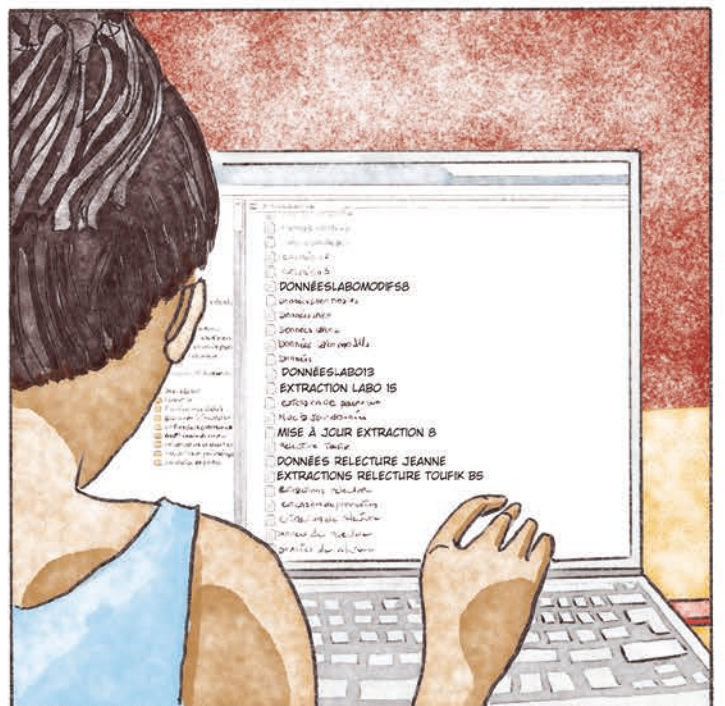
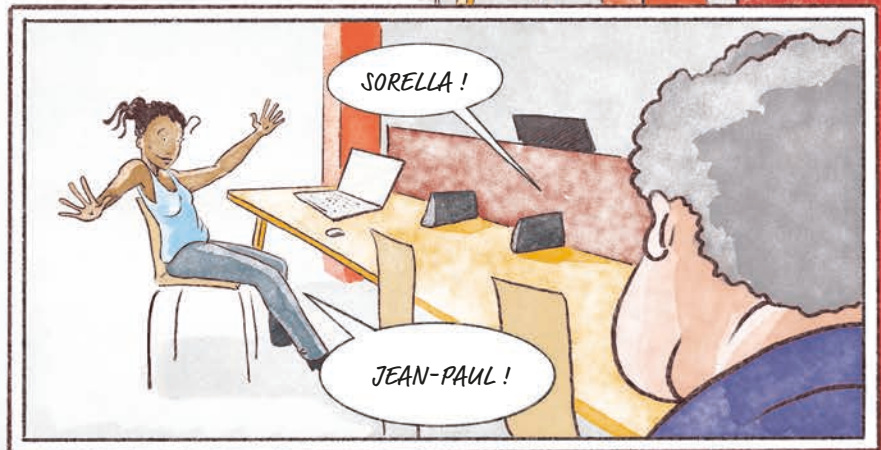


ON FAIT LE POINT
sur les **données**
de la recherche
avec **Sorella !**







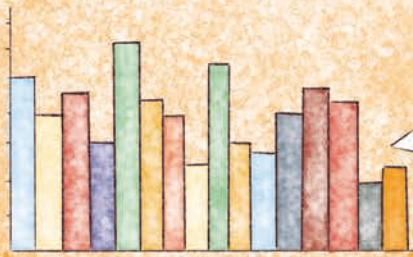
1. Qu'est-ce que des données de la recherche ?

EN FAIT, IL EXISTE UN GRAND NOMBRE DE TYPES DE DONNÉES, QUI NE SONT D'AILLEURS PAS LES MÊMES SELON LES THÈMES TRAITÉS ET LES DISCIPLINES COUVERTES. SELON LE CONTEXTE, ELLES PEUVENT ÊTRE :

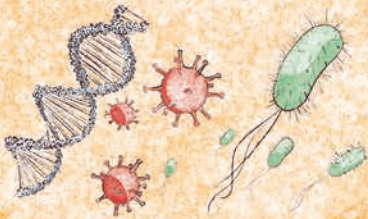
Des données de simulation numérique (modèles climatiques, modèles économiques, etc.)



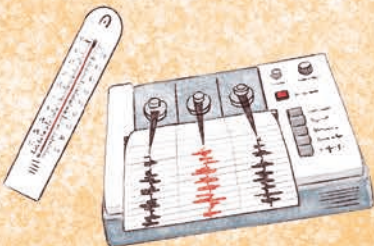
Des données dérivées ou compilées (bases de données compilées, feuilles de texte, statistiques de population, etc.)



Des données de référence (corpus textuels, séquences gènes TP53, structures chimiques, etc.)



Des données d'observation (relevés météo, mesures sismiques, images, enquêtes sociales, fouilles archéologiques, etc.)



Des données d'archives (plaques de verre, fonds de photographies, textes de lois, graines de plantes, etc.)

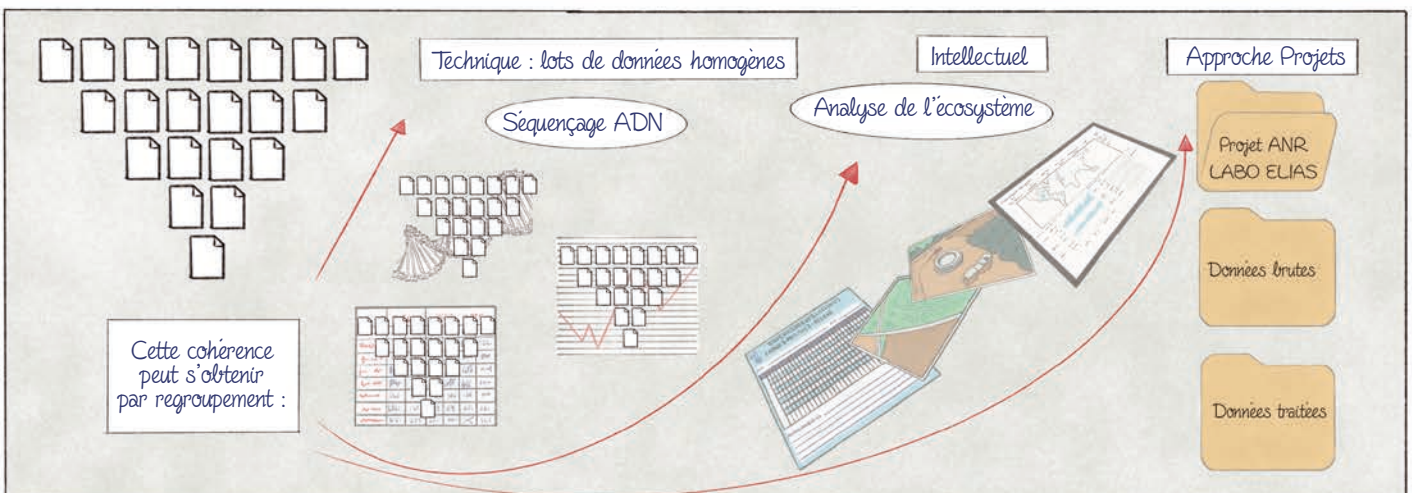
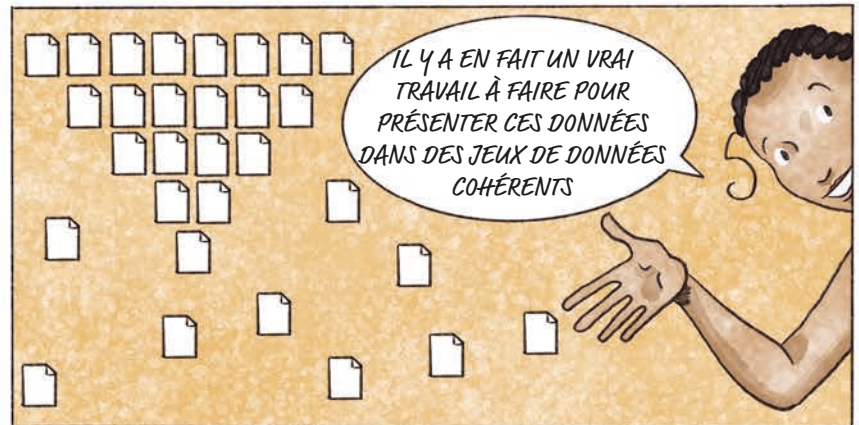
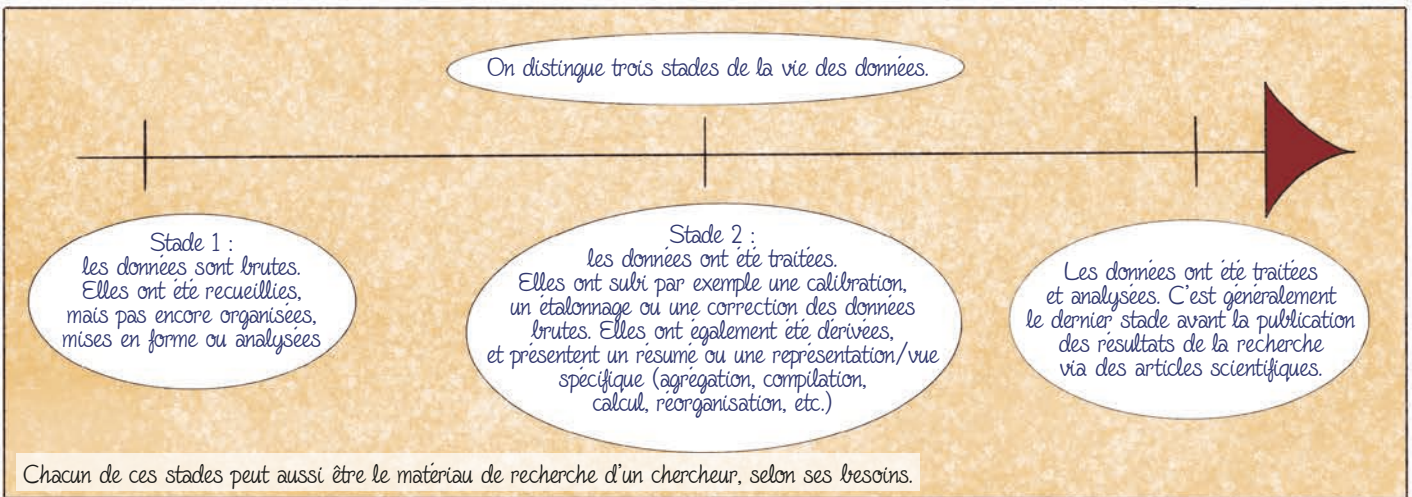


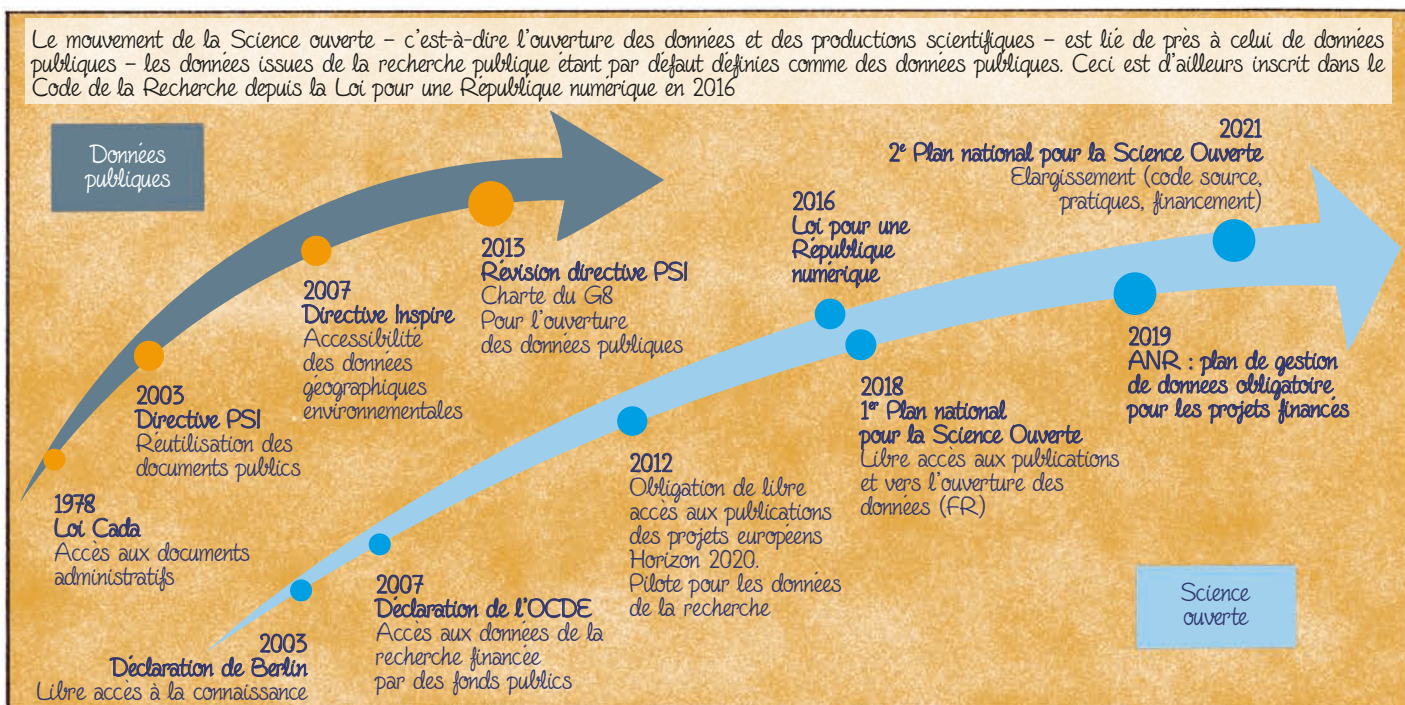
Des données expérimentales (poids biomasse, séquences peptides, tests en psychologie, etc.)



AH ! TU VOIS !
C'EST EXACTEMENT CE QUE
J'AI FAIT AU FINAL :
DU TABLEUR EXCEL !







PARCE QUE LA SCIENCE OUVERTE EST AU CŒUR DES STRATÉGIES DE CINQ ACTEURS CLÉS QUI Y TROUVENT LEURS INTÉRÊTS.

- Les financeurs et États
- Les éditeurs scientifiques
- Les organismes de recherche et les universités
- Les chercheurs
- La société civile

LES FINANCEURS ET ÉTATS
Il s'agit pour eux de favoriser la réutilisation des données afin de permettre un retour sur investissement et de faciliter l'innovation scientifique et technologique.

TIENS, LÉOPOLD DIAGNE A RÉALISÉ QUASIMENT LA MÊME ÉTUDE QUE MOI. ET ELOÏSE SAMPRO AUSSI, MAIS SUR LES ANNÉES 2005-2015. JE N'AI PAS ACCÈS À LEURS DONNÉES OU CELLES-CI SONT INUTILISABLES. TANT PIS, JE REDEMANDE UN FINANCEMENT POUR REFAIRE L'ÉTUDE !

C'EST PAS UN PEU BÊTE, ÇA ?

LES ÉDITEURS SCIENTIFIQUES
Outre le devoir de répondre à la demande des financeurs, les éditeurs scientifiques voient dans le partage des données un moyen scientifique d'aboutir à une meilleure validation des travaux en s'assurant notamment de leur reproductibilité et de renforcer la confiance dans les résultats présentés.

C'EST BIZARRE. RAOUL NEWTON A PUBLIÉ UNE ÉTUDE CERTIFIANT QUE LE PLOMB POUVAIT VOLER SOUS CERTAINES CONDITIONS ATMOSPHÉRIQUES TERRESTRES. MAIS PERSONNE N'A ENCORE RÉUSSI À REPRODUIRE L'EXPÉRIENCE AVEC LES DONNÉES QU'IL A PARTAGÉES...

LES ORGANISMES DE RECHERCHE ET LES UNIVERSITÉS
C'est pour eux l'occasion de promouvoir une bonne éthique et une reproductibilité des recherches au sein des laboratoires, en améliorant la gestion des ressources et en suscitant des collaborations entre institutions.

VOICI JEANNINE DUGLAS, SPÉCIALISTE DE LA REPRODUCTION DES DROSOPHILES EN NOUVELLE-GUINÉE, QUI NOUS PRÉSENTE SON NOUVEAU PROJET

AH ! NOUS, ON TRAVAILLE SUR L'ADN DES DROSOPHILES : TU CROIS QU'ON POURRAIT LUI DEMANDER LES DONNÉES ?

LES CHERCHEURS
Quant aux chercheurs, ils apprennent à mieux gérer les données, à les sécuriser, à les préserver, et produisent des travaux académiques complémentaires qui peuvent leur apporter une reconnaissance et du crédit scientifiques.

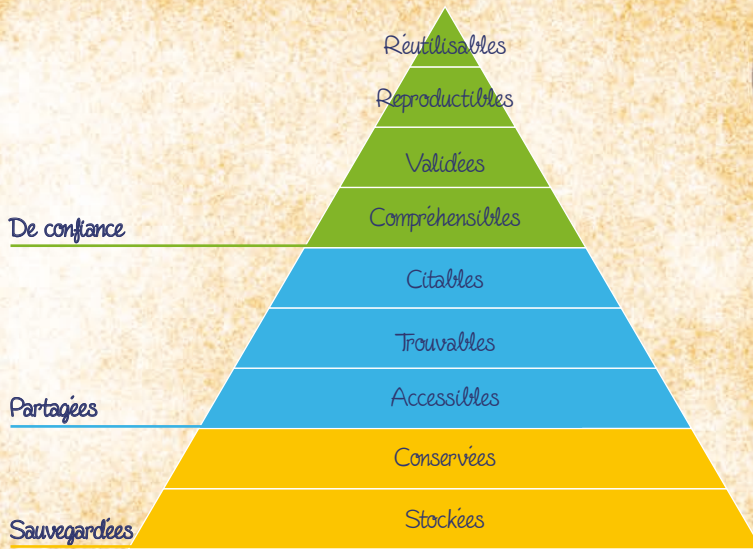
30 ANS DE RECHERCHE ! ET VA RETROUVER QUELQUE CHOSE MAINTENANT DANS CE BAZAR ! AH ! SI SEULEMENT J'AVAIS TRAITÉ ÇA AU FUR ET À MESURE...

LA SOCIÉTÉ CIVILE
Et enfin, la société civile peut ainsi profiter d'une source d'information fiable, favoriser l'innovation et la participation des citoyens à la recherche.

SUPER ! J'AI TOUTES LES DONNÉES DE GÉOLOCALISATION DES TOILETTES PUBLIQUES DE PARIS : GO POUR UNE APPLICATION DE RÉSERVATION DE PLACES !



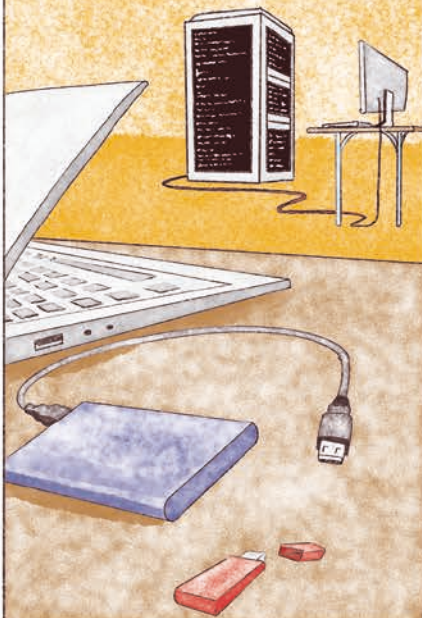
AVANT DE POUVOIR PARTAGER ET DIFFUSER DES DONNÉES - ET AFIN DE LE FAIRE DANS LES MEILLEURES CONDITIONS POSSIBLES - IL Y A DES ÉTAPES DE BONNE GESTION À RESPECTER, RENDANT AU FUR ET À MESURE LES DONNÉES PLUS FIABLES ET DE MEILLEURE QUALITÉ.



D'après Anita de Waard, Helena Cousijn et IJsbrand Jan Aalbersberg, CC BY NC 10 aspects of highly effective research data : www.elsevier.com/connect/10-aspects-of-highly-effective-research-data

Etape 1 : Données sauvegardées

Les données doivent d'abord être **stockées**, idéalement selon la règle du 3-2-1 : 3 copies identiques, stockées sur 2 supports différents + 1 copie hors site



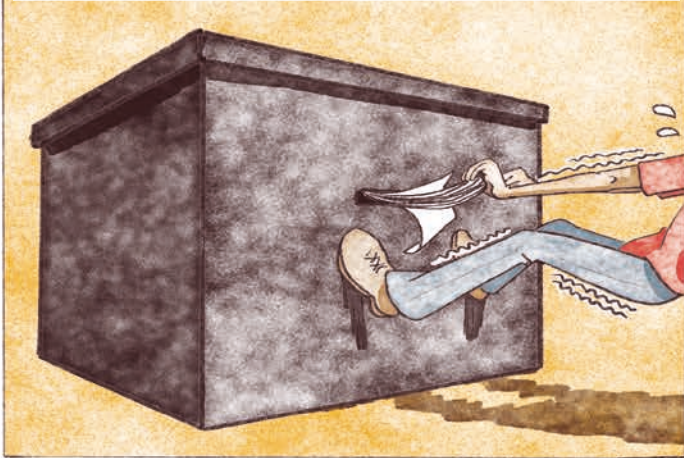
Mais cela ne suffit pas. Afin de prévenir l'obsolescence des technologies informatiques, les données doivent être **conservées**, c'est-à-dire sélectionnées sur différents critères et selon différentes durées de vie (conservation durant 6 mois, 3 ans, 10 ans, etc.)



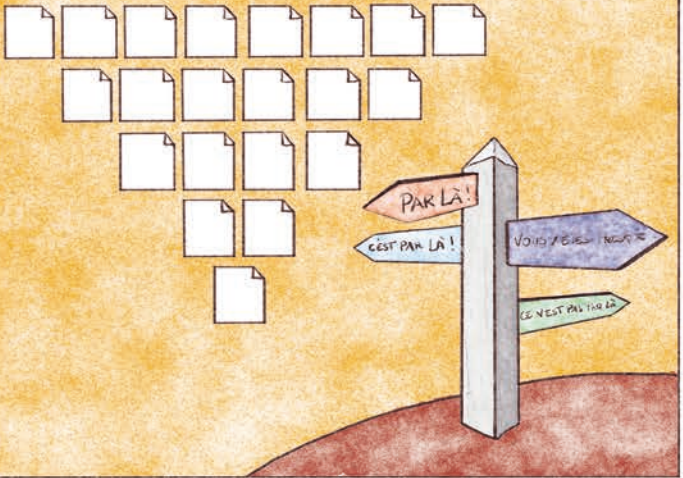
AH ! ÇA ME RAPPELLE ! LORSQUE J'ÉTAIS ENCORE JEUNE ET BEAU, J'AVAIS ENREGISTRÉ MA PREMIÈRE ÉTUDE SUR 52 DISQUETTES INFORMATIQUES. TU CROIS QUE JE POURRAIS DEMANDER À UN STAGIAIRE DE ME RESTAURER CES DONNÉES ?



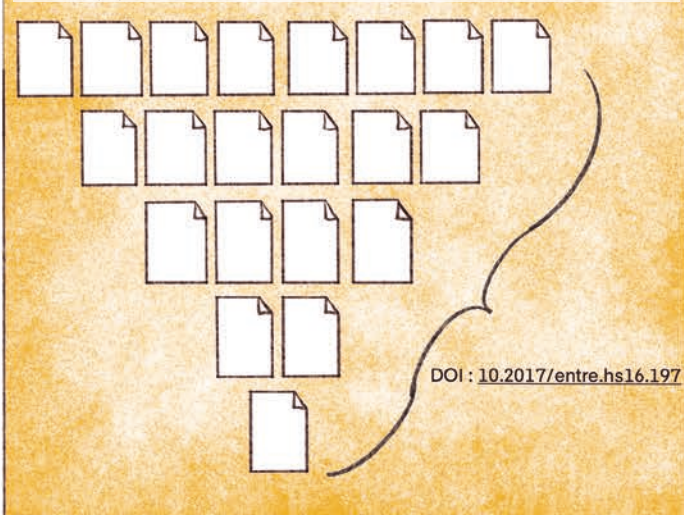
Même lorsque les données sont stockées et préservées, elles ne sont pour autant pas toujours disponibles pour les chercheurs et les machines qui voudraient les interroger. Il faut donc a minima les rendre **accessibles** en ligne.



Mais même en ligne, les données de la recherche ne sont pas toujours aisément **trouvables** par d'autres chercheurs. Il convient donc de les signaler - par exemple en améliorant la qualité de leur description.



Et afin de suivre la réutilisation de ces données et de s'assurer que les chercheurs qui les ont produites bénéficient du crédit scientifique qu'ils méritent, il est conseillé de les rendre **citables**, par exemple en leur attribuant un DOI (Digital Object Identifier) ou en les reliant à un *data-paper*.



Des données qui ont été collectées pour un usage interne ne sont pas forcément **compréhensibles** par un tiers. Il convient de documenter leur collecte : quelles unités de mesure ont été utilisées ? Quel est le contexte ? Quelles abréviations et paramètres ont été employés ? Il est également nécessaire de les décrire le mieux possible - notamment grâce à des métadonnées complètes et précises.



Afin d'apporter une validation scientifique à ces données, il peut être utile de les soumettre à une **validation** par ses pairs. Il existe en effet également des systèmes de **reviewing** pour les données de recherche (comme les *data papers*).



La **reproductibilité** des recherches permet d'accroître la crédibilité des résultats.



Compréhensibles, fiables et reproductibles, les données de recherche seront plus à même d'être **réutilisables** et citées par d'autres chercheurs. Il convient alors de leur appliquer une licence utilisateur qui encadrera strictement cette réutilisation par les autres chercheurs.



BIEN-SÛR, IL N'EST PAS TOUJOURS POSSIBLE D'APPLIQUER TOUTES CES ÉTAPES AUX JEUX DE DONNÉES EN RAISON DE CONTRAINTES DIVERSES ET VARIÉES. MAIS IL S'AGIT D'UN IDÉAL VERS LEQUEL TENDRE LE PLUS POSSIBLE

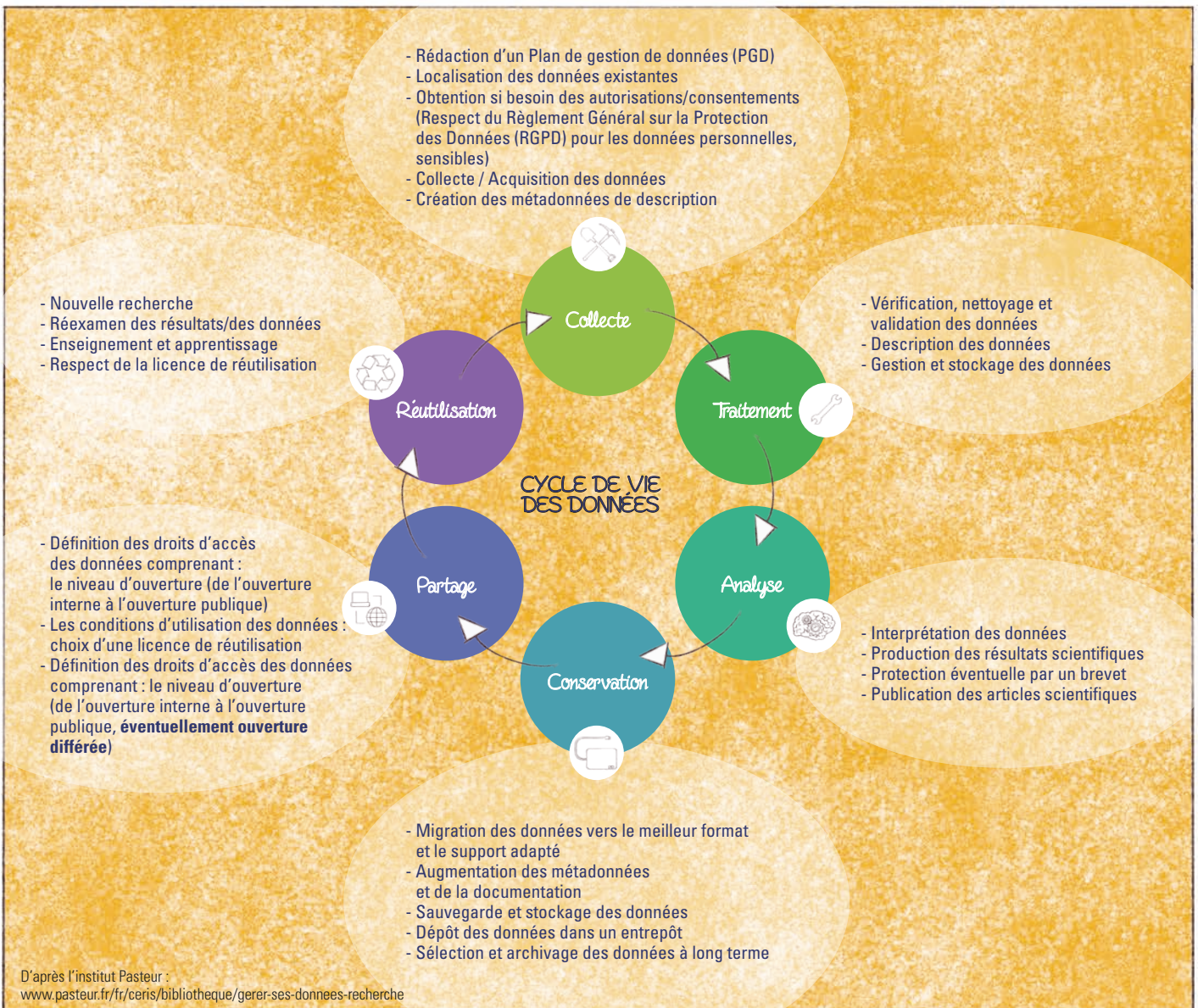


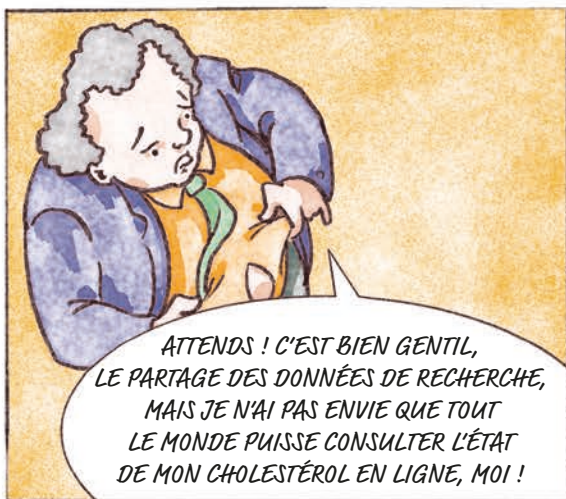
OH LÀ LÀ ! MAIS COMMENT JE VAIS FAIRE POUR SUIVRE TOUTES CES ÉTAPES, MOI ? JE N'Y COMPRENDS RIEN !



Pas de panique ! Il est temps de passer à la pratique !

2. Comment gérer et diffuser des données de recherche





QUESTIONS JURIDIQUES : DONNÉES OUVERTES OU NON OUVERTES ?

Données diffusables

- Communication libre si (cf. loi pour une République numérique, oct. 2016) :
 - données issues d'une activité de recherche financée au moins pour moitié par des fonds publics
 - non protégées par un droit spécifique
 - rendues publiques par le chercheur ou l'établissement (établissement décide quelles données seront ouvertes, le lieu et les conditions de dépôt)
- Communication obligatoire pour certaines données géographiques et environnementales (cf. convention Inspire et convention d'Arrhus)

Données diffusables sous conditions

- Données présentant des risques pour la protection du potentiel scientifique et technique de la nation (cf. laboratoire dit « unité protégée »)
- Zones à régime restrictif (ZRR) : accès physique et numérique soumis à autorisation
- Données protégées par le droit d'auteur et autres droits de propriété intellectuelle
- Données personnelles (cf. Règlement général sur la protection des données (RGPD))

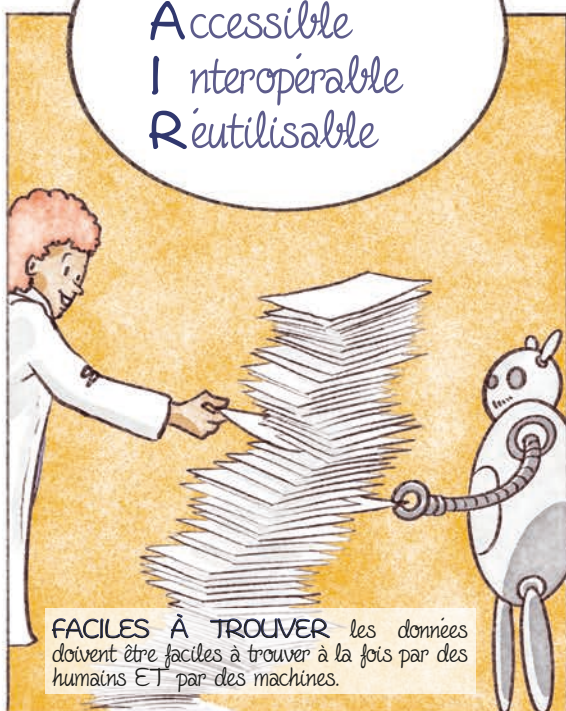
Données interdites à la diffusion

- Données présentant des risques pour la protection du secret de la défense nationale
- Données présentant des risques pour la sécurité de l'Etat, la sécurité publique, la sécurité de l'établissement
- Secret professionnel ou confidentialité (secret médical, secret de l'instruction, secret bancaire et fiscal)

Afin d'atteindre cet objectif, le partage des données doit respecter un ensemble de bonnes pratiques permettant leur découverte et leur utilisation par des humains – mais aussi par des machines. Cet engagement nécessaire est résumé par l'acronyme : **FAIR** ».

D'après Dominique L'Hostis, Du plan de gestion au data paper, juin 2019
<https://gricad-media.univ-grenoble-alpes.fr/video/plan-gestion-donnees-au-data-paper>

Facile à trouver
Accessible
Interopérable
Réutilisable



ACCESSIBLES : Les données et les métadonnées doivent être stockées durablement, avec des accès et/ou des téléchargements facilités, en spécifiant les conditions d'accès et d'utilisation.



Cela nécessite de rendre les données accessibles par leurs identifiants via un protocole de communication standardisé (ex : HTTPS, API REST). Il est conseillé d'utiliser au maximum des protocoles ouverts, libres, pouvant être implémentés de manière universelle.

AH ! AH ! ON FAIT DE LA RÉSISTANCE ?
VINGT-SIXIÈME FOIS QUE L'ON TENTE
D'IMPORTER LES DONNÉES ET ÇA PLANTE
ENCORE ? VAS-Y, RECOMMENCE, CHARLOTTE !
ON VA FINIR PAR LES AVOIR À L'USURE !

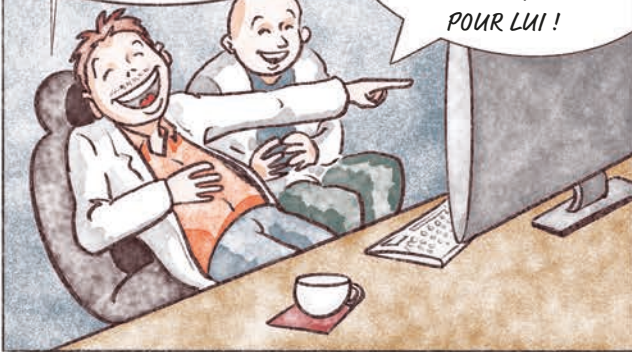
D'ACCORD !
MAIS TU ME PROMETS
QU'AU BOUT DE LA
CINQUANTIÈME TENTATIVE,
ON ARRÊTE ? JE N'AI PAS
ENVIE DE DORMIR AU
LABORATOIRE !...



Cela suppose également de prévoir des protocoles permettant l'authentification et l'autorisation si besoin, par exemple pour limiter ou restreindre la consultation de données sensibles, stratégiques ou confidentielles à un certain type d'utilisateurs identifiés.

HÉ ! CE SONT LES DONNÉES
SUR L'ÉVOLUTION DU
CHOLESTÉROL DE JEAN-PAUL !
ALLEZ, TÉLÉCHARGE !
ON VA SE MARRER !

OULALA ! C'EST PAS
BRILLANT ! VA FALLOIR
PENSER À LE REMETTRE
AU RÉGIME ! PAS DE POT
DE FIN D'ANNÉE
POUR LUI !



Enfin, il convient de rendre les métadonnées accessibles même quand les données ne le sont plus, ce qui suppose de mettre en place des protocoles d'archivage pérenne.

JE N'Y COMPRENDS
PLUS RIEN !
JE SUIS ALLÉE SUR LEUR
ENTREPÔT HIER ET IL Y AVAIT
BIEN TOUTES LEURS DONNÉES
EN LIGNE ! ET AUJOURD'HUI ?
RIEN. NADA. ERREUR 404.

HUM ! APRÈS LE FANTÔME
DE L'OPÉRA, VOICI LE FANTÔME
DE L'ACADÉMIE DES SCIENCES...



MOUAIS...
N'EMPÊCHE QUE
MON CHOLESTÉROL,
IL A BAISSÉ
LE MOIS DERNIER !



INTEROPÉRABLE : Les données et les métadonnées doivent être téléchargeables, utilisables, intelligibles et combinables avec d'autres données, par des humains ET par des machines



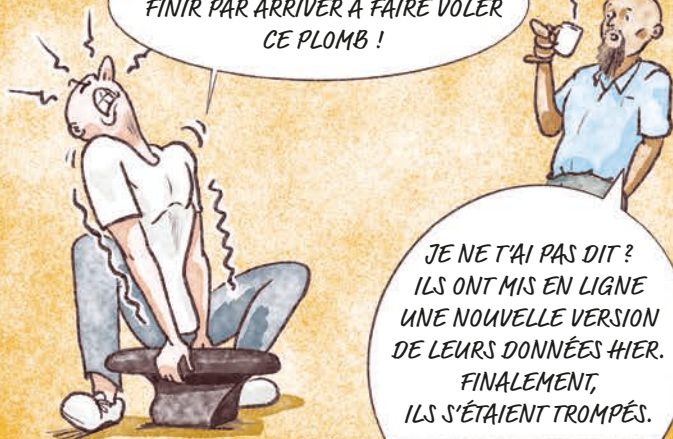
Il est souhaitable que les données et les métadonnées utilisent un langage formel, accessible et partage largement applicable à la représentation des connaissances, comme les technologies du Web sémantique. Il est conseillé de s'appuyer sur des ontologies et des vocabulaires contrôlés standards.

JE N'ARRIVE PLUS RIEN À RETROUVER DANS LES DONNÉES DU LABORATOIRE « COCCINELLES ». POUR DÉCRIRE LEUR ÉTUDE SUR LES VOITURES, ILS ONT UTILISÉ LES MOTS « AUTOMOBILE », « VÉHICULE », « BAGNOLE » ET « CHAR ». ILS LES ONT MÊME QUALIFIÉES À UN MOMENT DE « CAISSES » ET DE « TACOTS » ! DU COUP, JE NE TROUVE PLUS RIEN !



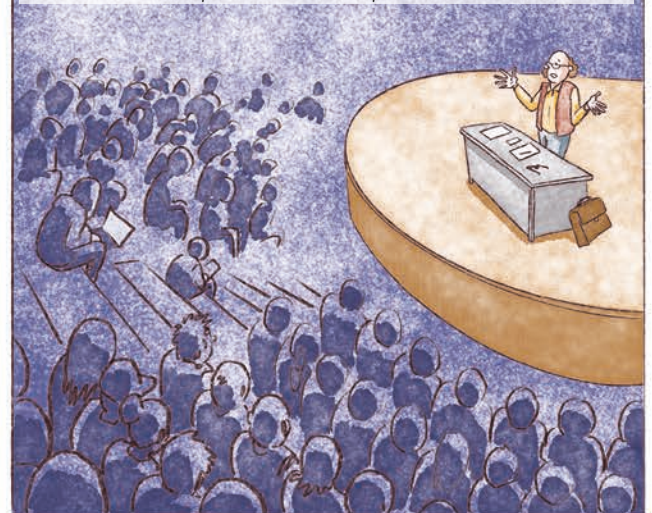
Les données et les métadonnées peuvent inclure des liens vers d'autres (méta)données, des versions antérieures ou plus récentes, ou encore des données complémentaires ou des articles citant les données.

GNNN !!! JE TE DIS QU'ON VA FINIR PAR ARRIVER À FAIRE VOLER CE PLOMB !

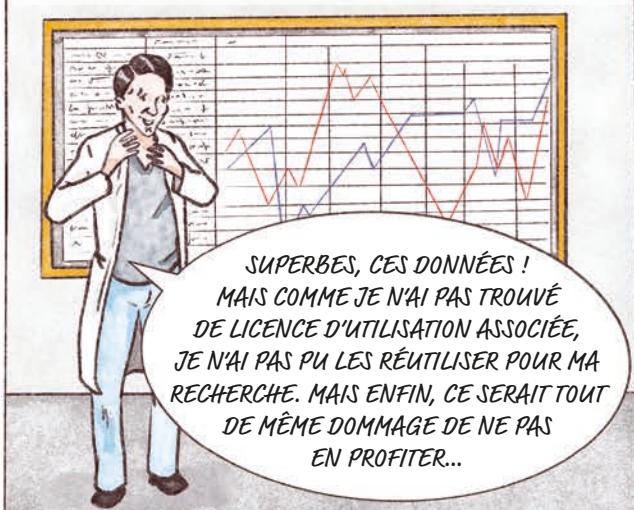


JE NE T'AI PAS DIT ? ILS ONT MIS EN LIGNE UNE NOUVELLE VERSION DE LEURS DONNÉES HIER. FINALEMENT, ILS S'ÉTAIENT TROMPÉS.

RÉUTILISABLES : Les données et les métadonnées doivent présenter des caractéristiques rendant les données réutilisables pour de futures recherches ou pour d'autres finalités (enseignement, innovation, reproduction et transparence de la science).



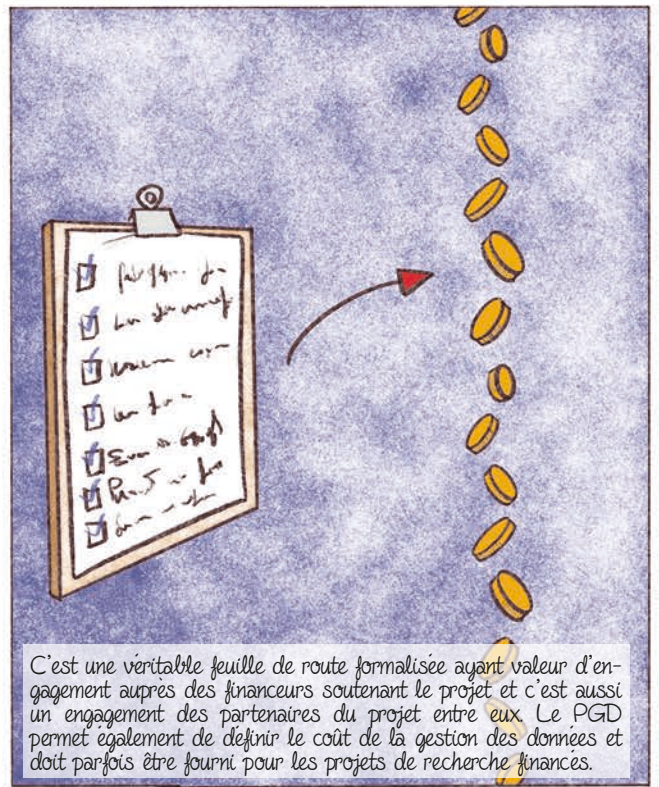
À cette fin, associer aux jeux de données une licence d'utilisation explicite et accessible, les associer à leur provenance en respectant les standards des communautés indiquées permet d'en faciliter la réutilisation ultérieure.



SUPERBES, CES DONNÉES ! MAIS COMME JE N'AI PAS TROUVÉ DE LICENCE D'UTILISATION ASSOCIÉE, JE N'AI PAS PU LES RÉUTILISER POUR MA RECHERCHE. MAIS ENFIN, CE SERAIT TOUT DE MÊME DOMMAGE DE NE PAS EN PROFITER...



OH LÀ LÀ ! IL Y A PLEIN DE CHOSES AUXQUELLES IL FAUT PENSER ! JE SUIS SÛR QUE JE VAIS EN OUBLIER LA MOITIÉ EN ROUTE... JE N'AI PLUS LA MÉMOIRE DE MES 20 ANS, MOI !



L'ÉA CHERCHEUR·SE



Le chercheur est responsable de la collecte, de la description et du découpage des données en différents jeux cohérents.

L'INGÉNIEUR·E PROJET



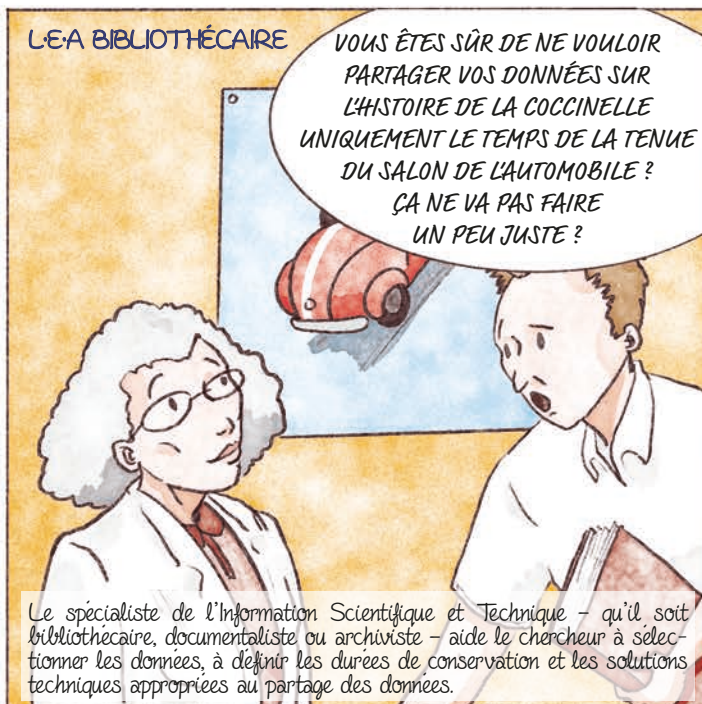
Quand un ingénieur-projet est nommé sur le projet, il devient responsable de la coordination des actions menées autour des données.

L'INFORMATICIENNE



L'informaticien est l'interlocuteur à privilégier pour les questions de stockage et de sécurisation des données, pour les aspects liés aux infrastructures et sur les coûts.

L'ÉA BIBLIOTHÉCAIRE



Le spécialiste de l'Information Scientifique et Technique – qu'il soit bibliothécaire, documentaliste ou archiviste – aide le chercheur à sélectionner les données, à définir les durées de conservation et les solutions techniques appropriées au partage des données.

L'ÉA JURISTE

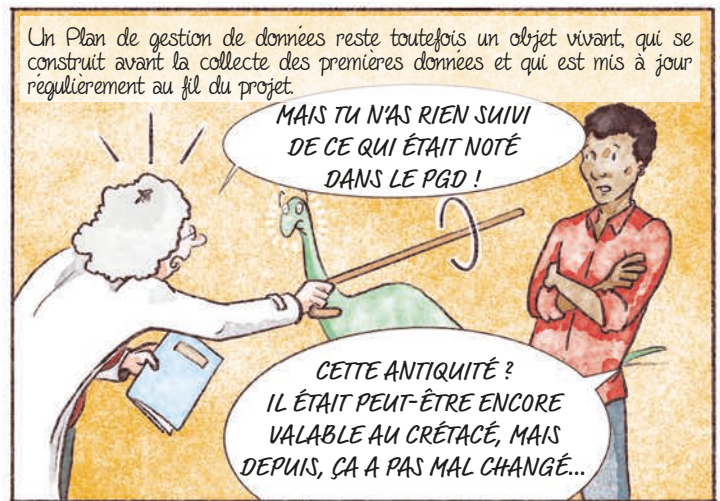
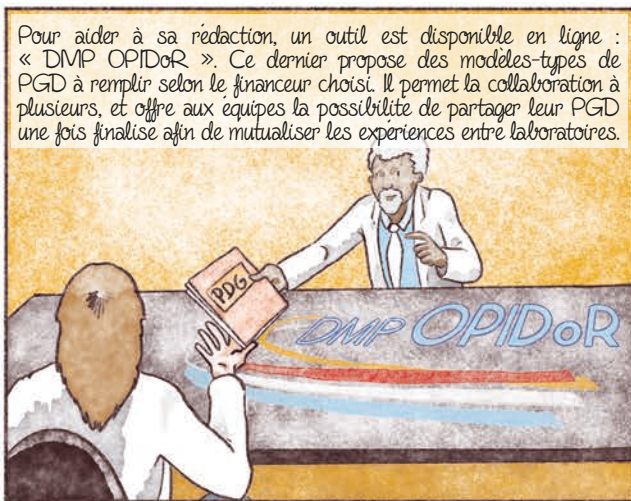


Le juriste et le DPD apportent leurs conseils sur les procédures d'anonymisation des données, sur les licences à utiliser, ou encore sur les permissions d'accès qui seront accordées aux différents groupes d'utilisateurs (accès ouvert, restreint ou sous embargo). Il conseille aussi sur la protection des données en vue de leur valorisation économique sous forme de brevets ou autres.

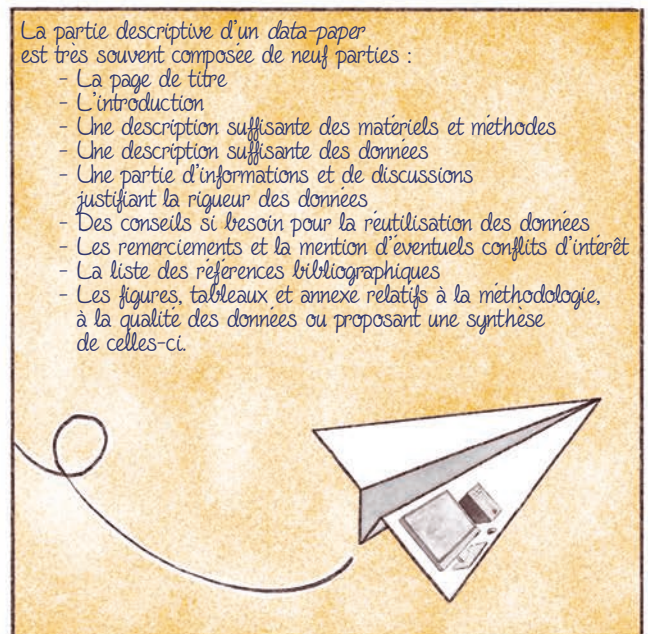
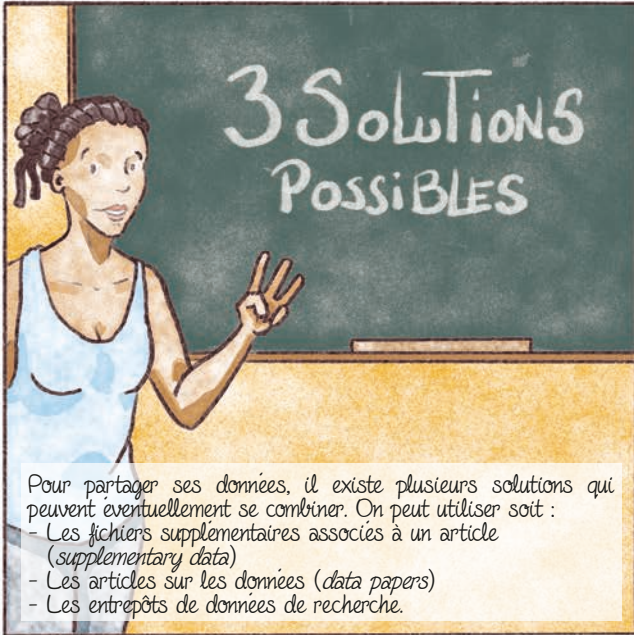
LES AUTRES PROFESSIONNEL·LES



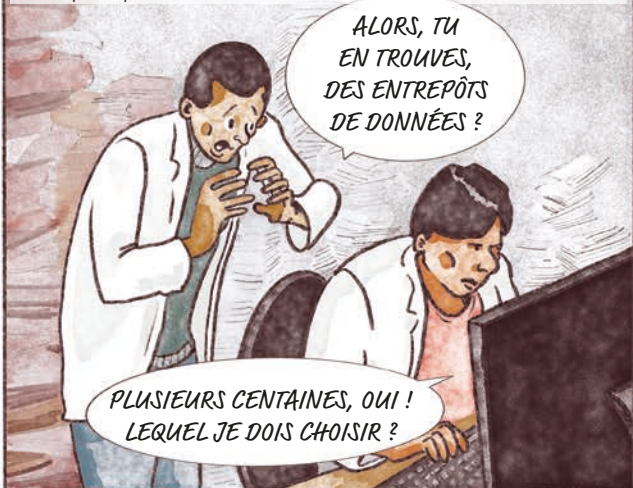
Cette liste de professionnels pouvant intervenir dans l'élaboration d'un plan de gestion de données (PGD) n'est bien-sûr pas exhaustive. Beaucoup d'autres compétences peuvent être convoquées.



3. Partager les données de recherche : quels outils ?



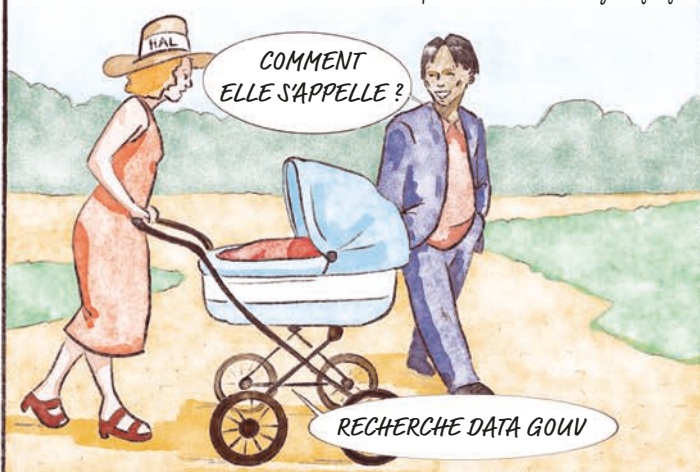
La troisième solution, qui est aussi la plus conseillée, consiste à déposer ses jeux de données dans un entrepôt dédié. Un entrepôt de données, c'est un réservoir de données de recherche, brutes ou dérivées, qui peuvent être retrouvées et réutilisées grâce à une description par des métadonnées.



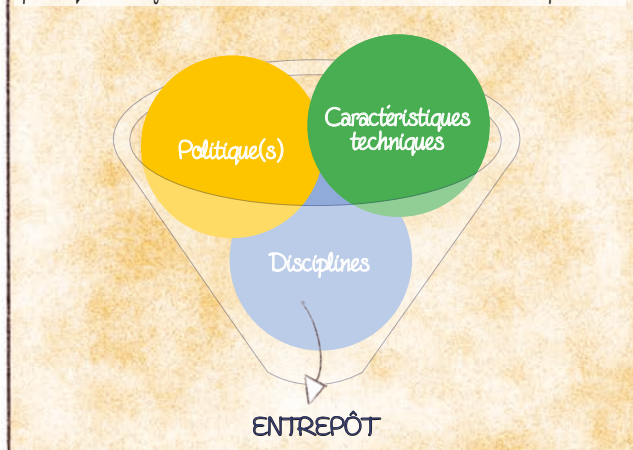
Eh oui ! Des entrepôts de données, dans le monde, il en existe plusieurs milliers. Ils se divisent en trois grandes catégories : les entrepôts disciplinaires, généralistes et institutionnels.



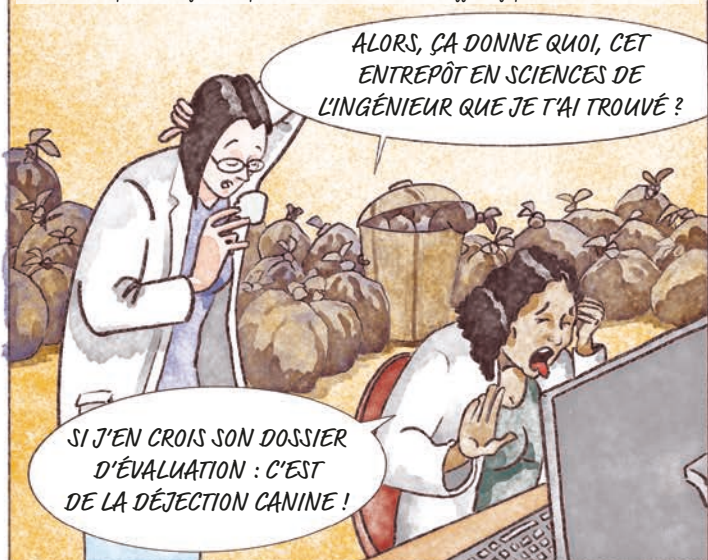
En France, il existe depuis fin 2021 un entrepôt national des données de la recherche, financé par le Ministère de l'enseignement supérieur et de la recherche, nommé Recherche Data Govu : <https://recherche.data.govu.fr/fr>



On l'a compris : le choix du bon entrepôt pour ses données est donc fondamental. Il obéit à trois critères principaux : les caractéristiques techniques de l'entrepôt, l'adéquation aux injonctions politiques des financeurs ou des éditeurs, et le critère disciplinaire.



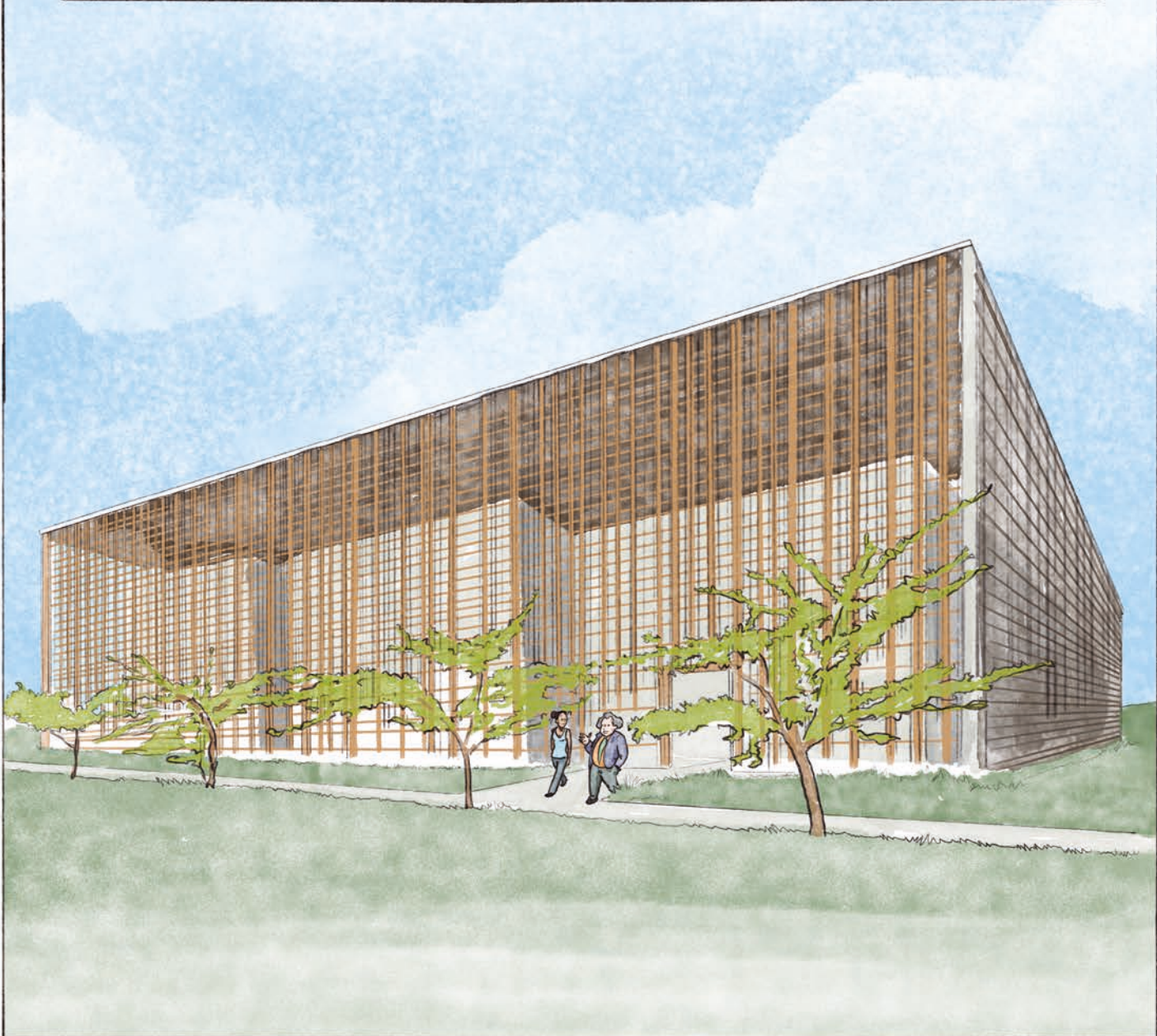
En ce qui concerne les caractéristiques techniques des entrepôts, un système de certification a été mis en place afin de permettre aux chercheurs de s'y retrouver. Ces certifications se basent sur l'évaluation de l'entrepôt par un jury indépendant sur une base déclarative. L'accès au dossier d'évaluation est rendu public afin de permettre un contrôle effectif par les utilisateurs



La certification permet notamment d'attester de la conformité de l'entrepôt avec les principes FAIR. Il en existe actuellement trois niveaux : le « Core Trust Seal », s'obtient après avoir répondu favorablement à 16 critères (134 centres étaient certifiés dans le monde en 2023). Le Nestor Seal requiert, lui, un engagement sur 34 critères (5 centres certifiés dans le monde en 2023). Enfin, pour décrocher l'ISO +, il faut se conformer à une centaine de critères (un seul centre certifié dans le monde en 2023) !







RÉFÉRENCES BIBLIOGRAPHIQUES :

Borgman, C. L. 2020. *Qu'est-ce que le travail scientifique des données ?* (traduit par C. Matoussowsky). OpenEdition Press. 10.4000/books.oep.14692

Callisto Formation. Fondation UNIT.
https://callisto-formation.fr/?theme=boostplus_c06&redirect=0

CoopIST : délégation à la formation scientifique et technique, CIRAD.
Gérer les données de la recherche.
<https://doi.org/10.18167/COPIST/0005>

DoRANum – Données de la recherche : Apprentissage Numérique.
<https://doranum.fr/>

Ouverture des données de recherche – Guide d'analyse
du cadre juridique en France – V2. (2017). Ouvrir la science !
Comité pour la science ouverte.
www.ouvrirlascience.fr/ouverture-des-donnees-de-recherche-guide-danalyse-du-cadrejuridique-en-france-v2

Partager les données liées aux publications scientifiques –
Guide pour les chercheurs. (2022). Ouvrir la Science !
Comité pour la science ouverte.
www.ouvrirlascience.fr/partager-les-donnees-liees-aux-publications-scientifiques-guidepour-les-chercheurs

AUTEUR

Marie Latour, directrice adjointe du
Service commun de la documentation de l'Université de Guyane

RESPONSABLE SCIENTIFIQUE

Annaïg Mahé, enseignante-chercheuse à l'URFIST de Paris

DESSINATEUR

Olivier Copin

GRAPHISME

Bénédicte Sauvage (BCOM)

RELECTEURS SCIENTIFIQUES

Cyril Heude (data librarian à SciencePo Paris),
Romain Féret (directeur de Média Normandie),
Amélie Barrio (Co-responsable de l'URFIST Occitanie)

**Projet co-financé par l'URFIST de Paris
(Ecole nationale des Chartes - PSL)**

