



**HAL**  
open science

## Principes et enjeux de la science ouverte

Gaëlle Leroux

► **To cite this version:**

Gaëlle Leroux. Principes et enjeux de la science ouverte. Master. Principes et enjeux de la science ouverte, Lyon, France. 2023, pp.38. hal-04137763

**HAL Id: hal-04137763**

**<https://hal.science/hal-04137763>**

Submitted on 22 Jun 2023

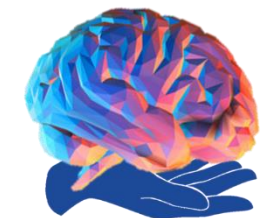
**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

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



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

## Principes et enjeux de la science ouverte



Soutien  
Méthodologique  
aux projets  
d'imagerie

Gaëlle Leroux, PhD

Research engineer

[gaelle.leroux@cnrs.fr](mailto:gaelle.leroux@cnrs.fr)

 <https://orcid.org/0000-0003-2729-4945>



Research data College

# 7 piliers de la science ouverte

- publications ouvertes (*open access*)
- données ouvertes (*open data*)
- sources et pratiques ouvertes (*open source & open materials*)
- éducation ouverte (*open educational resources*)
- science citoyenne (*citizen science*)
- évaluation responsable (*assessment*)
- intégrité scientifique (*scientific integrity*)

[Source](#)

# Mon parcours professionnel

## Formation initiale à Caen

DEUG de  
Biologie  
(2 ans)

Maîtrise des sci. &  
techniques –  
neurosciences  
(2 ans)

*Traitement de données  
en SPECT & MEG de  
patients épileptiques*

DEA de  
Neuropsychologie  
(1 an)

*Collecte & traitement de données  
comportementales, en EEG, MEG et  
IRMf chez des adultes sains*

Doctorat de  
neurosciences  
(3 ans)

## Expériences professionnelles

- **Post-doc (CDD) en neurosciences cognitives développementales** (2 ans) à Stockholm
- **Ingénieure de recherche biologiste en analyse de données**  
Univ. Paris Descartes (4 ans), Univ. Bordeaux (3 ans), CNRS (9 ans) Bordeaux puis Lyon

## Spécialité

- **Conception et implémentation de protocoles de recherche** (rédaction de protocoles, demandes éthiques, choix logiciels, prog. logicielle de chronométrie, design des stimuli)
- **Acquisition de données** (comportementales, en neuroimagerie non invasive)
- **Analyse des données et valorisation** (gestion de bases de données, pré- & post-traitements, statistiques, publications, diffusion de données)

# Ce qui m'a sensibilisé à la science ouverte

## Évolution des jeux de données étudiés

2001 6 participants avec 2 modalités d'imagerie + qqes données cliniques  
(labo 1) données stockées sur mon ordinateur du labo  
1 pers traitant données

2002-07 < 25 participants avec 1-3 modalités d'imagerie + quelques tests comportementaux  
(labo 2) données stockées sur CD + 1 serveur du labo  
(labo 3) 2 pers traitant données

2008-11 60 participants avec 1 modalité d'imagerie + batterie tests comportementaux  
(labo 4) données stockées sur 2 serveurs du labo  
2 pers traitant données

2011-18 450 participants avec 4 modalités d'imagerie + batterie de tests comportementaux  
(labo 5) + données génétiques  
données stockées sur 3 serveurs du labo  
6 pers traitant données

*Si je devais répliquer les résultats & figures dans les 15 articles publiés...aïe aïe aïe*

# Pourquoi on parle d'« ouvrir » la science ?

*Elle est issue d'un mouvement de contestation de la situation de monopoles des grands éditeurs par la promotion de l'open access.*

La notion de science ouverte s'est développée afin de caractériser un **double mouvement d'ouverture des résultats de la science** (publications, données de la recherche lorsque c'est possible), **mais également une nouvelle façon de faire de la recherche**, passant par l'ouverture des processus, des codes, des méthodes ou des protocoles.

Source : guide Couperin

« La science ouverte est un mouvement qui consiste à **ouvrir le processus de recherche à tous types d'acteurs** (partenaires, citoyens) et à **rendre accessibles et réutilisables les produits de la recherche** (publications scientifiques, données, logiciels, etc.) **par la communauté scientifique et la société.**

Partager les résultats de recherche, non seulement entre scientifiques, mais aussi avec la société, est en effet le meilleur moyen de **faire progresser la connaissance** et de **développer des relations de confiance avec les différents acteurs.** » INRAE. <https://www.inrae.fr/inrae-engage-science-ouverte>

# Contours de la science ouverte

**La science ouverte n'est pas différente de la science traditionnelle.**

**Elle signifie simplement que l'on effectue ses recherches de manière plus transparente et plus collaborative.**

**La science ouverte s'applique à toutes les disciplines de recherche.**

[Source](#)

# Un historique rapide de la science ouverte

## 1<sup>ères</sup> initiatives : 1990's

- 1991 Création de l'archive ouverte [ArXiv](#)
- 1992 Création des licences [Creative Commons](#)

## 3 textes fondateurs de l'Open Access au niveau international :

- 2002 Initiative de Budapest** sur l'Open Access
- 2003 Déclaration de Berlin** sur le libre accès à la connaissance
- 2003 Bethesda statement** on Open Access publishing



[Source](#)



# Un historique rapide de la science ouverte

1<sup>ères</sup> initiatives : 1990's



Dates clefs en France

2002

HAL : archives ouvertes des publications scientifiques



# Un historique rapide de la science ouverte

[Source](#)

1<sup>ères</sup> initiatives : 1990's

## Dates clefs en France

2002 HAL : archives ouvertes des publications scientifiques

2016 [Loi pour une République numérique](#) (ou “Loi Lemaire”)



## Publications

L'article 30 **garantit le droit aux chercheurs de diffuser la version acceptée pour publication** (postprint ou AAM) **de leur article** **au plus tard 6 mois** après la première date de publication par l'éditeur **en sciences, techniques et médecine (STM)** **et 12 mois** pour les chercheurs **en sciences humaines et sociales (SHS)**.

**Cette loi s'applique, quel que soit le contrat signé avec l'éditeur de la revue, que la revue soit française ou étrangère.**

# Un historique rapide de la science ouverte

[Source](#)

1<sup>ères</sup> initiatives : 1990's

## Dates clefs en France

2002 HAL : archives ouvertes des publications scientifiques

2016 [Loi pour une République numérique](#) (ou “Loi Lemaire”)



## Données

« II. Dès lors que les données issues d'une activité de recherche **financée au moins pour moitié par des dotations de l'Etat**, des collectivités territoriales, des établissements publics, des subventions d'agences de financement nationales ou par des fonds de l'Union européenne **ne sont pas protégées par un droit spécifique ou une réglementation particulière** et qu'elles **ont été rendues publiques par le chercheur**, l'établissement ou l'organisme de recherche, leur réutilisation est libre. »

« III. L'éditeur d'un écrit scientifique mentionné au I ne peut limiter la réutilisation des données de la recherche rendues publiques dans le cadre de sa publication. »

« IV. Les dispositions du présent article sont d'ordre public et toute clause contraire à celles-ci est réputée non écrite. »

Mise en place du principe d'**Open Data** par défaut.

# Un historique rapide de la science ouverte

[Source](#)

1<sup>ères</sup> initiatives : 1990's

## Dates clefs en France

2002 HAL : archives ouvertes des publications scientifiques

2016 [Loi pour une République numérique](#) (ou “Loi Lemaire”)

2018-21 [1<sup>er</sup> plan national pour la science ouverte](#) (“Plan S”)

2021-24 [2<sup>nd</sup> plan national pour la science ouverte](#)

Codes &  
logiciels



3<sup>ème</sup> axe :  
« ouvrir et promouvoir les codes sources produits par la recherche »



Et en + :



meilleur taux de citation pour les articles et données de la recherche librement accessibles

[Source](#)

# Crise de la reproductibilité des résultats

RESEARCH ARTICLE

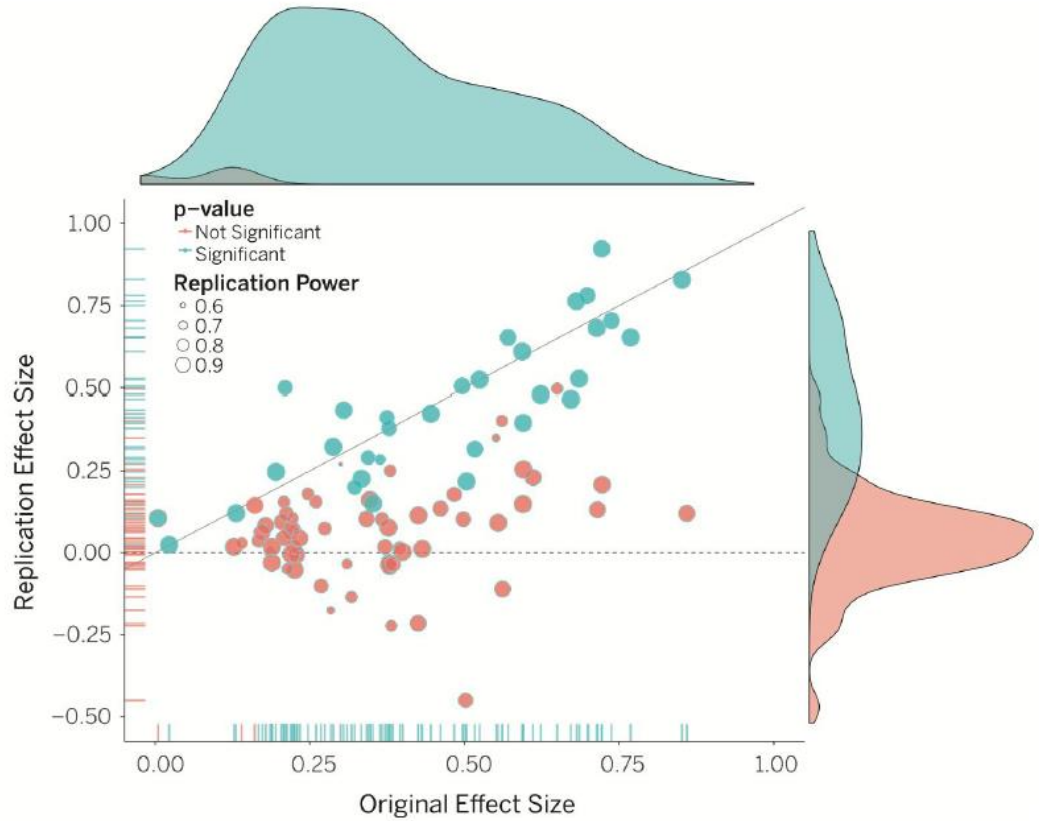
## Estimating the reproducibility of psychological science

Open Science Collaboration<sup>\*,†</sup>  
+ See all authors and affiliations

Science 28 Aug 2015  
Vol. 349, Issue 6251, aac4716  
DOI: 10.1126/science.aac4716

The *Reproducibility project* set out to replicate 100 experiments published in high-impact psychology journals.

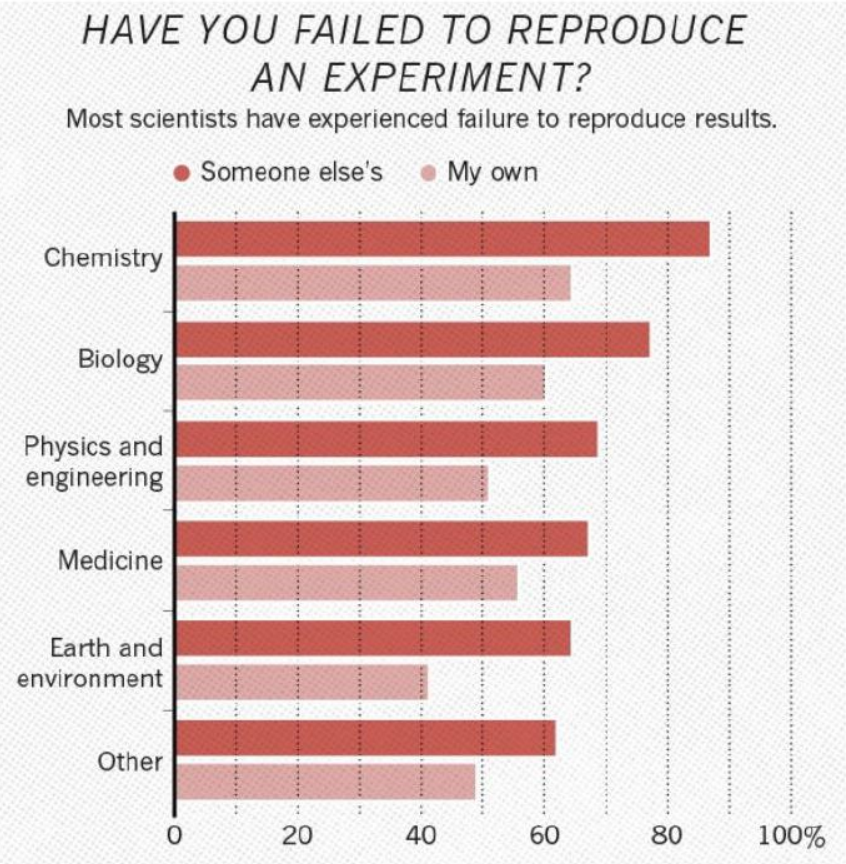
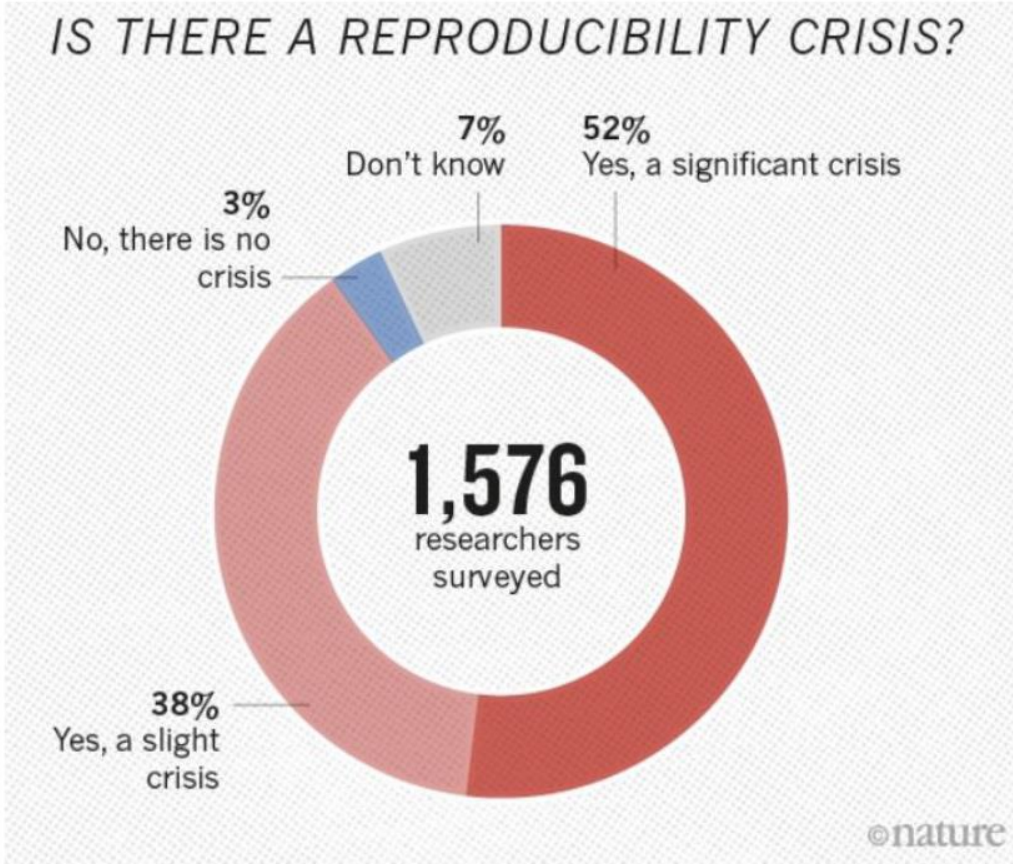
About one-half to two-thirds of the original findings could not be observed in the replication study.



Extrait de la [formation IFB "Science ouverte et PGD"](#). Nov. 2021



# Crise de la reproductibilité des résultats



"1,500 scientists lift the lid on reproducibility". Nature. 533: 452-454 - 2016

[Source](#)



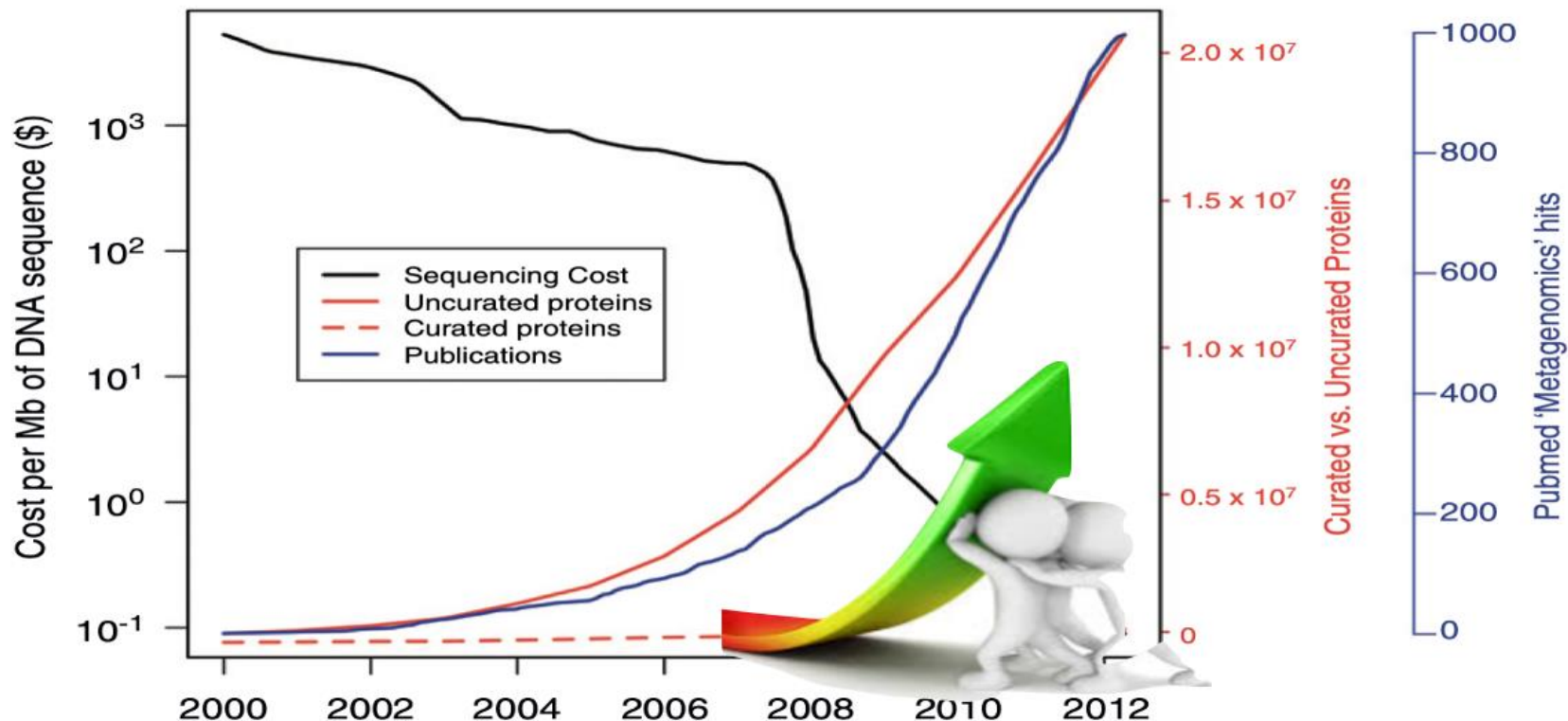
# Le déluge de données en science

Les techniques à haut débit, une révolution qui provoque un déluge de données

Génome humain :

en 1990 = **13 ans et 3 Milliards \$**

en 2015 = **quelques heures et 1000 \$**



Par conséquent :

- La quantité de données à stocker et analyser explose
- Le rendement d'analyse chute

[Source](#)

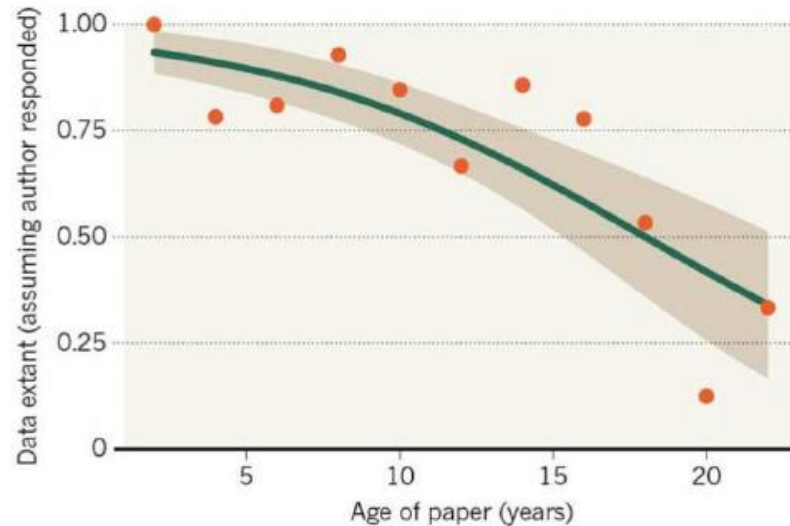


# Les ravages du temps

## Data Entropy

### MISSING DATA

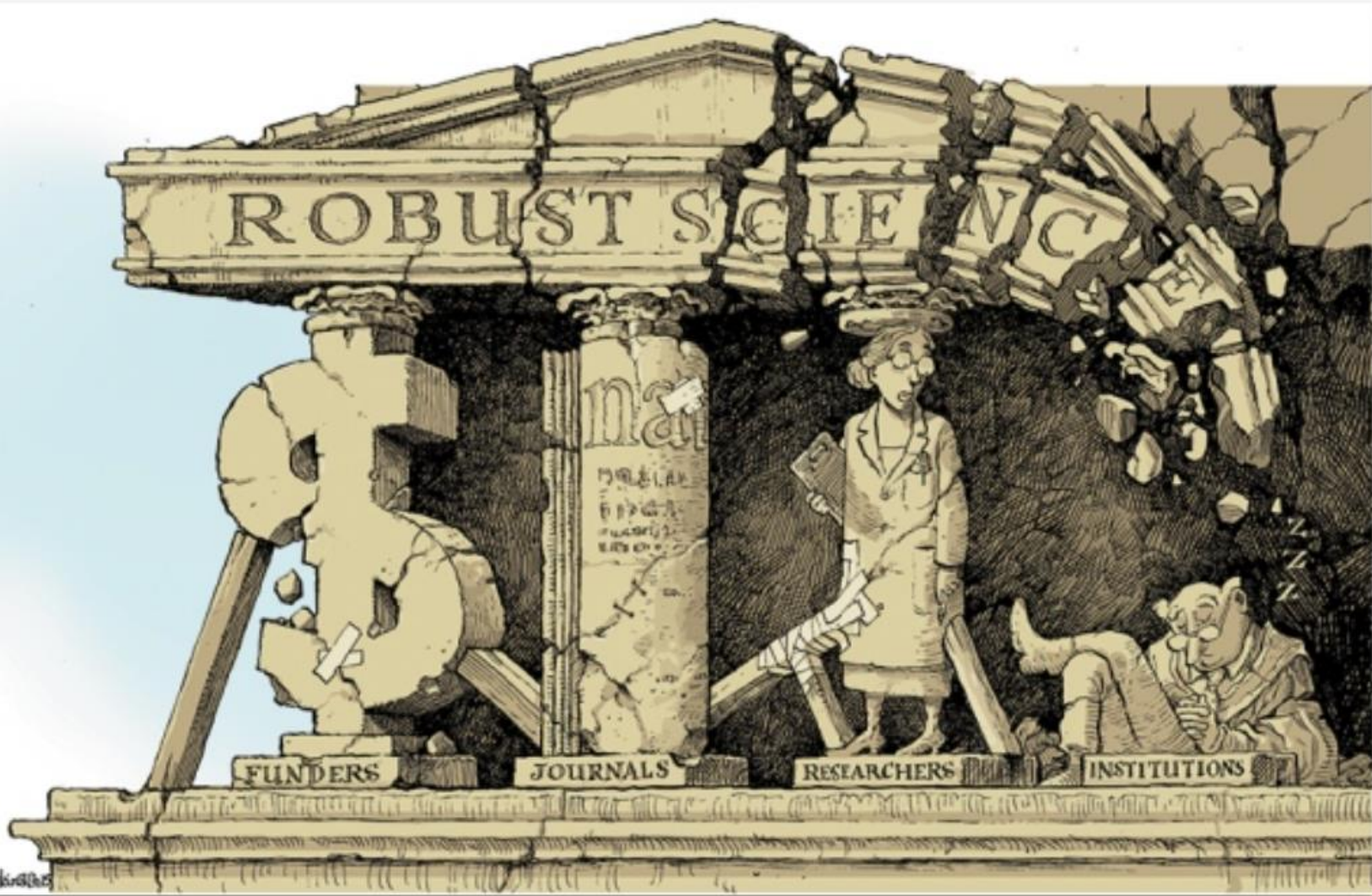
As research articles age, the odds of their raw data being extant drop dramatically.



Vines, T. H. et al. *Curr. Biol.* <http://dx.doi.org/10.1016/j.cub.2013.11.014> (2013).

[Source](#)

DataONE

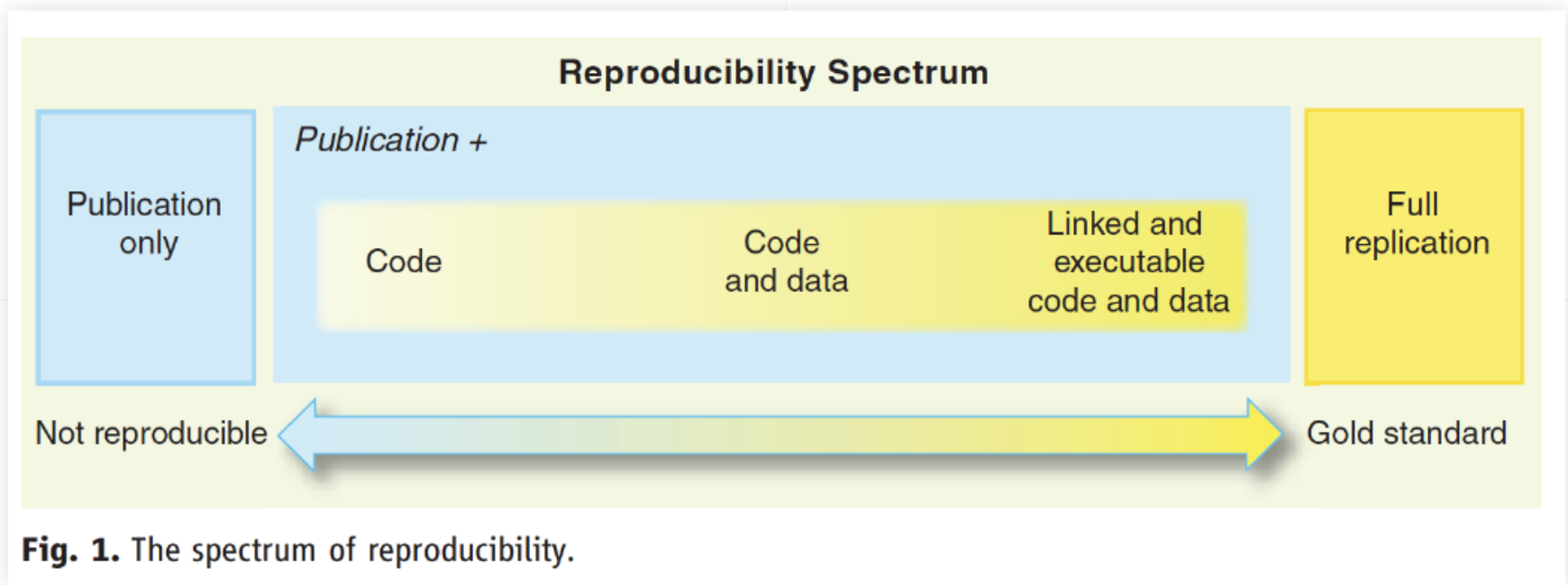


Reproduced by courtesy of Nature Publishing Group.  
Nature, September 1, 2015.

[Source](#)

# Pas si simple pour autant

*Reproducibility is not black an white, but rather a continuum*



Peng, Science, 2011

[Source](#)

# Pourquoi une si faible reproductibilité des résultats ?

- **Absence d'une "culture de la reproductibilité" profondément enracinée**
- **Manque de formation aux méthodes de recherche robustes**
- **Pratiques de recherche**
  - Mauvaise organisation
  - Mauvaise documentation
  - Mauvaise automatisation
  - Mauvaise diffusion
  - Données ne sont pas disponibles

[Source](#)

# Comment l'organisation aide à la reproductibilité

## Principes d'organisation

1. Mettre chaque projet dans son propre répertoire de projet
2. Garder ensemble ce qui doit l'être
3. Gardez les données source
4. Rendez les noms de fichiers lisibles par des humains
5. Rendez les noms de fichiers lisibles par des machines
6. Rendre les noms de fichiers faciles à trier
7. Sauvegardez vos données

[Source](#)

# Comment l'organisation aide à la reproductibilité

Ziemann *et al. Genome Biology* (2016) 17:177  
DOI 10.1186/s13059-016-1044-7

Genome Biology

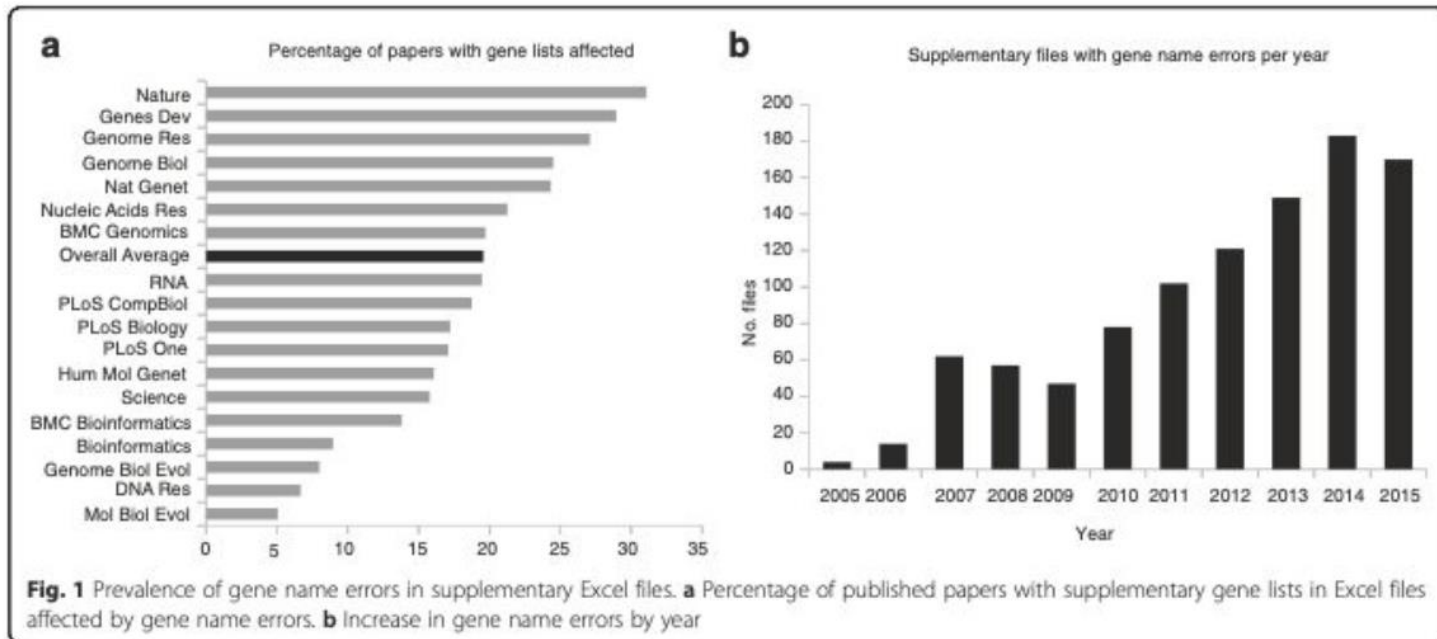
COMMENT

Open Access



## Gene name errors are widespread in the scientific literature

Mark Ziemann<sup>1</sup>, Yotam Eren<sup>1,2</sup> and Assam El-Osta<sup>1,3\*</sup>



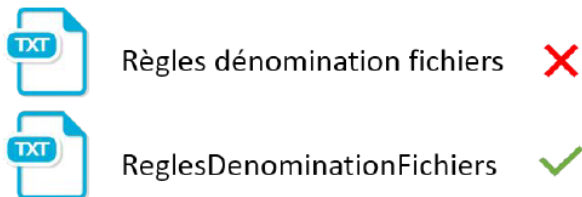
Source



# Comment l'organisation aide à la reproductibilité

## DONNER UN NOM BREF ET EXPLICITE et ...

Pas d'espace Ni de caractères spéciaux (& / + > : ? % ...)



Dates au format **AAAAMMJJ** (année, mois, jour)



Versionnez



Rangez



Et documentez vos règles !

REGLES DE NOMMAGE DES FICHIERS	
EGE-10-Sec7.2.2a-v0.7	Domaine: Systèmes Information
Page: 1/13	



REPUBLIQUE ET CANTON DE GENEVE  
Collège spécialisé des systèmes d'information

DIRECTIVE TRANSVERSALE

REGLES DE NOMMAGE DES FICHIERS	
EGE-10-Sec7.2.2a-v0.7	Domaine : Systèmes Information
Date : 26.11.2012	Entrée en vigueur : Immédiate
Rédacteur(s): Groupe Records management-archives définitives (RM-Archdét)	Direction/Service transversal(e): CSSI
Responsable(s) de la mise en œuvre: Archivistes de département et d'institution	Approbateur : Collège spécialisé Systèmes d'Information
Date: 21.11.2012	Date: 21.11.2012 /mise à jour de l'annexe : décembre 2015

[Source](#)

# Comment l'organisation aide à la reproductibilité

Type	Format conseillé	Format non conseillé
Document texte	PDF, TXT, ODT	MS Word, RTF
Feuille de calcul	ODS, CSV	MS Excel, PDF, OOXML
Base de données	SQL, SIARD, DB tables (.CSV)	MS Access, dBase (.dbf), HDF5
Données statistiques	SPSS Portable, STATA, XML, CSV, TXT	SAS et R
Images	JPEG, TIFF, PNG	DICOM
Audio	BWF, MXF, Matroska (.mka), FLAC, OPUS	<u>WAVE</u> , <u>MP3</u> , <u>AAC</u> , <u>AIFF</u> , <u>OGG</u>
Video	MXF, MKV	MPEG-4, MPEG-2, AVI, QuickTime (.mov, .qt)
Information géographique	GML, MIF/MID	ESRI Shapefiles, MapInfo, KML
Images géoréférencées Raster	GeoTIFF (.tif, .tiff) ASCII GRID (.asc, .txt)	TIFF World File ESRI GRID

<https://facile.cines.fr/> service de validation des formats

En pratique, on peut souvent travailler avec un format fermé populaire et le **convertir** en format ouvert. **Mais il faut vérifier si la conversion altère les informations, et prendre des mesures de compensation si nécessaire.**

Ex : la conversion XLSX -> CSV perd les mises en forme.

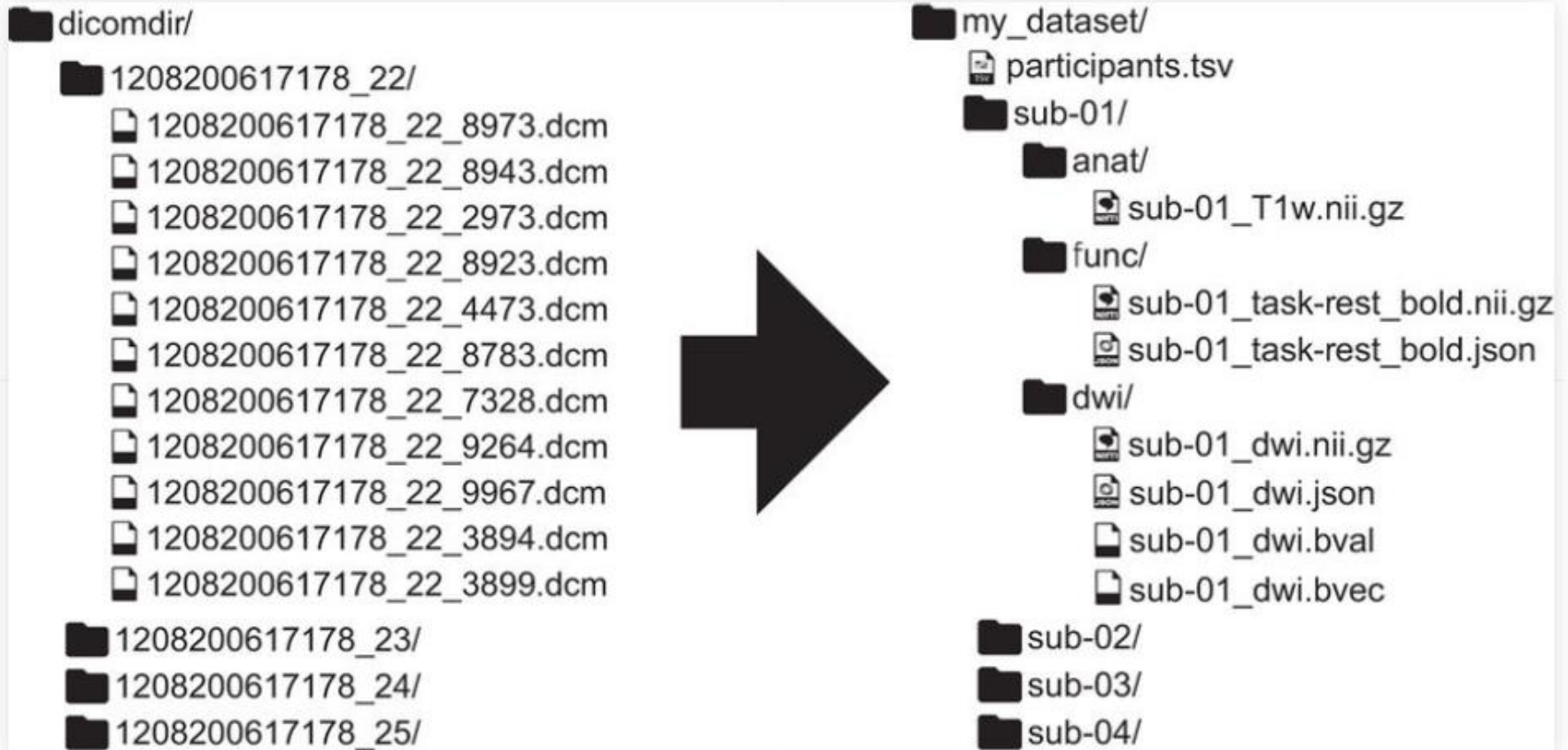
51

[Source](#)



# Comment l'organisation aide à la reproductibilité

## The Brain Imaging Data Structure (BIDS)



Source & further reading: Gorgolewski et al. Sci Data, 2016, <http://bids.neuroimaging.io/>

[Source](#)

# Comment l'organisation aide à la reproductibilité

## Stocker et sécuriser : quels compromis ?

### Comparatif de systèmes de stockage des données

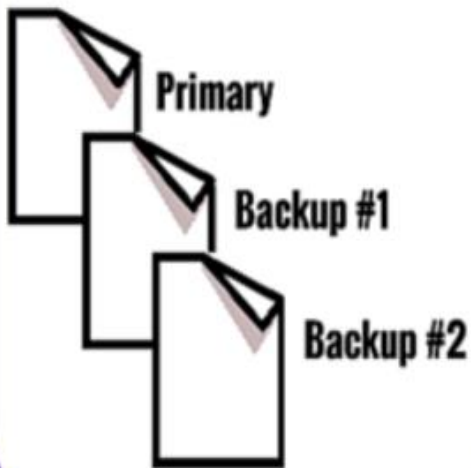
Support de stockage	Sécurité	Accès	Coût	Remarque d'utilisation
 <p>Ordinateur professionnel</p>	<p>★★★★☆</p> <p>Sujet au piratage informatique, aux détériorations et pannes</p>	<p>★★★★☆</p> <p>Pas adapté au partage, nécessite l'utilisation d'un support externe ou d'Internet (mail, cloud...)</p>	<p>★★★★☆</p> <p>Pas de coût supplémentaire ou coût peu important</p>	<ul style="list-style-type: none"> <li>- Pour un stockage temporaire</li> <li>- Nécessité de crypter les données confidentielles et sensibles</li> </ul>
 <p>Support externe</p>	<p>★★★★☆</p> <ul style="list-style-type: none"> <li>- Sujet au vol, à la perte du support</li> <li>- Durée de vie limitée (dégradation du matériel)</li> </ul>	<p>★★★★☆</p> <p>Facilement transportable, il permet de transférer les données vers un autre ordinateur</p>	<p>★★★★☆</p> <p>Pas de coût supplémentaire ou coût peu important</p>	<ul style="list-style-type: none"> <li>- Pour un stockage temporaire</li> <li>- Nécessité de crypter ou de sécuriser physiquement les données confidentielles et sensibles</li> </ul>
 <p>Serveur institutionnel</p>	<p>★★★★☆</p> <p>Stockage fiable, durable et sécurisé (contre le vol, le piratage, les incendies...)</p>	<p>★★★★☆</p> <p>La connexion au serveur institutionnel ne facilite pas le travail avec des personnes extérieures</p>	<p>★★★★☆</p> <p>Coût assez important mais pas forcément répercuté sur l'usager</p>	<ul style="list-style-type: none"> <li>- Pour un stockage plus pérenne</li> <li>- Adapté pour le stockage de données sensibles et des versions « stables » de vos données</li> <li>- Toutes les institutions ne proposent pas ce service</li> </ul>
 <p>Serveur Cloud</p>	<p>★★★★☆</p> <p>On ne sait pas vraiment où sont stockées les données, ni ce qu'elles deviennent</p>	<p>★★★★☆</p> <p>Permet un travail synchronisé avec toutes les personnes ayant été autorisées au partage</p>	<p>★★★★☆</p> <p>Payant à partir d'une certaine limite de stockage</p>	<ul style="list-style-type: none"> <li>- Pour un partage avec des personnes externes à l'institution</li> <li>- Ne pas y mettre de données sensibles ou confidentielles</li> <li>- Pas de contrôle sur la procédure de sauvegarde des données</li> </ul>

Tableau tiré de <http://doranum.fr/le-stockage-des-donnees/>

Source

# Comment l'organisation aide à la reproductibilité

Keep **3** copies of any important file



Store files on **2** different media types

Secure Server



External HD    Secure Cloud



OR



Keep at least **1** copy offsite



Source: NYU libraries

[Source](#)

# Comment la documentation aide à la reproductibilité

## Principes de documentation :

8. Créez une vue d'ensemble de votre projet
9. Rendre les dépendances logicielles explicites
10. Décrire les variables dans un livre de codes
11. Documentez votre code et son utilisation
12. Donner des noms significatifs aux fonctions et aux variables
13. Gardez la trace de la façon dont chaque résultat a été produit

[Source](#)

# Comment la documentation aide à la reproductibilité

Variable Name	Description	Type	Values or Characteristics
patientid	Patient ID	Numeric	Integers
sex	Patient Gender	Numeric	1=Female, 2=Male
proceduredate	Date of Procedure	Date, M/D/YYYY	
dob	Date of Birth	Date, M/D/YYYY	
patientrace	Patient Race	Numeric	0=Not Latino, 1=Latino,
tx	Treatment Group	Numeric	1-Tx 1, 2=Tx 2, 3=Tx 3
dischargeppt	Pt at Discharge	Numeric	2 decimal places
notes	Notes	Character	
location	Location	Numeric	1=local, 2=regional, 3=distant
patientdied	Patient is Dead	Numeric	0=No, 1=Yes
deathdate	Date of Death	Date, M/D/YYYY	
clinicaloutcome	Clinical Outcome	Character	
outcomecomments	Outcome Comments	Character	

Source: <http://domstat.med.ucla.edu/pages/codebook2>

Further reading: Codebook cookbook - A guide to writing a good codebook for data analysis projects in medicine, Creating a codebook - Document, Discover, and Interoperate alliance

Tools: codebook generator for R: memisc

[Source](#)

# Comment la documentation aide à la reproductibilité



# Comment l'automatisation aide à la reproductibilité

## Principes d'automatisation :

14. Éviter les étapes de manipulation manuelle des données
15. Modulariser le code plutôt que de copier et coller.
16. Réutiliser le code plutôt que de le réécrire.
17. Ne commentez pas et ne décommentez pas les sections de code pour contrôler le comportement du programme
18. Ajouter des assertions aux programmes pour vérifier leur fonctionnement
19. Enregistrer la graine du générateur de nombres aléatoires
20. Enregistrer tous les résultats intermédiaires

# Comment la diffusion aide à la reproductibilité

## Principes de diffusion :

21. Stocker les données dans des **formats ouverts**
22. Traiter les données avec des **logiciels ouverts**
23. Éliminer les chemins codés en dur dans votre code
24. Relier les déclarations textuelles aux résultats sous-jacents
25. Soumettez le code et les données à un dépôt émettant des DOI et faites-y référence dans votre article.
26. Rendre explicite la **licence du code et des données**



# Les principes FAIR

F indable A ccessible I nteroperable R eusable



[Source](#)

**Des données/logiciels peuvent être « FAIR » sans être librement accessibles.**

Suivre les principes FAIR permet de s'assurer que ses données sont réutilisables, qu'elles soient partagées ou non.

## Facile à trouver :

- Dans un **entrepôt**
- **Identifiant** unique et pérenne (PID)
- **Métadonnées** riches

## Accessible :

- Définir les conditions d'accès aux données
  - rendre les données accessibles librement
  - sinon : rendre accessibles les métadonnées pour signaler l'existence des données

## Interopérable :

- **Formats ouverts**
- **Mettre à disposition le code source** du logiciel nécessaire pour lire, traiter, analyser les données
- **Privilégier les standards de métadonnées**
- **Indiquer des liens vers d'autres ressources** (autres données, publication...)

## Réutilisable :

- **Licence** de diffusion aux données
- **Documentation** pour décrire les données de façon détaillée, les contextualiser, les rendre compréhensibles...



**Et en + :**

meilleur taux de citation pour les articles et données de la recherche librement accessibles

[Source](#)

# Science ouverte et intégrité scientifique

Principes fondamentaux de l'intégrité scientifique :

- **Fiabilité**, autrement dit garantir la qualité de la recherche, qui transparaît dans la conception, la méthodologie, l'analyse et l'utilisation des ressources
- **Honnêteté**, autrement dit élaborer, entreprendre, évaluer, déclarer et faire connaître la recherche d'une manière transparente, juste, complète et objective
- **Respect** envers les collègues, les participants à la recherche, la société, les écosystèmes, l'héritage culturel et l'environnement
- **Responsabilité** [partagée entre les acteurs de la recherche] **assumée pour les activités de recherche**, de l'idée à la publication, leur gestion et leur organisation, pour la formation, la supervision et le mentorat, et pour les implications plus générales de la recherche.

Code de conduite européen pour l'intégrité en recherche, 2018, page 4

[Source](#)



Et en + :




meilleur taux de citation pour les articles et données de la recherche librement accessibles

[Source](#)

# Formalisation des efforts dans un document

Plan de Gestion des Données (PGD) / *Data Management Plan (DMP)*



 ANR

Modèle  
de Plan de gestion  
des données (PGD)

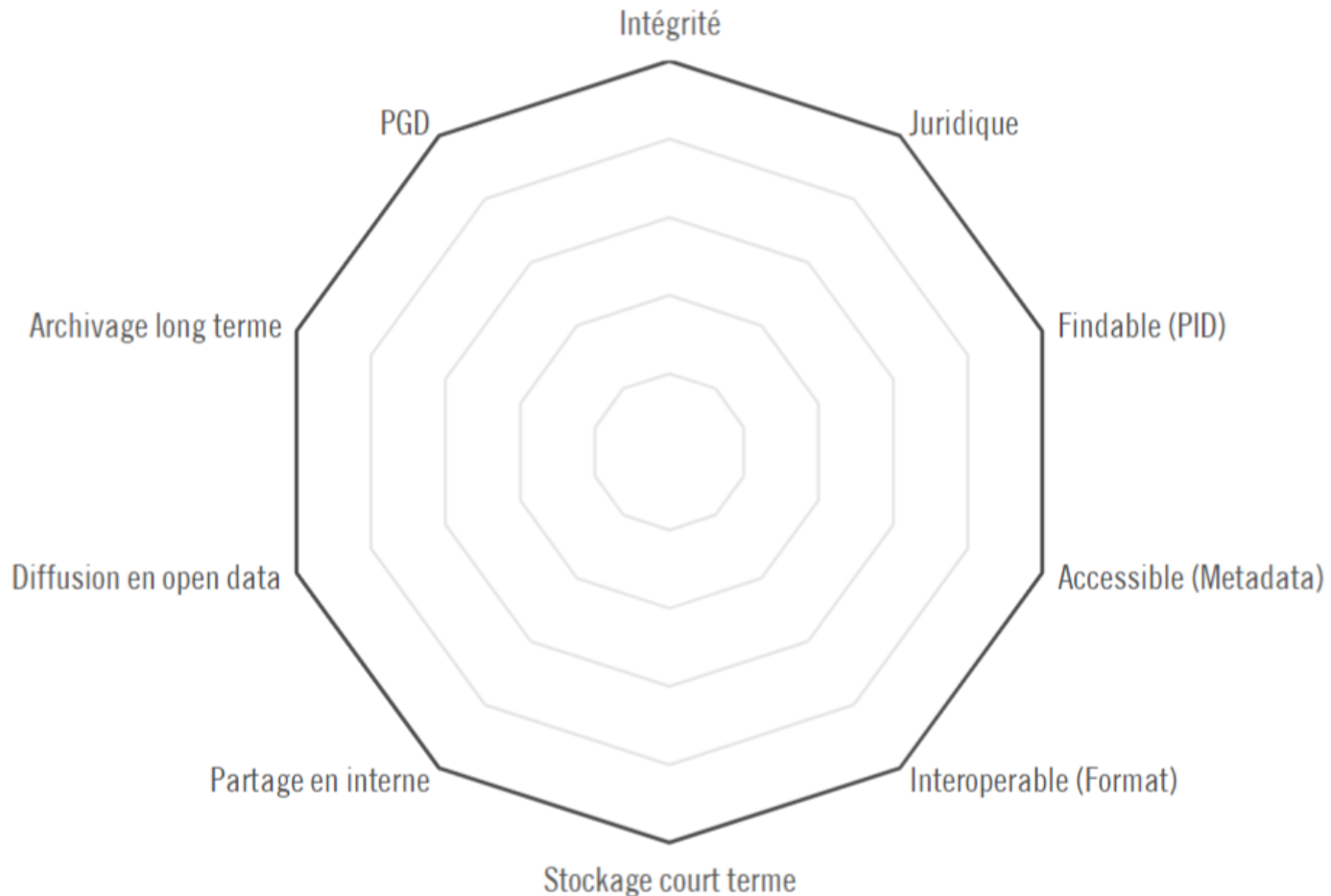
Agence nationale de la recherche - 50 avenue Daumesnil - 75012 Paris

[www.anr.fr](http://www.anr.fr)

[Source](#)

# Formalisation des efforts dans un document

## Radar Gestion des données de recherche



[Source](#)

# Ressources

## **Guides pratiques :**

<https://scienceouverte.couperin.org/>

<https://www.ouvrirlascience.fr/accueil/>

**Publications ouvertes :** <https://hal.science/>

**Données ouvertes :** <https://recherche.data.gouv.fr/fr>

**La science ouverte dans l'ED NSCo** <https://nsco.universite-lyon.fr/science-ouverte/>

## **Quelques outils :**

Cahier de laboratoire électronique <https://datacc.elab.one/login.php>

Plan de gestion de données <https://dmp.opidor.fr/>

Aide à la gestion & partage des données <https://dorum.fr/>

Exemple de feuille de route dans un labo : <https://www.crnl.fr/fr/page-base/science-ouverte>