <http://rs.tdwg.org/dwc/>
- **EML (Ecological Metadata Language)** : très développé en écologie. En grande partie conçu pour décrire des ressources numériques. Il peut également être utilisé pour décrire des ressources non numériques telles que des cartes papier ou d'autres médias. <https://knb.ecoinformatics.org/external//emlparser/docs/index.html>
- **MIDAS-Heritage** : domaine de l'architecture. <https://historicengland.org.uk/images-books/publications/midas-heritage/>

Ressources pour trouver des standards de métadonnées :

- DoRANum. Les standards de métadonnées : pourquoi, lequel ? 13 avril 2021. https://doranum.fr/metadonnees-standards-formats/standard-metadonnees_10_13143_y5py-w521/
- FAIRsharing.org. Répertoire de standards de métadonnées en Sciences de la Vie. <https://fairsharing.org/search?fairsharingRegistry=Standard>
- RDA, Research Data Alliance. Metadata Standards Catalog. <https://rdamsc.bath.ac.uk/>
- DCC, Digital Curation Centre. Disciplinary Metadata. <https://www.dcc.ac.uk/guidance/standards/metadata>

Consulter également les **informations fournies par les entrepôts de données sur les standards de métadonnées.**

Exemple : Le Dublin Core

Élément simple	Définitions	Exemple d'élément spécifique
Title	nom de la ressource	
Subject	thème du contenu de la ressource	
Description	résumé, table des matières...	
Creator	auteur principal de la ressource	
Publisher	entité responsable de la diffusion de la ressource	
Contributor	co-auteurs associés à l'élaboration de la ressource	
Date	date de création ou mise à disposition	<i>Date Created, Date copyrighted, Date Valid, Date Available, Date Modified, Date Accepted, Date Submitted, Date Issued</i>
Type	nature du contenu : image, son texte...	
Format	format ou taille de la ressource	
Identifiant	référence univoque, DOI, URL, ISSN...	
Source	référence à une ressource à partir de laquelle la ressource actuelle a été dérivée ou créée	<i>Has Format (Les relations de transformation de format sont celles où une ressource a été dérivée d'une autre à l'aide d'une technologie de reproduction ou de reformatage qui n'est pas fondamentalement une interprétation mais une représentation.)</i>
Language	langue originale de la ressource	
Relation	référence à une ressource apparentée	<i>isVersionOf (Les relations de version sont celles où une ressource est un état ou une parution historique d'une autre ressource par le même créateur)</i>
Coverage	périmètre spatial et temporel	
Rights	informations sur les droits associés à la ressource	

Le Dublin Core est un **standard international et multidisciplinaire**.

Il **comporte 15 éléments** (= minimum exigé) ayant trait au contenu, à la propriété intellectuelle, à la version.

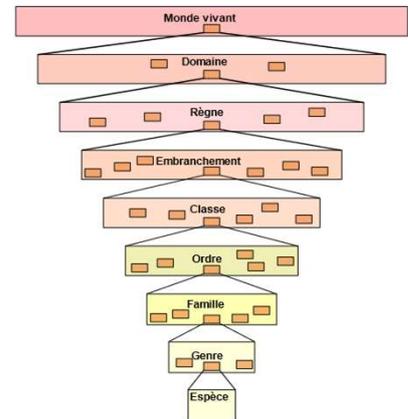
Ces 15 éléments restent cependant insuffisants pour représenter finement certaines données dans un univers dédié pour des utilisateurs et applications spécifiques. C'est pourquoi une quarantaine d'éléments plus spécifiques, qui constituent le Dublin Core étendu, viennent préciser les éléments simples.

Sources :

- DoRANum. *Métadonnées, standards, formats : fiche synthétique*. 27 novembre 2017. https://doranum.fr/metadonnees-standards-formats/metadonnees-standards-formats-fiche-synthetic/10_13143_vbjs-6288/
- Wikipédia. *Dublin Core*. https://fr.wikipedia.org/wiki/Dublin_Core

Enrichissement des métadonnées

- Utiliser des vocabulaires contrôlés disciplinaires utilisés couramment : (lexiques, thesaurus, ontologies...)
- Exemples :
 - Codex de médicaments
 - Classifications taxonomiques
 - Nomenclature internationale des formules chimiques
- Cela augmentera la capacité des données à être combinées avec d'autres données (interopérabilité)



« Le **Thésaurus** est un ensemble organisé de termes contrôlés et normalisés représentant les concepts d'un domaine de connaissance. Les termes sont reliés entre eux par des relations de synonymie (terme équivalent), de hiérarchie (terme générique et terme spécifique) et d'association (terme associé) ; chaque terme appartient à une catégorie ou domaine. » (Wikipedia. *Thésaurus documentaire*.)

https://fr.wikipedia.org/wiki/Th%C3%A9saurus_documentaire

« L'**ontologie** est un ensemble organisé de termes/concepts avec des relations sémantiques variées décrivant un domaine de connaissance. Exemple : FOAF (FriendOfAFriend). (Ecole thématique E-Envir 2021. *Interopérable & Réuse. Introduction aux concepts clefs et immersion. 2 au 5 novembre 2021.* https://e-envir-21.sciencesconf.org/data/pages/6_E_ENVIR21_INTEROP_REUSE.pdf)

« La **taxonomie** ou taxinomie est une branche des sciences naturelles qui a pour objet l'étude de la diversité du monde vivant. Cette activité consiste à décrire et circonscrire en termes d'espèces les organismes vivants et à les organiser en catégories hiérarchisées appelées taxons. (Wikipedia. *Taxonomie*.)

<https://fr.wikipedia.org/wiki/Taxonomie>

« Par extension, [la taxonomie est une] classification, suite d'éléments formant des listes qui concernent un domaine, une science. » (Dictionnaire Larousse.)

<https://www.larousse.fr/dictionnaires/francais/taxinomie/76893>

Source de l'image : Ariel Provost. File:Taxonomic hierarchy.svg (figure renversée : espèce en bas), CC BY-SA 4.0. <https://commons.wikimedia.org/w/index.php?curid=114679881>

Utilité des métadonnées



Les métadonnées sont utiles pour :

Comprendre l'origine des données et leur contexte de création ou de collecte

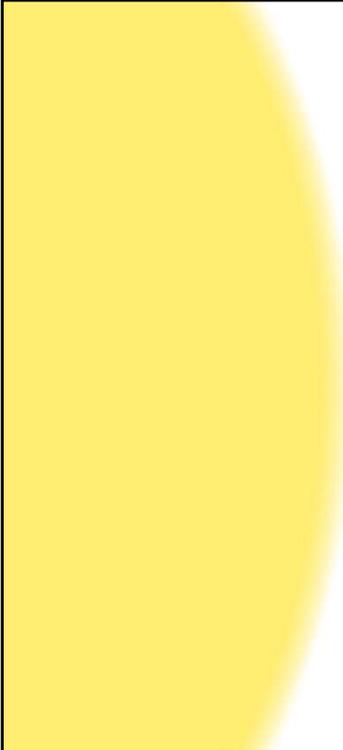
Améliorer le moissonnage par les machines (moteur de recherche)

Garantir l'interopérabilité

Connaitre les conditions de réutilisation et de partage des données

Accéder à des informations très utiles lorsqu'on ne peut pas partager ses données ou lors de la suppression des données obsolètes

Quoiqu'il arrive aux données, les métadonnées resteront toujours accessibles. Il est important de mentionner que les données ont été supprimées et pourquoi.

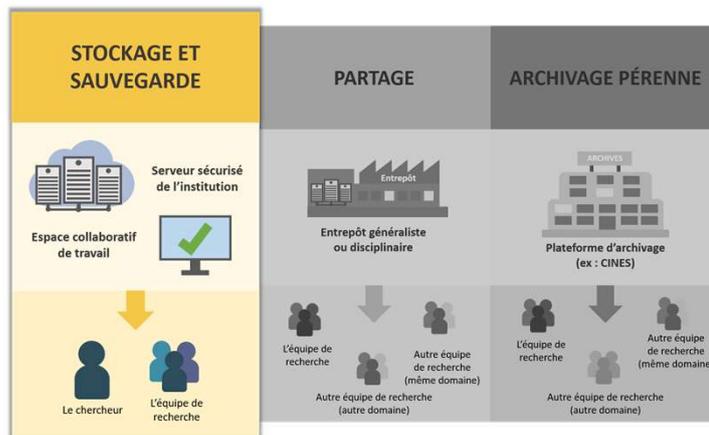


Étape 3 :
Traitement et analyse
Stockage sécurisé durant le projet

Place dans le cycle de vie des données



Stockage et sauvegarde des données



Le **stockage** et la **sauvegarde** sécurisés des données se font durant le projet. L'objectif est de garantir la sécurité des données et d'en faciliter l'accès pour l'ensemble des collaborateurs du projet.

Ressource :

- DoRANum. *Stockage, partage et archivage : quelles différences ? 1^{er} juillet 2021.*
https://doranum.fr/stockage-archivage/stockage-partage-archivage-quelles-differences_10_13143_5dax-qp58/
- DoRANum. *Stocker ses données de façon sécurisée. 26 juillet 2022.*
https://doranum.fr/stockage-archivage/stockage-donnees_10_13143_z0qe-nc29/

Mesures de sauvegarde (stockage)

- Sélectionner les données
- Appliquer la règle du 3-2-1
 - Garder 3 exemplaires des données
 - sur 2 supports ou technologies différents,
 - dont 1 se trouve hors site
- Déterminer l'hébergement
- Estimer la volumétrie des données et la durée de conservation

Organiser et planifier la sauvegarde des données :

- À chaque point d'étape du projet, sélectionner les données à sauvegarder, à supprimer

Dans l'idéal, dupliquer et stocker les données à différents endroits sur différents supports en veillant à bien gérer les versions :

- Les différents états des données sont conservés en corrélation avec les différentes étapes de traitement
- Permet de revenir à une version antérieure si besoin.

Le mieux est de suivre la règle du 3-2-1, c'est-à-dire : garder 3 exemplaires des données, sur 2 supports ou technologies différents dont 1 se trouve hors site.

Déterminer l'hébergement : serveurs locaux (machines virtuelles), cloud institutionnel avec accès sécurisé...

Estimer la volumétrie des données et la durée de conservation

Contactez les services de votre institution pour l'hébergement, la volumétrie, la durée de conservation.

Source :

- *Université de Lille Sciences et Technologies. Sauvegardes ?*
https://indico.math.cnrs.fr/event/2317/contributions/1314/attachments/692/773/problematique_des_sauvegardes-journees_mathrice-20170928.pdf

Exemples d'espaces de stockage



À privilégier

- Plateformes institutionnelles
- Infrastructures disciplinaires



À limiter

- Clé USB
- Disque dur
- Solutions commerciales

Exemples **Plateformes institutionnelles ou infrastructures disciplinaires** :

- Huma-Num Box : <https://documentation.huma-num.fr/humanum-box/>
- Stockage données recherche (PETA) à l'Université de Lorraine : <https://numerique.univ-lorraine.fr/catalogue-des-services/stockage-donnees-recherche-peta>
- Espaces de stockage du CC-IN2P3 : <https://doc.cc.in2p3.fr/fr/Data-storage/storage-areas.html>

Exemples de **solutions commerciales** : Google Drive, Dropbox...

Formats de fichiers

Formats ouverts et non propriétaires

Opter pour des formats de fichiers les plus ouverts possible (non propriétaires), standardisés et pérennes

Exemples :

- Privilégier .csv à .xls
- Privilégier .odt à .doc
- Privilégier .jpg à .tif
- Privilégier .zip à .rar

Choix du format

Le choix d'un format peut être guidé par :

- les recommandations de son institution
- les usages de la communauté scientifique de la discipline
- les logiciels ou équipements utilisés

Ressource :

- DoRANum. Format ouvert ou fermé ? https://doranum.fr/stockage-archivage/quiz-format-ouvert-ou-ferme_10_13143_mcwq-qs64/

Nommage des fichiers

La fiabilité d'accès passe par un nommage unique et précis des fichiers de données :



Bonnes pratiques

- 30 caractères maximum
- Noms de partenaires insérables si leur graphie est harmonisée entre les fichiers
- Numéros de versions le cas échéant
- Dates au format ISO : AAAA-MM-JJ

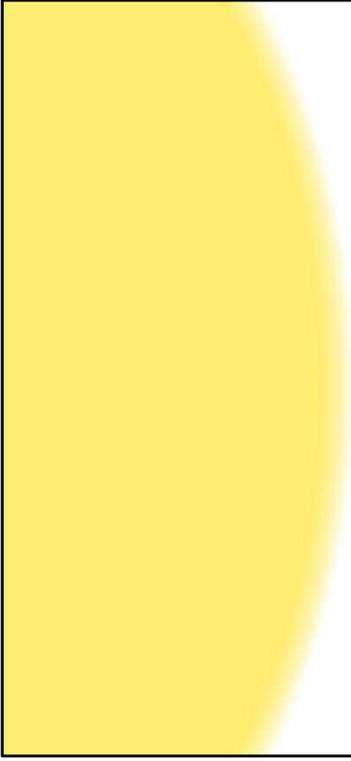


A éviter

- Pas de caractères spéciaux ou accentués du type ùéàç+'@[] :</> */ »& !\$...
- Séparateurs : pas d'espace, pas de mots vides, éventuellement Majuscules ou underscore _
- Pas de dénomination vague : divers, autres, à classer...

Source :

- *Arnould Pierre-Yves, Jacquemot-Perbal Marie-Christine. Guide de bonnes pratiques : Gestion et valorisation des données de recherche. 23 février 2016. <https://hal.science/hal-01275841/>*

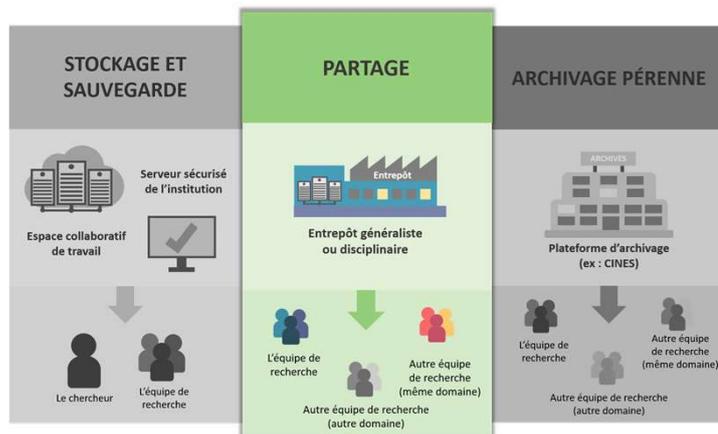


Étape 4 :
Partage, diffusion des données
Dépôt dans un entrepôt

Place dans le cycle de vie des données



Partage, diffusion des données

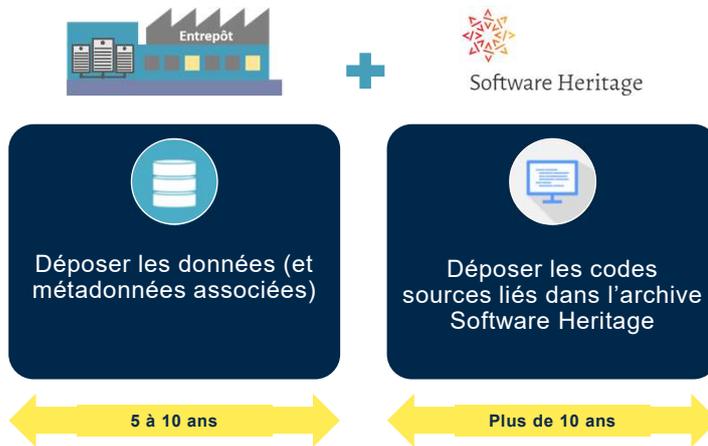


Le **partage** consiste à déposer les données dans un entrepôt de données afin de les rendre accessibles facilement et de permettre leur réutilisation par des chercheurs du même domaine ou d'un autre domaine, selon les principes FAIR et sur le court et le moyen terme. Le partage se fait souvent à l'issue du projet.

Ressource :

- DoRANum. *Stockage, partage et archivage : quelles différences ? 1^{er} juillet 2021.* https://doranum.fr/stockage-archivage/stockage-partage-archivage-queelles-differences_10_13143_5dax-qp58/

Dépôt dans un entrepôt



Pour les logiciels, **déposer les codes sources dans HAL** (<https://hal.archives-ouvertes.fr/>), en lien avec **Software Heritage**, archive de logiciels (<https://www.softwareheritage.org/?lang=fr>).

Préparation des données pour le partage



Check-list

- Sélectionner les données à partager
- Estimer la volumétrie des données
- Vérifier la compatibilité et l'interopérabilité des formats de données
- Migrer si possible vers un format ouvert
- Préparer si nécessaire les codes sources
- Compléter et enrichir les métadonnées
 - Si ce n'est pas déjà fait, choisir un standard de métadonnées
 - S'il n'en existe pas d'adapté, créer un schéma de métadonnées
 - Compléter les champs pour chaque jeu de données

Tout ces éléments contribuent aux principes FAIR :

- Sélectionner les données à partager
 - Estimer la volumétrie des données et éventuellement le budget.
 - Vérifier la compatibilité et l'interopérabilité des formats de données (I/R)
 - Migrer si possible vers un format ouvert (I/R)
 - Préparer si nécessaire les codes sources qui permettront de lire et traiter les données (A/I/R)
 - Compléter et enrichir les métadonnées (en fonction de l'entrepôt choisi) (F/A/I/R)
- :
- Si ce n'est pas déjà fait, choisir un standard de métadonnées
 - S'il n'en existe pas d'adapté, créer un schéma de métadonnées
 - Compléter les champs pour chaque jeu de données, suivant le standard adopté.

Ressource :

DoRANum. Vérifier ses données de recherche. 6 septembre 2022.

https://doranum.fr/depot-entrepots/verifier-donnees-recherche_10_13143_5rs6-4r06/

Choix de l'entrepôt de données

F A I R

1

Privilégier un entrepôt **disciplinaire**

Exemples :

- [PANGAEA](#) (Environnement)
- [GenBank](#) (Génétique)
- [Ortolang](#) (Linguistique)
- ...



2

Si aucun entrepôt disciplinaire n'existe, choisissez un entrepôt **institutionnel**

Exemples :

- [CNRS Research Data](#)
- [Data INRAE](#)
- [DOREL](#)
- ...



3

En alternative, vous pouvez déposer vos données dans l'entrepôt national pluridisciplinaire **Recherche Data Gouv**

Exemples d'entrepôts disciplinaires :

- **PANGAEA** (Environnement) : <https://www.pangaea.de/>
- **GenBank** (Génétique) : <https://www.ncbi.nlm.nih.gov/genbank/>
- **Ortolang** (Linguistique) : <https://www.ortolang.fr/>

Exemples d'entrepôts institutionnels :

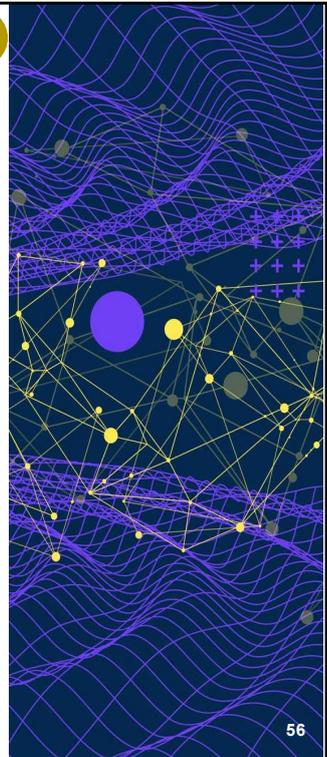
- **CNRS Research Data** : <https://entrepot.recherche.data.gouv.fr/dataverse/cnrs>
- **Data INRAE** : <https://entrepot.recherche.data.gouv.fr/dataverse/inrae>
- **DOREL** : <https://dorel.univ-lorraine.fr/>

Entrepôt pluridisciplinaire :

Recherche Data Gouv : <https://entrepot.recherche.data.gouv.fr/>

Principaux critères de choix d'un entrepôt

- Type de données acceptées
- Qualité des métadonnées
- Entrepôt de confiance - Certification**
- Pérennité des métadonnées et des données
- Génération d'un identifiant unique pérenne
- Gestion des versions
- Gestion des licences



Un entrepôt digne de confiance devrait permettre et assurer :

- le repérage et l'identification des données
- la recherche, la citation et le téléchargement des données
- la gestion des versions des jeux de données
- le référencement d'informations pertinentes complémentaires, telles que d'autres jeux de données et publications
- l'accès ouvert à des informations mises à jour, y compris sur des données non publiées, protégées, rétractées ou supprimées : métadonnées archivées sur le long terme, même si les données correspondantes ne sont plus disponibles
- la récupération des métadonnées par les machines
- l'accès aux données dans des conditions bien définies (licences)
- l'authenticité et l'intégrité des données
- la confidentialité et le respect des droits des personnes et créateurs de données
- la pérennité des données et métadonnées

Sources :

RDA. Entrepôts de données de confiance : critères de conformité. <https://www.rda-alliance.org/system/files/documents/CoretrustsealFR.pdf>

Science Europe. Practical guide to the international alignment of research data management – Extended Edition. 27 janvier 2021. <https://doi.org/10.5281/zenodo.4915862>

Ressource :

DoRANum. Les critères pour choisir un entrepôt de données. 30 août 2023. <https://doranum.fr/depot-entrepots/criteres-pour-choisir-entrepot-de-donnees-10-13143-zqpb-9449/>

Annuaire et catalogues d'entrepôts

Les annuaires et catalogues permettent d'identifier un entrepôt disciplinaire qui conviendrait à vos besoins :



- **CatOPIDoR** : Catalogue pour une Optimisation du Partage et de l'Interopérabilité des Données de la Recherche. Wiki des services dédiés aux données de recherche. Il recense des entrepôts français.
[https://cat.opidor.fr/index.php/Entrep%C3%B4t de donn%C3%A9es](https://cat.opidor.fr/index.php/Entrep%C3%B4t_de_donn%C3%A9es)
- **Re3data** : annuaire recensant des entrepôts au niveau international et permettant de filtrer selon plusieurs critères (discipline, attribution d'un identifiant pérenne, types de données ou formats acceptés, préservation à long terme des données, choix de la licence, certification, etc.). <https://www.re3data.org/>
- **OAD** : Open Access Directory – Data repositories : liste d'entrepôts et de bases de données pour les données ouvertes.
https://oad.simmons.edu/oadwiki/Data_repositories
- **OpenDOAR** : annuaire mondial d'entrepôts ou d'archives ouvertes en libre accès. La recherche et la navigation peuvent s'effectuer sur différents critères tels que le nom, la région du monde ou le pays et le logiciel. <https://v2.sherpa.ac.uk/opensoar/>
- **FAIRsharing** : ressource permettant d'identifier et citer des standards, des bases de données ou des entrepôts existant pour leurs données et leur discipline.
<https://fairsharing.org/>

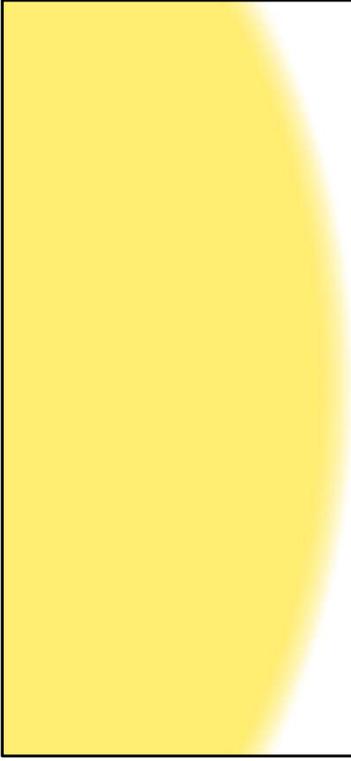
Inist-CNRS. Cat OPIDoR. <https://cat.opidor.fr/>

Recense et décrit les services français dédiés aux données scientifiques. Proposé sous forme d'un wiki, cet outil collaboratif ouvert à tous permet de repérer et ajouter des services utiles dans le cadre d'un projet de recherche

Cat OPIDoR présente par domaine scientifique :

- des sites d'information,
- de formation,
- des outils de gestion,
- des plateformes,
- Des entrepôts de données

pour accompagner les chercheurs sur l'ensemble des étapes clés de la gestion, collecte, stockage, conservation et ouverture des données.



Identifiants pérennes

Identifiants pérennes pour les données

- Attribuer un identifiant pérenne à chacun des jeux de données
- Le plus utilisé est le DOI
- Un identifiant pérenne facilite le suivi, la localisation, l'accès et la citation des données lors de leur publication ou à des fins de réutilisation
- Le plus souvent, un identifiant pérenne est attribué aux données de manière automatique lors du dépôt dans un entrepôt !



Une équipe dédiée de l'Inist-CNRS se tient à disposition pour conseiller dans l'attribution de DOI aux données de recherche et fournit :

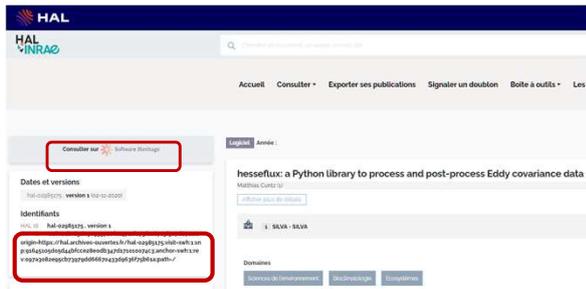
- un accès à un espace test pour enregistrer temporairement des DOI et vérifier la compatibilité de ce service avec ses propres workflows
- un préfixe unique de DOI
- un accès à la plateforme Metadata Store de DataCite (MDS) pour commencer à créer les DOI
- une assistance à la création et à la conversion de métadonnées...
- un accompagnement dans l'utilisation des différents services proposés par DataCite.

Ressources :

- DoRANum. Zoom sur le DOI. 21 mai 2021. https://doranum.fr/identifiants-perennes-pid/zoom-doi_10_13143_j5xt-6j41/
- Inist-CNRS. PID OPIDoR. <https://opidor.fr/identifier/>

Identifiants pérennes pour les logiciels

SWHID est l'identifiant pérenne dédié aux logiciels, créé par Software Heritage



Il est possible de déposer les codes sources dans HAL, avec transfert vers Software Heritage

Créé par Software Heritage, SWHID (Software Hash Identifier, anciennement Software Heritage Identifier) est un système d'identifiant pérenne sous forme de code alphanumérique qui identifie de manière unique les codes sources des logiciels.

Ressources :

- DoRANum. Zoom sur SWHID. 2 juillet 2021. https://doranum.fr/identifiants-perennes-pid/zoom-swhid_10_13143_3qqg-yx41/
- Software Heritage. <https://www.softwareheritage.org/?lang=fr>
- HAL Documentation. Déposer le code source d'un logiciel. <https://doc.archives-ouvertes.fr/deposer/deposer-le-code-source/>

Identifiants pérennes pour les auteurs et les institutions

F A I R

Attribuer un **identifiant auteur** (ORCID)

- Fait le lien avec ses productions scientifiques
- Permet d'être bien identifié et cité
- Augmente la visibilité internationale

Pensez à vérifier si votre institution bénéficie d'un **identifiant pérenne ROR**

ORCID

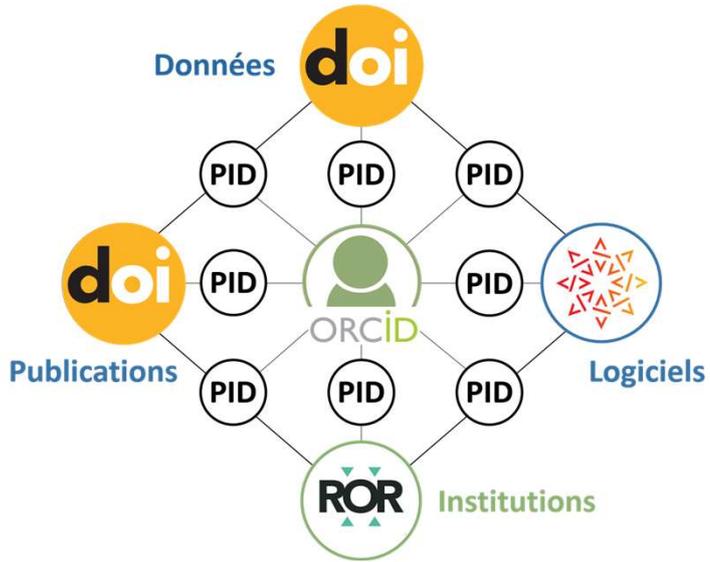
ROR

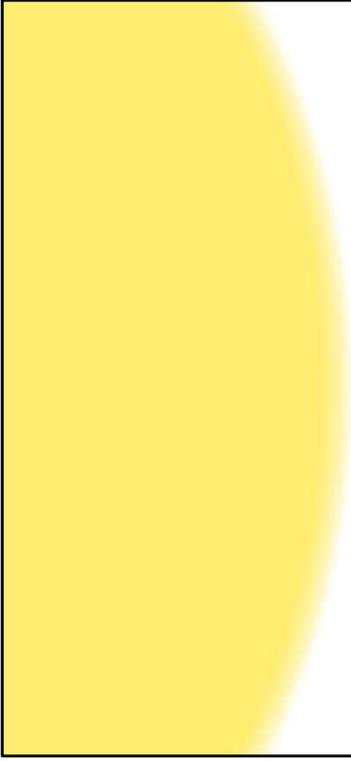
Ressources :

DoRANum. Zoom sur ORCID. 5 décembre 2022. https://doranum.fr/identifiants-perennes-pid/zoom-orcid_10_13143_c6rx-9w77/

DoRANum. Zoom sur ROR. 8 juin 2022. https://doranum.fr/identifiants-perennes-pid/zoom-sur-ror_10_13143_h5zy-bn73/

Identifiants pérennes – Pour résumer





Modalités de protection

Accès et licences



- Accès
 - Dispositif d'accès contrôlé : mot de passe...
 - Accès limité dans le temps par un embargo
 - Accès limité à certaines personnes



- Chiffrement des données sensibles pour éviter des intrusions malveillantes



- Licences pour éviter une mauvaise (ré)utilisation des données par autrui

Accès :

- Accès limité dans le temps par un embargo
 - Pour disposer du temps nécessaire au dépôt de brevets...
 - En fonction de la discipline
 - Peut être déterminé par les éditeurs
- Accès limité à certaines personnes
 - Ex : limitation aux membres du consortium ou à une communauté scientifique

Licences de diffusion

Attribuer une licence de diffusion à chaque jeu de données permet d'afficher clairement les modalités de réutilisation

- En France → **Licence ouverte** Etalab (équivalente de la CC-BY) : l'auteur doit obligatoirement être cité
- Licences Creative Commons (CC-BY)
 - avec 3 possibilités de combinaisons :
 - NC : pas d'utilisation commerciale
 - SA : partage dans les mêmes conditions
 - ND : pas de modifications
- Licences spécifiques pour les logiciels et les bases de données



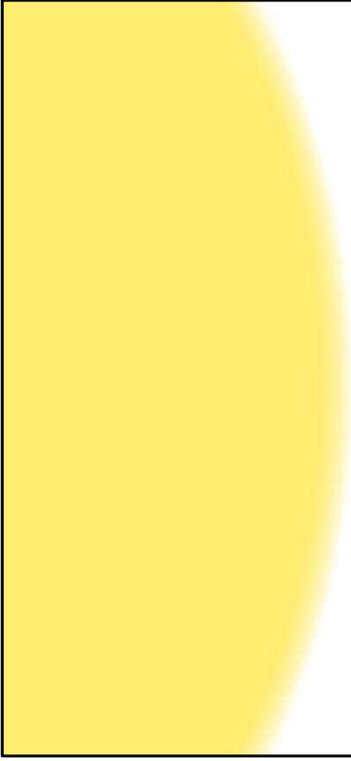
En France, la licence la plus adaptée pour les données publiques est la **Licence ouverte** Etalab (équivalente de la CC-BY) : l'auteur doit obligatoirement être cité. Sinon, il est possible d'utiliser les licences Creative Commons (CC-BY), avec 3 possibilités de combinaisons :

- NC : pas d'utilisation commerciale
- SA : partage dans les mêmes conditions
- ND : pas de modifications

Il existe des licences spécifiques pour les logiciels et les bases de données

Ressources :

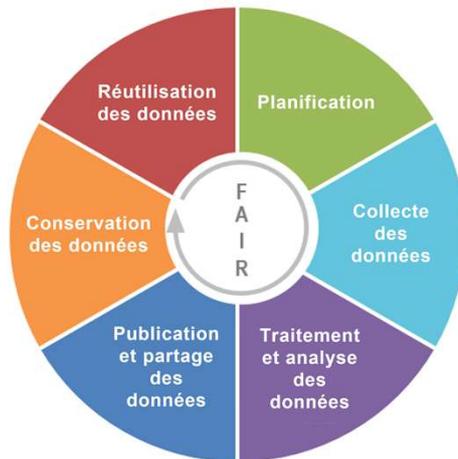
- DoRANum. Guide des licences ouvertes. https://dorum.fr/aspects-juridiques-ethiques/guide-des-licences-ouvertes_10_13143_tv6f-sv31/
- DoRANum. Les licences de réutilisation dans le cadre de l'open data et de la loi pour une République numérique. https://dorum.fr/aspects-juridiques-ethiques/les-licences-de-reutilisation-dans-le-cadre-de-lopen-data-2_10_13143_ssh2-zd93/
- Data.gouv.fr. Licences de réutilisation. <https://www.data.gouv.fr/fr/licences>
- Licentia by Inria. License your Data. <http://licentia.inria.fr/>
- License Selector. <http://ufal.github.io/public-license-selector/>
- Choose an open source license : <https://choosealicense.com/>

A large, semi-transparent yellow shape on the left side of the slide, resembling a stylized 'C' or a partial circle.

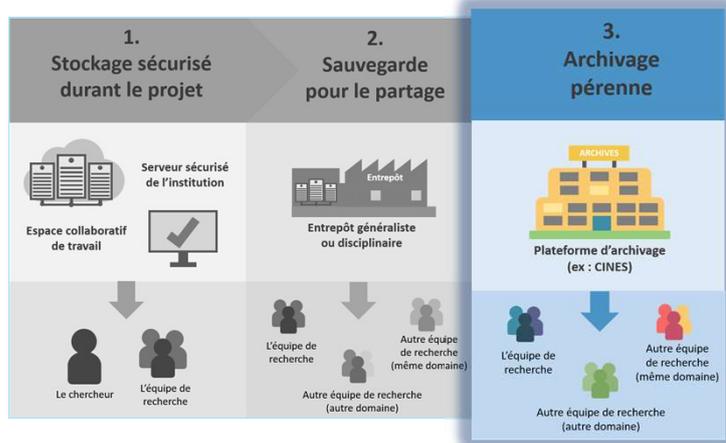
Étape 5 : Archivage pérenne

Place dans le cycle de vie des données

Étape 3
Archivage
pérenne



Archivage pérenne



L'**archivage pérenne** a pour objectif de conserver les données, d'en garantir l'accès et d'en préserver l'intelligibilité sur le long terme, c'est-à-dire plus de 30 ans. En réalité, l'archivage pérenne concerne peu de données. Seulement celles qui présentent une grande valeur scientifique reconnue par la communauté dont elles proviennent. Soit parce qu'elles sont très coûteuses, soit parce qu'elles sont uniques, non reproductibles.

Ressource :

- *DoRANum. Stockage, partage et archivage : quelles différences ? 1^{er} juillet 2021.*
<https://doranum.fr/stockage-archivage/stockage-partage-archivage-quelles-differences-10-13143-5dax-qp58/>

Définition et objectifs

L'archivage pérenne a pour objectif de :

- Conserver uniquement les **données à forte valeur scientifique, impossibles ou coûteuses** à (re)produire
- Conserver les données dans leur **aspect physique** comme dans leur **aspect intellectuel**
- En préserver l'**intelligibilité** sur le très long terme
- De manière à ce qu'elles soient en permanence **accessibles et compréhensibles**



Archivage pérenne coûteux et soumis à la décision de l'institution

Valeur scientifique des données :

Sont-elles uniques, non reproductibles (ou à des coûts trop élevés) ?

Ont-elles une valeur historique (représentent-elles un point de repère dans les découvertes scientifiques ?)

...

Attention, l'archivage pérenne est une opération très coûteuse. Elle est soumise à une décision de votre institution/laboratoire. Dans le cadre d'un projet financé, ce coût peut être éligible.

Source :

NERC Data Value Checklist. <https://www.ukri.org/publications/nerc-data-value-checklist/>

Plateforme d'archivage du CINES

- Le **CINES** est l'opérateur mandaté par le Ministère pour opérer la mission d'archivage pérenne pour l'Enseignement Supérieur et la Recherche.
- Propose l'outil FACILE
- Selon son institution, sa discipline ou l'entrepôt choisi, il existe déjà des partenariats avec le CINES, proposant un accompagnement pour l'archivage.

Ex : Huma-Num en SHS

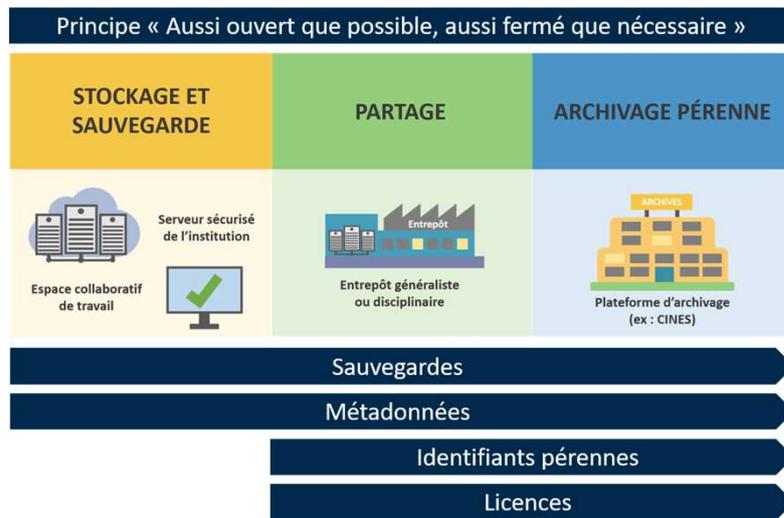
Outil FACILE : pour tester la validité et/ou l'éligibilité de ses formats de fichiers selon les règles du CINES.

Ressources :

CINES. Nos solutions d'archivage. <https://www.cines.fr/archivage/nos-solutions-darchivage/>

Programme VITAM. <https://www.programmevitam.fr/>

En résumé



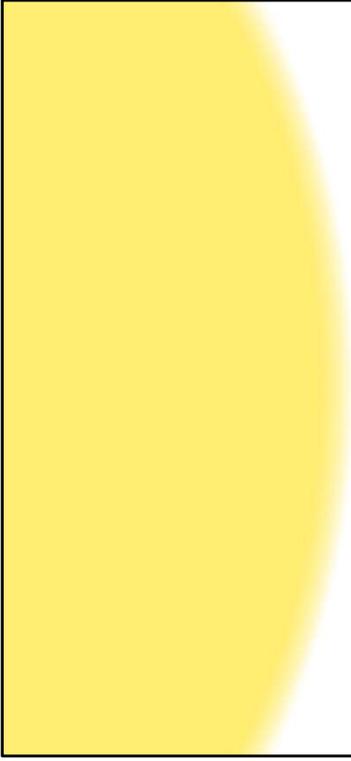
- Le **stockage** et la **sauvegarde** sécurisés des données se font durant le projet. L'objectif est de garantir la sécurité des données et d'en **faciliter l'accès** pour l'ensemble des collaborateurs du projet.
- Le **partage** consiste à déposer les données dans un entrepôt de données afin de les rendre **accessibles facilement** et de permettre leur **réutilisation** par des chercheurs du même domaine ou d'un autre domaine, selon les principes FAIR et sur le court et le moyen terme. Le partage se fait souvent à l'issue du projet.
- L'**archivage pérenne** a pour objectif de conserver les données, d'en garantir l'**accès** et d'en préserver l'intelligibilité sur le long terme.

Les **bonnes pratiques** de **sauvegarde**, d'ajout de **métadonnées**, attribution d'**identifiants pérennes** et de **licences** participent tous au respect des **principes FAIR**.

Ressource :

DoRANum. *Stockage, partage et archivage : quelles différences ? 1^{er} juillet 2021.*

<https://doranum.fr/stockage-archivage/stockage-partage-archivage-quelles-differences-10-13143-5dax-qp58/>

A large yellow shape on the left side of the slide, resembling a quarter-circle or a soft-edged rectangle, with a gradient effect.

Étape 6 : Réutilisation et valorisation des données

Place dans le cycle de vie des données



Du côté du chercheur

- Plus de visibilité et meilleure valorisation de ses travaux
- Mieux retrouver ses propres données
- Conformité aux demandes des financeurs
- Meilleure accessibilité de la science aux citoyens



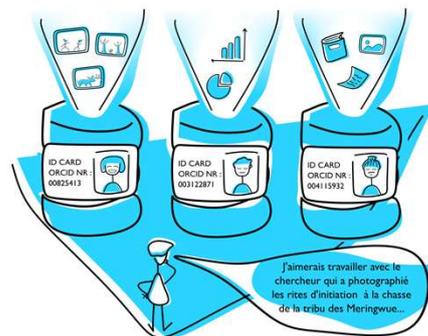
- Plus de visibilité et meilleure valorisation de ses travaux : si vos publications sont reliées à vos données et que vous donnez accès à celles-ci, cela renforce la confiance en vos publications et leur apporte davantage de crédit. Cela permet aussi d'éviter de dupliquer des données ou des résultats déjà publiés.
- Des données mal documentées ou mal conservées sont des données perdues ! Les principes FAIR facilitent la bonne gestion des données, pendant et après le travail de recherche. L'objectif d'un travail de recherche est d'être lu, vérifié, cité, réutilisé. Ceci est impossible si personne n'a accès aux informations, si le format choisi pour le stockage des données n'est plus lisible, si l'information est perdue, ou si l'adresse url d'accès n'est plus disponible !
- Science citoyenne : vos publications, vos données gagnent à être mise à la disposition de tous, mais pour ce faire, il faut qu'elles soient compréhensibles et manipulables, et accompagnées de métadonnées.

Source :

Bibliothèque de l'Observatoire de Paris. GT Science ouverte. Des principes ? Pour quoi FAIR ?
13 avril 2021. <https://www.observatoiredeparis.psl.eu/IMG/pdf/pour-quoi-fair-3.pdf>

Du côté du réutilisateur

- Trouver des jeux de données déjà existants (dans les entrepôts et data papers)
- Découverte de potentiels partenaires
- Reproductibilité et répliquabilité des résultats



La **reproductibilité** est la capacité, par une équipe différente, de reproduire une expérience, sans se fier au dispositif expérimental et aux codes logiciels développés par l'équipe d'origine.

La **répliquabilité** est la capacité, par une équipe différente, de reproduire une expérience en ré-utilisant le même dispositif expérimental décrit (y compris les codes logiciels).

Ressources pour rechercher des jeux de données :

- *DataCite Commons.* <https://commons.datacite.org/>
- *OpenAIRE EXPLORE.* <https://explore.openaire.eu/search/find/datasets>
- *Google Dataset Search.* <https://datasetsearch.research.google.com/>

Data journal – data paper

- Data paper = publication qui décrit des jeux de données de recherche et les métadonnées associées
- Article à part entière, suivant le même processus éditorial que les articles scientifiques classiques
- Deux possibilités de publication d'un data paper :
 - dans un data journal (revue dédiée à ce type de publication)
 - dans une revue classique



Voici un exemple de valorisation des données : le data paper.

Structure d'un data paper :

Partie descriptive :

- Éléments communs aux articles classiques : titre, résumé, mots-clés...
- Éléments spécifiques aux données : types de données, formats, processus et méthodes de production, métadonnées, réutilisation...

Accès aux données : déposées dans un entrepôt. L'identifiant des données (exemple DOI) permet d'établir le lien du data paper vers les données.

Ressources :

- Exemples de data paper :

Richardson Andrew D., Hufkens Koen, Milliman Tom et al. Tracking vegetation phenology across diverse North American biomes using PhenoCam imagery.

<https://www.nature.com/articles/sdata201828>

Richardson Andrew D., Hufkens Koen, Milliman Tom et al. Tracking vegetation phenology across diverse biomes using Version 2.0 of the PhenoCam Dataset.

<https://www.nature.com/articles/s41597-019-0229-9>

- Exemple de data journal : Journal of Physical and Chemical Reference Data.

<https://pubs.aip.org/aip/jpr>

- Exemples de revues publiant des data papers : Cirad. CoopIST. Publier un Data paper.

<https://coop-ist.cirad.fr/gerer-des-donnees/publier-un-data-paper/4-choisir-la-revue>

- DoRANum. Data papers et data journals. <https://doranum.fr/data-paper-data-journal/>

Exposition et visualisation des données

L'outil Omeka est un logiciel libre couramment utilisé pour l'exposition et la visualisation des données

Exemple 2 : Bibliothèque numérique réalisée avec Omeka

Exemple : Bibliothèque numérique CoReA (Corpus et Ressources Archéologiques) pour la documentation archéologique du Centre Camille Julian, réalisée avec Omeka. Permet de naviguer dans des corpus et ressources en archéologie. <http://ccj-corea.cnrs.fr/>

Omeka est un outil qui s'inscrit pleinement dans la science ouverte et qui permet la création de bases de données au plus près des principes FAIR.

À partir de données de recherche brutes, l'outil permet de créer des bases de données éditorialisées, autrement dit structurées, accessibles, et visibles sur le web. L'outil offre une grande modularité des fonctionnalités grâce à de nombreux plugins, et traite les divers objets multimédia (textes, images, sons, vidéos).

l'outil offre plusieurs avantages techniques :

- l'interface est simple et intuitive ;
- les métadonnées sont moissonnables, permettant notamment le référencement dans d'autres bases ;
- une base Omeka peut être connectée à d'autres services grâce à une API REST.

Ressources :

DoRANum. Accès et visualisation. <https://doranum.fr/acces-visualisation/>

04

À retenir



En résumé,

- Vous rendrez vos données **Faciles à trouver** en les déposant dans un entrepôt de confiance, vérifiant qu'un identifiant pérenne leur est attribué et en les décrivant précisément avec des métadonnées.
- Vous rendrez vos données **Accessibles** en déterminant qui peut y accéder et comment. Ne pas oublier que si vos données ne sont pas ouvertes, les métadonnées doivent l'être.
- Vous rendrez vos données **Interopérables** en utilisant des standards de votre discipline et utilisant des formats ouverts.
- Vous rendrez vos données **Réutilisables** en documentant vos données de manière à ce que les autres utilisateurs puissent les comprendre et en leur attribuant une licence.

Si vous rédigez votre **Plan de Gestion de Données** en pensant aux bonnes pratiques à mettre en œuvre en suivant les principes FAIR, vous atteindrez une gestion optimale de vos données.

Merci de votre attention

contact@dorandum.fr

**Pour en savoir plus, rendez-vous sur
<https://dorandum.fr>**

