

# Science Ouverte et PGD

## *Contexte et Opportunités*



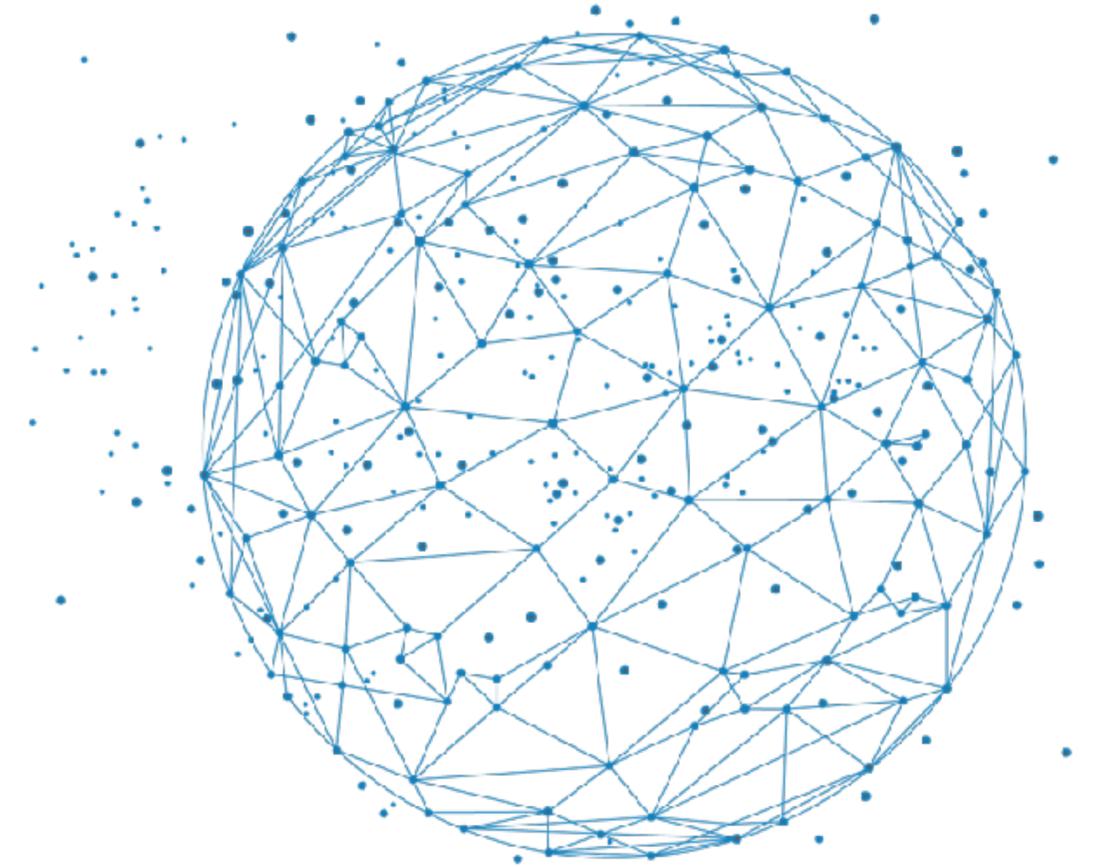
- Début d'un processus
  - Vous avez déjà des compétences
    - Spécialistes
      - Tout le monde va progresser grâce à nos échanges
  
- Interactif !
  - N'hésitez pas à partager vos expériences sur le sujet
  
- Dans une ambiance détendue et bienveillante

[frederic.de-lamotte@inrae.fr](mailto:frederic.de-lamotte@inrae.fr)

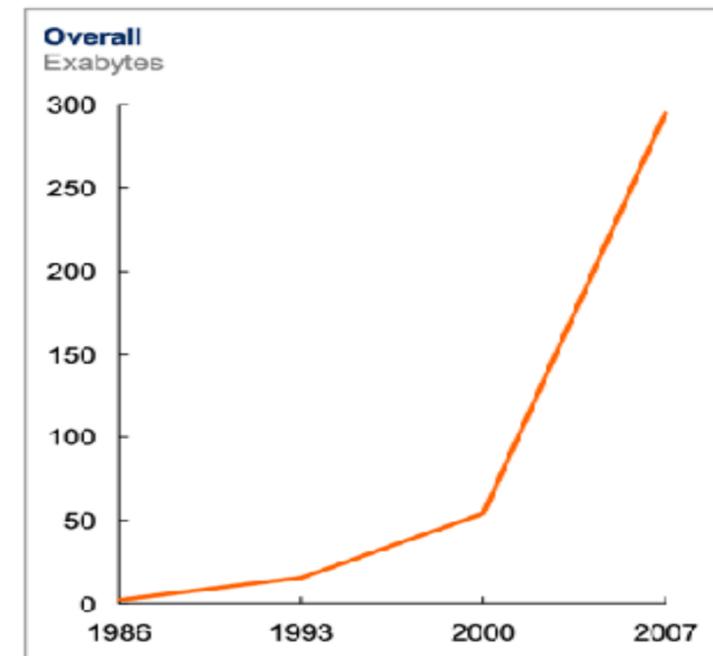
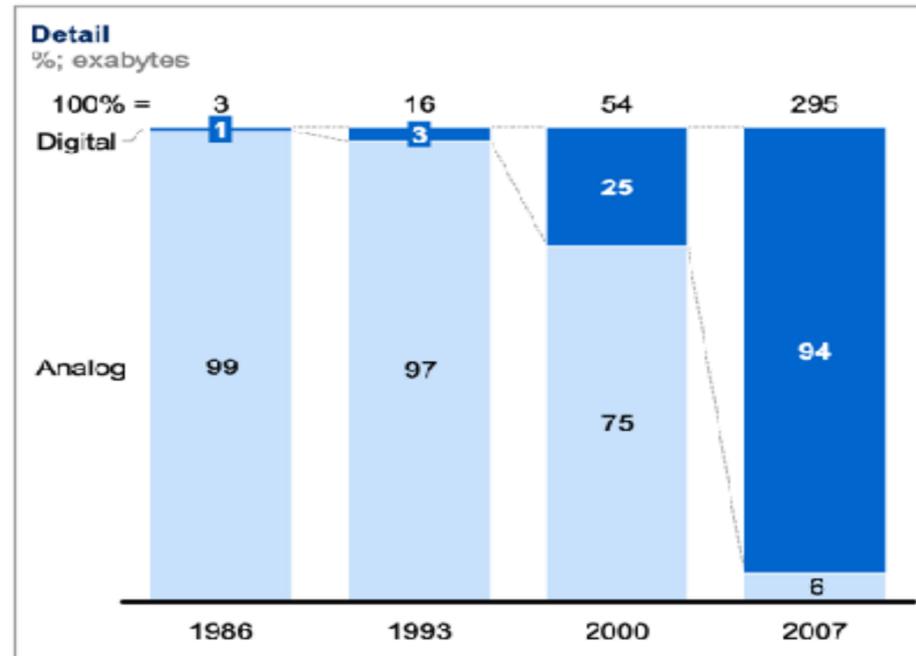
[oana.vigy@fpp.cnrs.fr](mailto:oana.vigy@fpp.cnrs.fr)

# FAIR pourquoi faire ?

Pourquoi maintenant ?



# La disruption numérique, une bascule brutale



The World's Technological Capacity to Store, Communicate, and Compute Information. Science. 2011;332.



Une musique  
4 Mo



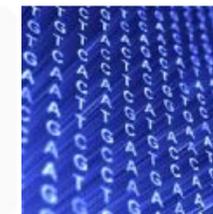
Une photo  
6 Mo



Un document  
50 Ko



Un film  
700 Mo



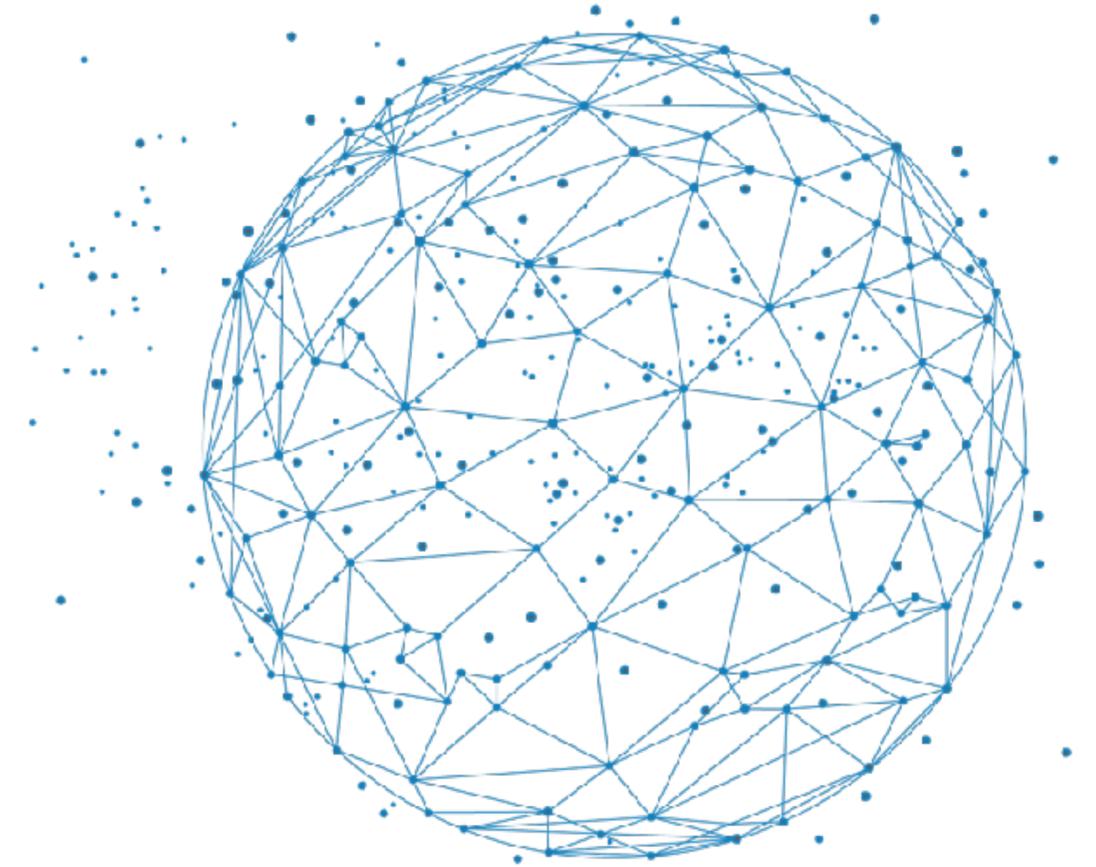
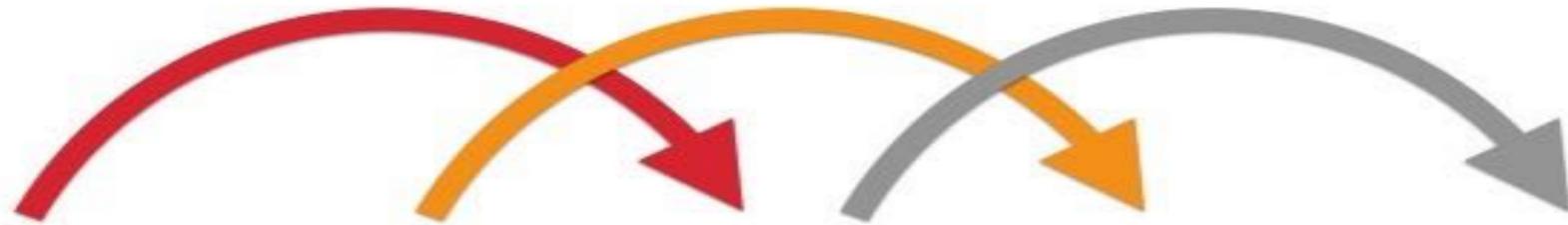
Un Génome  
100 Go

- La première compagnie de taxi n'en possède aucun (Uber)
- Le premier fournisseur de logement n'en possède pas (AirBnB)
- La première compagnie de téléphonie ne possède pas de standard (Skype)
- Le premier fournisseur d'info ne crée pas de contenu (Facebook)
- Le premier diffuseur de film ne possède pas de salle de cinéma (Netflix)



# Les technologies numériques ont transformé la biologie et santé

**Waves of Digital Disruption**



**Et pour vous, d'où est venue la disruption numérique ?**

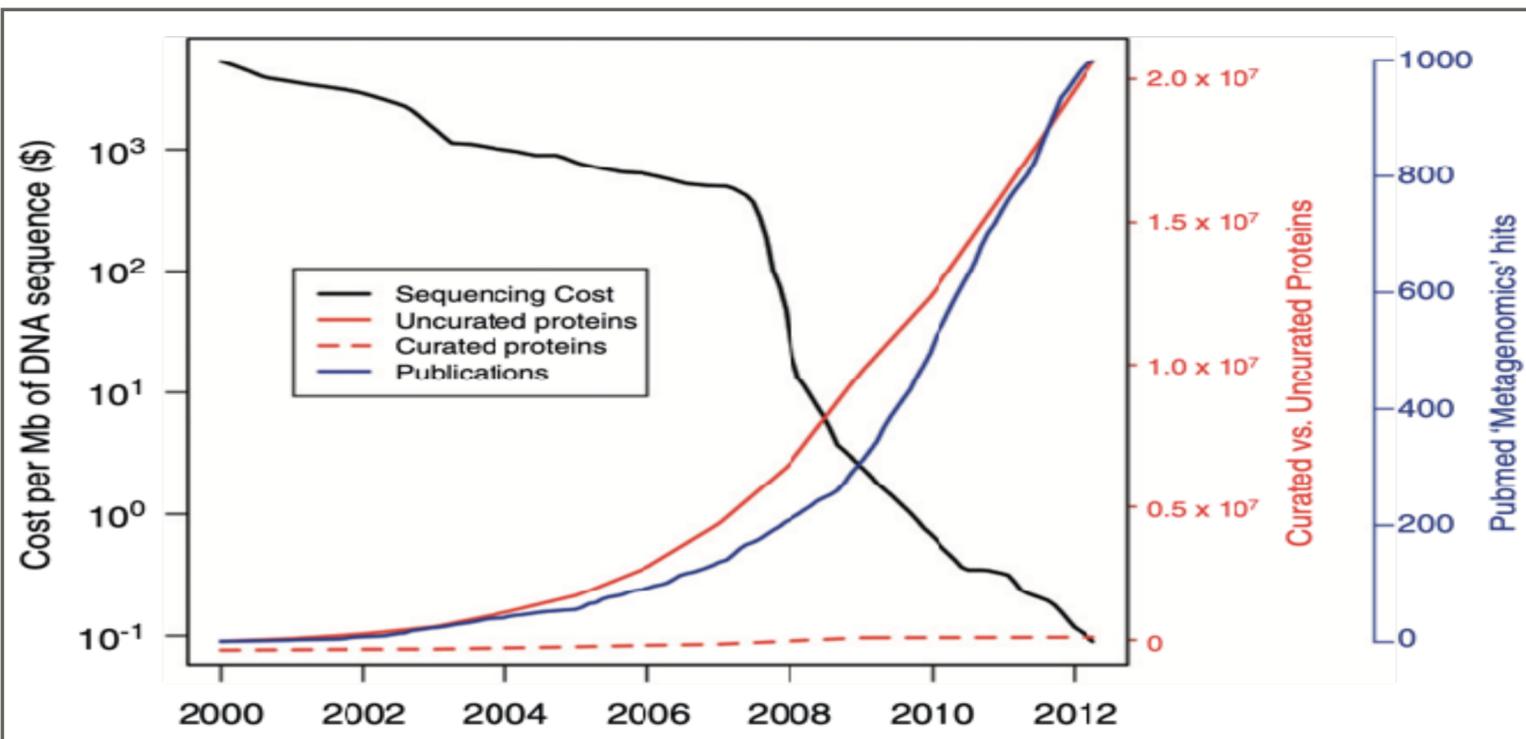
# Le déluge des données en Science



Les techniques à haut débit induisent une chute des coûts et une explosion de la production de données  
Génome humain :

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

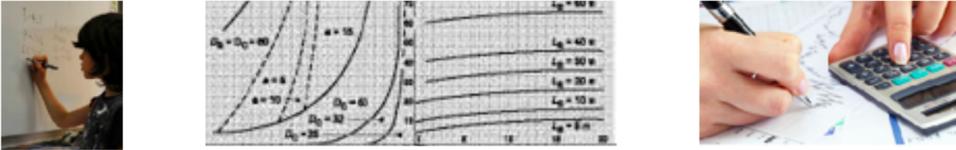
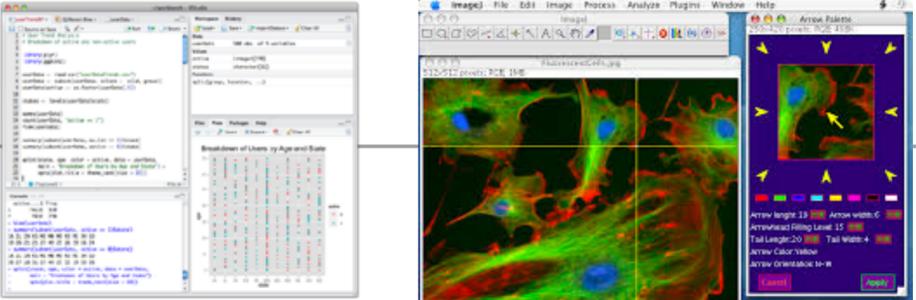
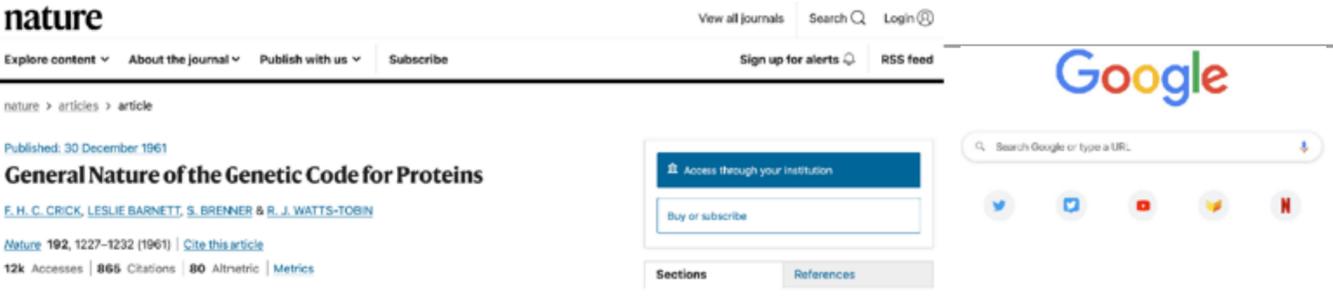
en 2015 = **quelques heures** et **1000 \$**

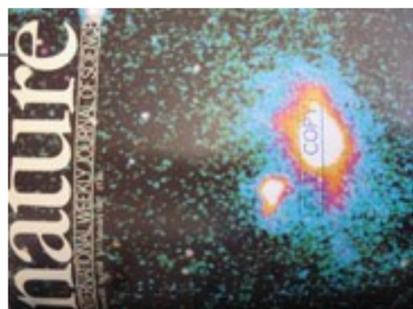


→ La quantité de données à **stocker** et **analyser** explose

→ Le *rendement* d'analyse chute

# Un changement d'échelle qui impacte profondément la science

Les étapes	<i>Avant</i>	<i>Maintenant</i>
Concevoir l'expérimentation	Une connaissance des dispositifs expérimentaux existants accessible à un ou quelques individus. Un volume d'observations attendues à taille humaine	Une matrice de technologies qui échappe à un expérimentateur individuel Possibilité d'utiliser et de générer une quantité massive de données expérimentales
Collecter des résultats		
Analyser des résultats		
Diffuser le savoir		





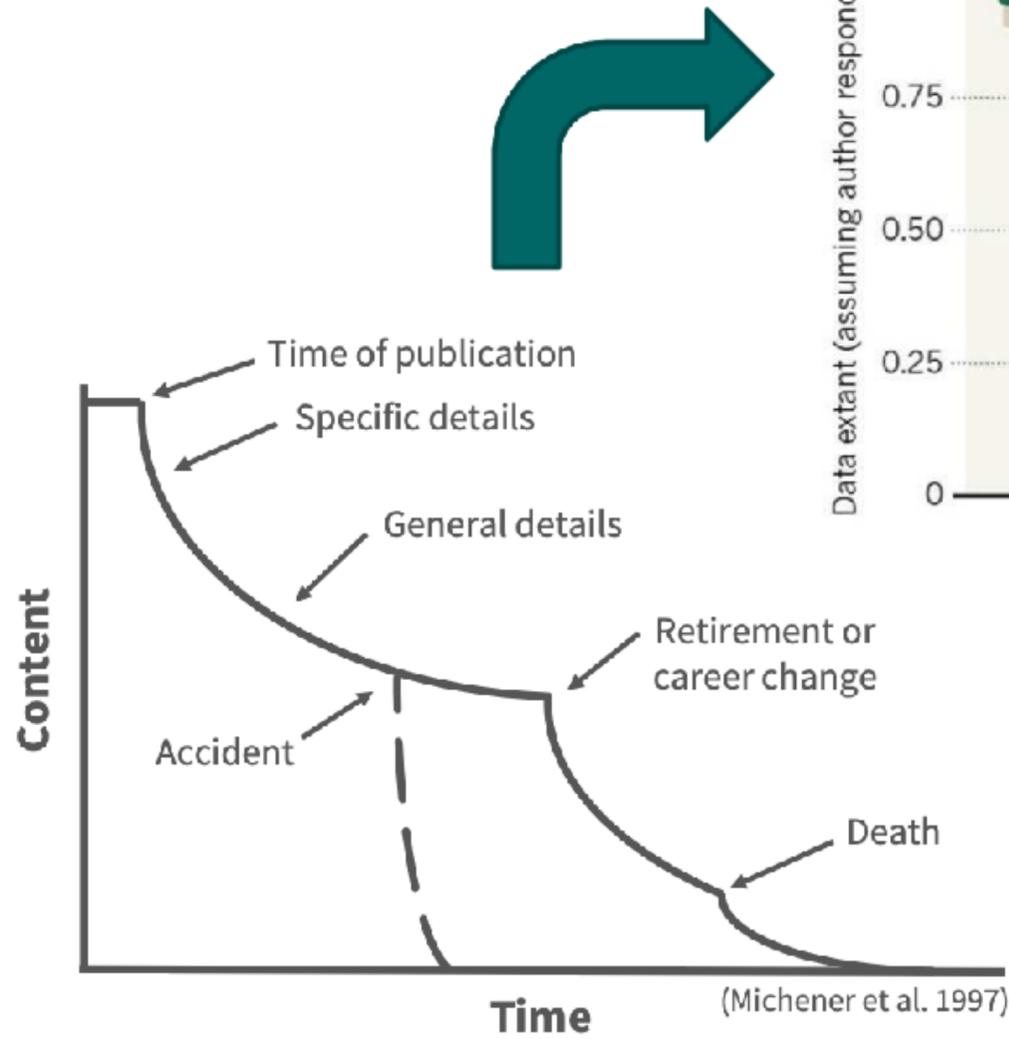
## THE HUMAN GENOME

### The Sequence of the Human Genome

J. Craig Venter,<sup>1\*</sup> Mark D. Adams,<sup>1</sup> Eugene W. Myers,<sup>1</sup> Peter W. Li,<sup>1</sup> Richard J. Mural,<sup>1</sup> Granger G. Sutton,<sup>1</sup> Hamilton O. Smith,<sup>1</sup> Mark Yandell,<sup>1</sup> Cheryl A. Evans,<sup>1</sup> Robert A. Holt,<sup>1</sup> Jeannine D. Gocayne,<sup>1</sup> Peter Amanatides,<sup>1</sup> Richard M. Ballew,<sup>1</sup> Daniel H. Huson,<sup>1</sup> Jennifer Russo Wortman,<sup>1</sup> Qing Zhang,<sup>1</sup> Chinnappa D. Kodira,<sup>1</sup> Xiangqun H. Zheng,<sup>1</sup> Lin Chen,<sup>1</sup> Marian Skupski,<sup>1</sup> Gangadharan Subramanian,<sup>1</sup> Paul D. Thomas,<sup>1</sup> Jinghui Zhang,<sup>1</sup> George L. Gabor Miklos,<sup>2</sup> Catherine Nelson,<sup>3</sup> Samuel Broder,<sup>1</sup> Andrew G. Clark,<sup>4</sup> Joe Nadeau,<sup>5</sup> Victor A. McKusick,<sup>6</sup> Norton Zinder,<sup>7</sup> Arnold J. Levine,<sup>7</sup> Richard J. Roberts,<sup>8</sup> Mel Simon,<sup>9</sup> Carolyn Stayman,<sup>10</sup> Michael Hunkapiller,<sup>11</sup> Randall Bolanos,<sup>1</sup> Arthur Delcher,<sup>1</sup> Ian Dew,<sup>1</sup> Daniel Fasulo,<sup>1</sup> Michael Flanigan,<sup>1</sup> Liliana Florea,<sup>1</sup> Aaron Halpern,<sup>1</sup> Sridhar Hannenhalli,<sup>1</sup> Saul Kravitz,<sup>1</sup> Samuel Levy,<sup>1</sup> Clark Mobarry,<sup>1</sup> Knut Reinert,<sup>1</sup> Karin Remington,<sup>1</sup> Jane Abu-Threideh,<sup>1</sup> Ellen Beasley,<sup>1</sup> Kendra Biddick,<sup>1</sup> Vivien Bonazzi,<sup>1</sup> Rhonda Brandon,<sup>1</sup> Michele Cargill,<sup>1</sup> Ishwar Chandramouliswaran,<sup>1</sup> Rosane Charlab,<sup>1</sup> Kabir Chaturvedi,<sup>1</sup> Zuoming Deng,<sup>1</sup> Valentina Di Francesco,<sup>1</sup> Patrick Dunn,<sup>1</sup> Karen Elbeck,<sup>1</sup> Carlos Evangelista,<sup>1</sup> Andrei E. Gabrielian,<sup>1</sup> Weiniu Gan,<sup>1</sup> Wangmao Ge,<sup>1</sup> Fangcheng Gong,<sup>1</sup> Zhiping Gu,<sup>1</sup> Ping Guan,<sup>1</sup> Thomas J. Heiman,<sup>1</sup> Maureen E. Higgins,<sup>1</sup> Rui-Ru Ji,<sup>1</sup> Zhaoxi Ke,<sup>1</sup> Karen A. Ketchum,<sup>1</sup> Zhongwu Lai,<sup>1</sup> Yiding Lei,<sup>1</sup> Zhenya Li,<sup>1</sup> Jiayin Li,<sup>1</sup> Yong Liang,<sup>1</sup> Xiaoying Lin,<sup>1</sup> Fu Lu,<sup>1</sup> Gennady V. Merkulov,<sup>1</sup> Natalia Milshina,<sup>1</sup> Helen M. Moore,<sup>1</sup> Ashwinikumar K Naik,<sup>1</sup> Vaibhav A. Narayan,<sup>1</sup> Beena Neelam,<sup>1</sup> Deborah Nusskern,<sup>1</sup> Douglas B. Rusch,<sup>1</sup> Steven Salzberg,<sup>12</sup> Wei Shao,<sup>1</sup> Bixiong Shue,<sup>1</sup> Jingtao Sun,<sup>1</sup> Zhen Yuan Wang,<sup>1</sup> Aihui Wang,<sup>1</sup> Xin Wang,<sup>1</sup> Jian Wang,<sup>1</sup> Ming-Hui Wei,<sup>1</sup> Ron Wides,<sup>13</sup> Chunlin Xiao,<sup>1</sup> Chunhua Yan,<sup>1</sup> Alison Yao,<sup>1</sup> Jane Ye,<sup>1</sup> Ming Zhan,<sup>1</sup> Weiqing Zhang,<sup>1</sup> Hongyu Zhang,<sup>1</sup> Qi Zhao,<sup>1</sup> Liansheng Zheng,<sup>1</sup> Fei Zhong,<sup>1</sup> Wenyan Zhong,<sup>1</sup> Shiaoqing C. Zhu,<sup>1</sup> Shaying Zhao,<sup>12</sup> Dennis Gilbert,<sup>1</sup> Suzanna Baumhueter,<sup>1</sup> Gene Spier,<sup>1</sup> Christine Carter,<sup>1</sup> Anibal Cravchik,<sup>1</sup> Trevor Woodage,<sup>1</sup> Feroze Ali,<sup>1</sup> Huijin An,<sup>1</sup> Aderonke Awe,<sup>1</sup> Danita Baldwin,<sup>1</sup> Holly Baden,<sup>1</sup> Mary Barnstead,<sup>1</sup> Ian Barrow,<sup>1</sup> Karen Beeson,<sup>1</sup> Dana Busam,<sup>1</sup> Amy Carver,<sup>1</sup> Angela Center,<sup>1</sup> Ming Lai Cheng,<sup>1</sup> Liz Curry,<sup>1</sup> Steve Danaher,<sup>1</sup> Lionel Davenport,<sup>1</sup> Raymond Desilets,<sup>1</sup> Susanne Dietz,<sup>1</sup> Kristina Dodson,<sup>1</sup> Lisa Doup,<sup>1</sup> Steven Ferreira,<sup>1</sup> Neha Garg,<sup>1</sup> Andres Gluecksmann,<sup>1</sup> Brit Hart,<sup>1</sup> Jason Haynes,<sup>1</sup> Charles Haynes,<sup>1</sup> Cheryl Heiner,<sup>1</sup> Suzanne Hladun,<sup>1</sup> Damon Hostin,<sup>1</sup> Jarrett Houck,<sup>1</sup> Timothy Howland,<sup>1</sup> Chinyere Ibegwam,<sup>1</sup> Jeffery Johnson,<sup>1</sup> Francis Kalush,<sup>1</sup> Lesley Kline,<sup>1</sup> Shashi Koduru,<sup>1</sup> Amy Love,<sup>1</sup> Felecia Mann,<sup>1</sup> David May,<sup>1</sup> Steven McCawley,<sup>1</sup> Tina McIntosh,<sup>1</sup> Ivy McMullen,<sup>1</sup> Mee Moy,<sup>1</sup> Linda Moy,<sup>1</sup> Brian Murphy,<sup>1</sup> Keith Nelson,<sup>1</sup> Cynthia Pfannkoch,<sup>1</sup> Eric Pratts,<sup>1</sup> Vinita Puri,<sup>1</sup> Hina Qureshi,<sup>1</sup> Matthew Reardon,<sup>1</sup> Robert Rodriguez,<sup>1</sup> Yu-Hui Rogers,<sup>1</sup> Deanna Romblad,<sup>1</sup> Bob Ruhfel,<sup>1</sup> Richard Scott,<sup>1</sup> Cynthia Sitter,<sup>1</sup> Michelle Smallwood,<sup>1</sup> Erin Stewart,<sup>1</sup> Renee Strong,<sup>1</sup> Ellen Suh,<sup>1</sup> Reginald Thomas,<sup>1</sup> Ni Ni Tint,<sup>1</sup> Sukyee Tse,<sup>1</sup> Claire Vech,<sup>1</sup> Gary Wang,<sup>1</sup> Jeremy Wetter,<sup>1</sup> Sherita Williams,<sup>1</sup> Monica Williams,<sup>1</sup> Sandra Windsor,<sup>1</sup> Emily Winn-Deen,<sup>1</sup> Keriellen Wolfe,<sup>1</sup> Jayshree Zaveri,<sup>1</sup> Karena Zaveri,<sup>1</sup> Josep F. Abril,<sup>14</sup> Roderic Guigó,<sup>14</sup> Michael J. Campbell,<sup>1</sup> Kimmen V. Sjolander,<sup>1</sup> Brian Karlak,<sup>1</sup> Anish Kejariwal,<sup>1</sup> Huaiyu Mi,<sup>1</sup> Betty Lazareva,<sup>1</sup> Thomas Hatton,<sup>1</sup> Apurva Narechania,<sup>1</sup> Karen Diemer,<sup>1</sup> Anushya Muruganujan,<sup>1</sup> Nan Guo,<sup>1</sup> Shinji Sato,<sup>1</sup> Vineet Bafna,<sup>1</sup> Sorin Istrail,<sup>1</sup> Ross Lippert,<sup>1</sup> Russell Schwartz,<sup>1</sup> Brian Walenz,<sup>1</sup> Shibu Yooseph,<sup>1</sup> David Allen,<sup>1</sup> Anand Basu,<sup>1</sup> James Baxendale,<sup>1</sup> Louis Blick,<sup>1</sup> Marcelo Caminha,<sup>1</sup> John Carnes-Stine,<sup>1</sup> Parris Caulk,<sup>1</sup> Yen-Hui Chiang,<sup>1</sup> My Coyne,<sup>1</sup> Carl Dahlke,<sup>1</sup> Anne Deslattes Mays,<sup>1</sup> Maria Dombroski,<sup>1</sup> Michael Donnelly,<sup>1</sup> Dale Ely,<sup>1</sup> Shiva Esparham,<sup>1</sup> Carl Foster,<sup>1</sup> Harold Gire,<sup>1</sup> Stephen Glanowski,<sup>1</sup> Kenneth Glasser,<sup>1</sup> Anna Glodek,<sup>1</sup> Mark Gorokhov,<sup>1</sup> Ken Graham,<sup>1</sup> Barry Gropman,<sup>1</sup> Michael Harris,<sup>1</sup> Jeremy Heil,<sup>1</sup> Scott Henderson,<sup>1</sup> Jeffrey Hoover,<sup>1</sup> Donald Jennings,<sup>1</sup> Catherine Jordan,<sup>1</sup> James Jordan,<sup>1</sup> John Kasha,<sup>1</sup> Leonid Kagan,<sup>1</sup> Cheryl Kraft,<sup>1</sup> Alexander Levitsky,<sup>1</sup> Mark Lewis,<sup>1</sup> Xiangjun Liu,<sup>1</sup> John Lopez,<sup>1</sup> Daniel Ma,<sup>1</sup> William Majoros,<sup>1</sup> Joe McDaniel,<sup>1</sup> Sean Murphy,<sup>1</sup> Matthew Newman,<sup>1</sup> Trung Nguyen,<sup>1</sup> Ngoc Nguyen,<sup>1</sup> Marc Nodell,<sup>1</sup> Sue Pan,<sup>1</sup> Jim Peck,<sup>1</sup> Marshall Peterson,<sup>1</sup> William Rowe,<sup>1</sup> Robert Sanders,<sup>1</sup> John Scott,<sup>1</sup> Michael Simpson,<sup>1</sup> Thomas Smith,<sup>1</sup> Arlan Sprague,<sup>1</sup> Timothy Stockwell,<sup>1</sup> Russell Turner,<sup>1</sup> Eli Venter,<sup>1</sup> Mei Wang,<sup>1</sup> Meiyuan Wen,<sup>1</sup> David Wu,<sup>1</sup> Mitchell Wu,<sup>1</sup> Ashley Xia,<sup>1</sup> Ali Zandieh,<sup>1</sup> Xiaohong Zhu<sup>1</sup>

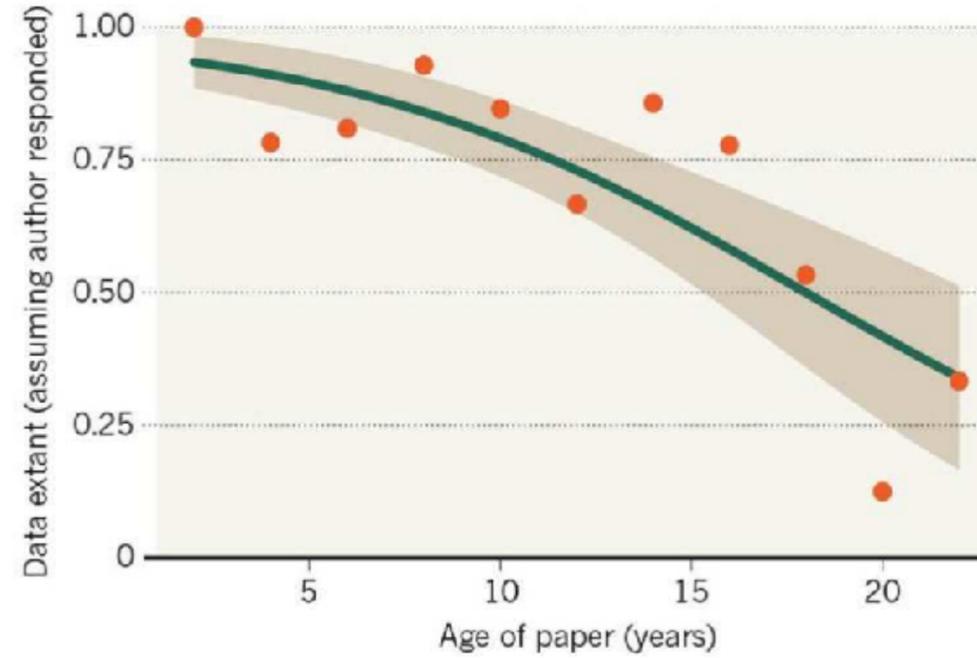


## Data Entropy



### MISSING DATA

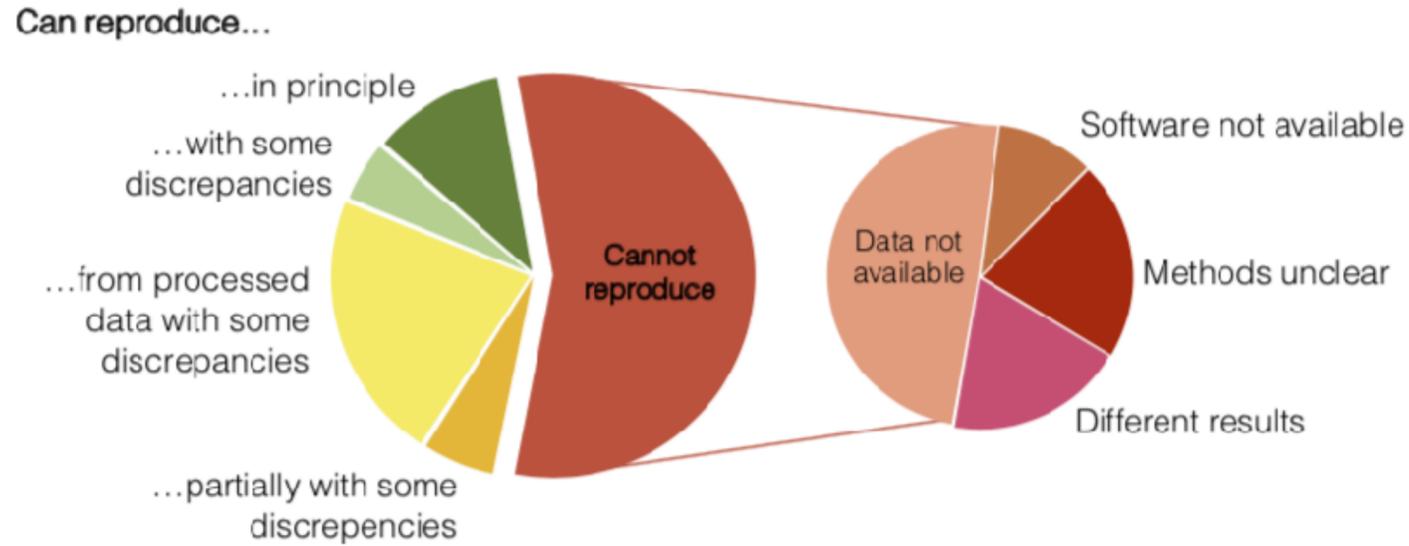
As research articles age, the odds of their raw data being extant drop dramatically.



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

# La crise de reproductibilité

Replication of data analyses in 18 articles on microarray-based gene expression profiling published in Nature Genetics in 2005–2006:



Repeatability of published microarray gene expression analyses. Nat Genet. **2009**;41: 149–155. doi:10.1038/ng.295

Step	Reference																		
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
Mapping	█	█	█	█	█	█	█	█	█	█	█	█	█	█	█	█	█	█	█
Realignment		█								█	█					█	█		█
Recalibration				█						█	█						█	█	
Initial variant detection	█	█	█	█	█	█	█	█	█	█	█			█		█	█	█	█

Next-generation sequencing data interpretation: enhancing reproducibility and accessibility. Nat Rev Genet. **2012**;13: 667–672.

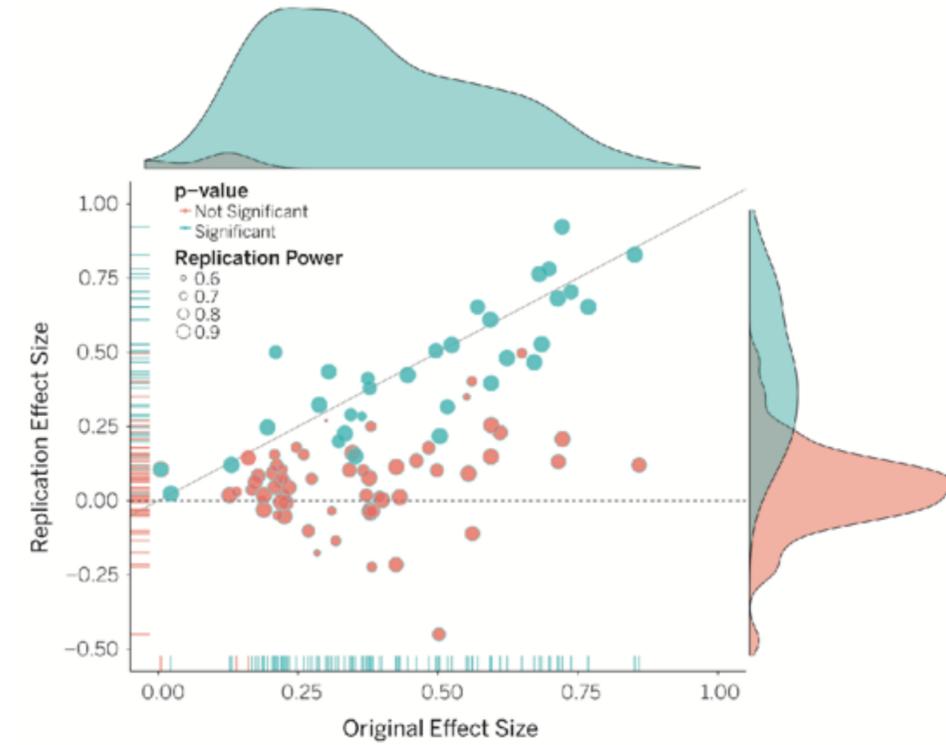
RESEARCH ARTICLE

Estimating the reproducibility of psychological science

Open Science Collaboration<sup>1,2</sup>  
 • See all authors and affiliations  
 Science 28 Aug 2015;  
 Vol. 349, Issue 6251, aac4716  
 DOI: 10.1126/science.aac4716

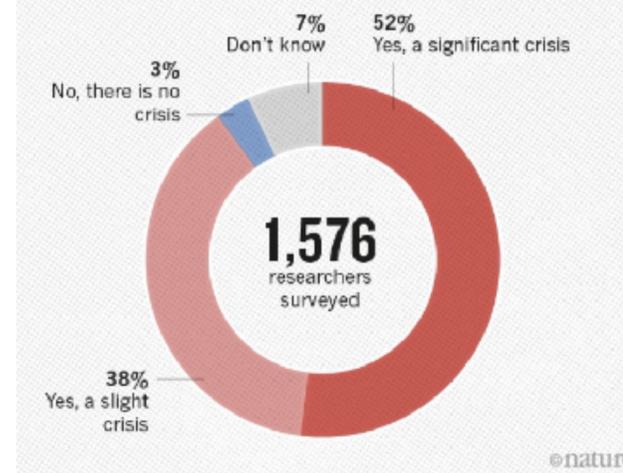
The *Reproducibility project* set out to replicate 100 experiments published in high-impact psychology journals.

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



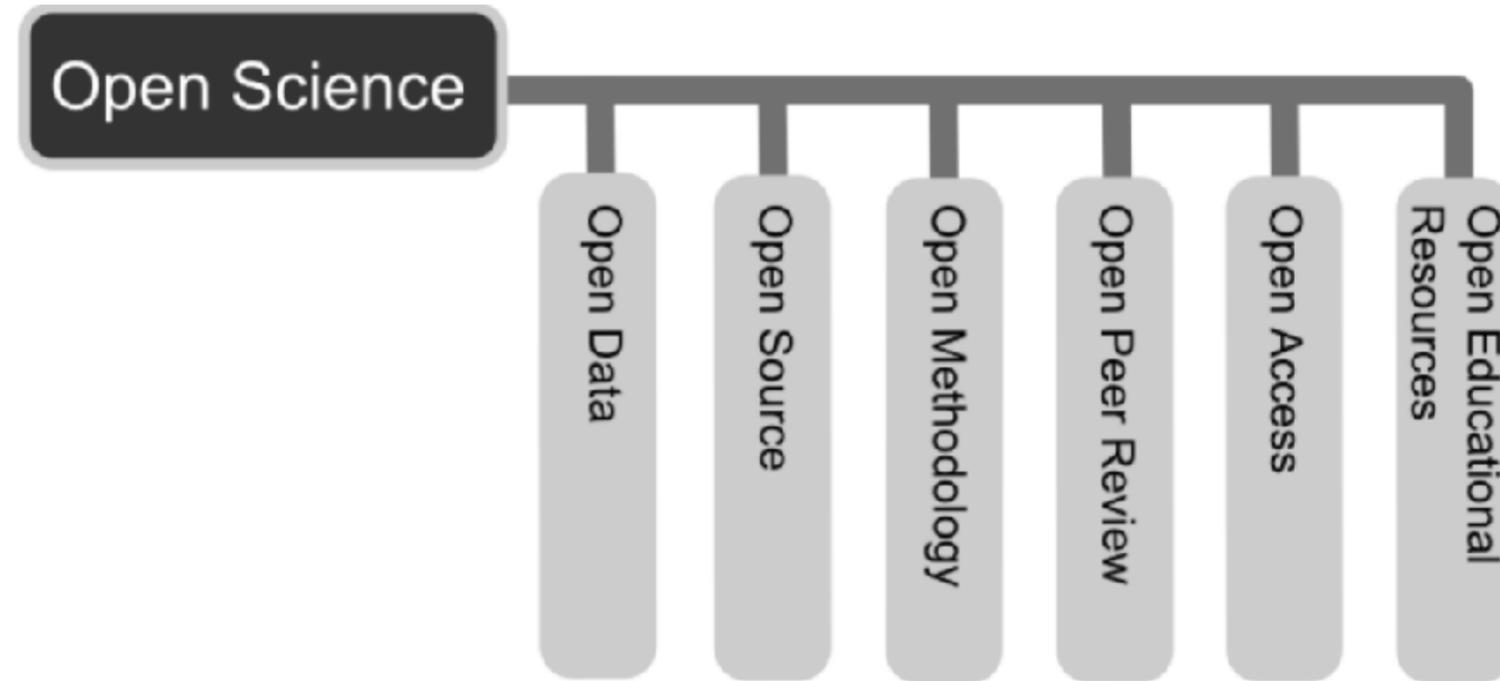
Estimating the reproducibility of psychological science. Science. **2015**;349: aac4716. doi:10.1126/science.aac4716

IS THERE A REPRODUCIBILITY CRISIS?



Is there a reproducibility crisis in science? Nature. **2016**. doi:10.1038/d41586-019-00067-3

- Une solution ?





- Accept that the computational component is becoming an integral component of biomedical research
- Always provide access to primary data
- Record versions of all auxiliary data sets used during the analysis
- Note the exact versions of software used
- Record all parameters even if defaults are used
- Provide all custom scripts
- Do not reinvent the wheel

**Anton Nekrutenko and James Taylor**  
Nat Rev Genet. **2012**;13: 667–672.

#### Box 3 | Guidelines for reproducibility

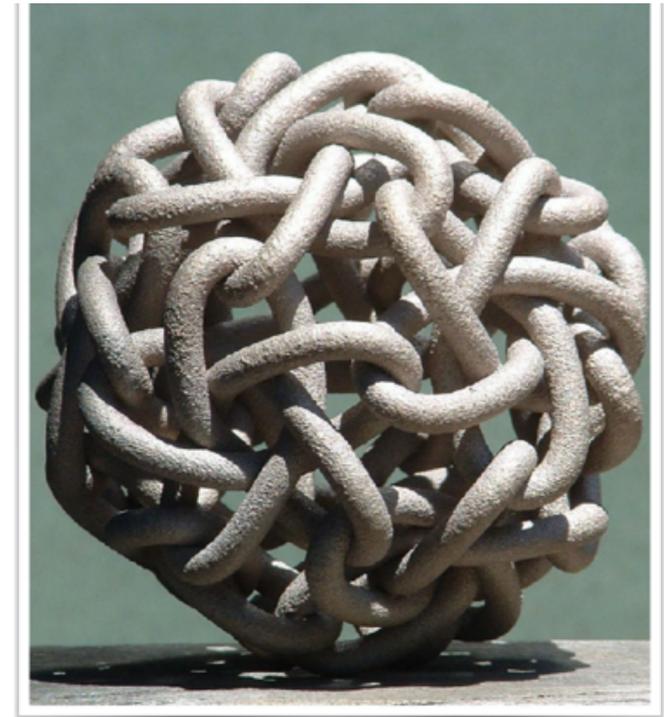
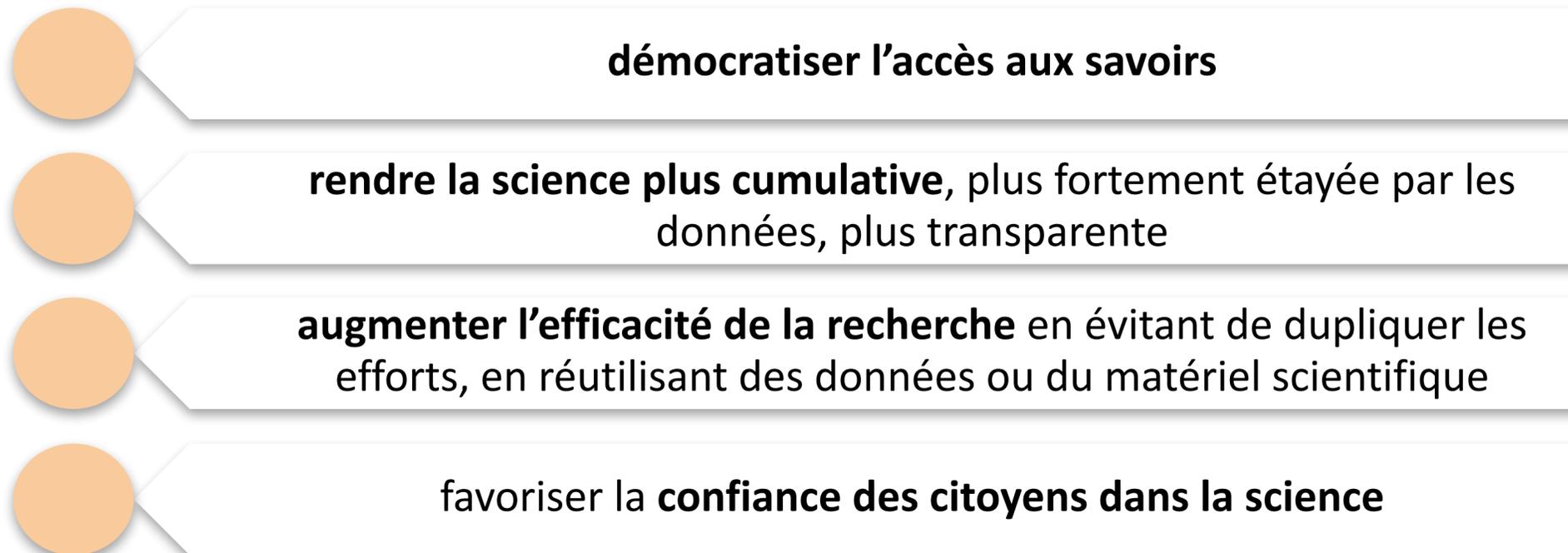
- **Accept that the computational component is becoming an integral component of biomedical research.** As the life sciences are becoming increasingly data-driven, there will be no escape from computation and data handling. Familiarize yourself with best practices of scientific computing using existing educational resources, such as the Software Carpentry project<sup>41</sup>. Implementing good computational practices in your group will automatically take care of many of the points listed below.
- **Always provide access to primary data.** It is obvious that without access to the original data sets, any claims made in a publication cannot be verified. In situations in which the data cannot be made public (for example, clinical data sets under Institutional Review Board protection), they should be deposited in controlled access repositories (such as dbGaP<sup>42</sup>), where they can be retrieved by authorized users. One potential issue with this point is the fact that there is currently a debate on what constitutes primary data. Storing images generated by some next-generation sequencing (NGS) machines on a large scale has long been unfeasible. Public sequencing archives, such as those at the National Center for Biotechnology Information (NCBI) and the European Bioinformatics Institute (EBI), are still accepting sequencing reads as submissions and should be used. Going forward, other formats, such as aligned data in BAM format, are likely to be used (as is already done by the 1000 Genomes Project).
- **Record versions of all auxiliary data sets used during the analysis.** For example, in most NGS analyses, such as variant discovery detailed here, sequencing reads are compared against a reference genome. It is crucial to record which reference genome was used because, just as software has versions and cars have model years, genomes have build identifiers. For example, the latest human genome build distributed by the UCSC Genome Browser is called hg19 (it is derived from the GRC37 build prepared by the Genome Reference Consortium) and has the highest number of functional annotations (7,330 annotation types) and should be the preferred version to be used. Note that the latest version may not always be the best choice. The latest mouse genome build (mm10) has only a fraction of annotations (258 tracks) compared with its predecessor (mm9, which has 2,096 tracks). Thus, it would be easier to interpret results of an NGS experiment mapped to the mm9 build even though mm10 has an additional 48 megabases of actual sequence.
- **Note the exact versions of software used.** Different versions of the same software often produce different results, and important bug fixes may have implications to results produced with a particular version of a tool.
- **Record all parameters, even if defaults are used.** Although the reason to record all parameters requires no explanation, we emphasize the importance of explaining default settings for reproducibility. A clause ‘software was used with default settings’ is found in many publications. However, the meaning of default settings often changes between versions of software and can be quite difficult to track down when a substantial amount of time has passed since publication. Thus, record what the default settings actually are.
- **Provide all custom scripts.** With the complexity of NGS analysis, it is often unavoidable to create simple scripts that carry out such straightforward tasks as, for example, changing data formats. Such scripts must be made accessible as any other part of the analysis.
- **Do not reinvent the wheel.** It pays to reuse existing software. Integrative frameworks and associated application stores already house hundreds of tools (for example, as of May 2012, Galaxy ToolShed contains ~1,700 tools). It is likely that a script for a particular problem has been already written. Ask around through existing resources such as SEQanswers<sup>43</sup> and BioStar<sup>44</sup>.

# Les racines du mouvement “#OpenScience”

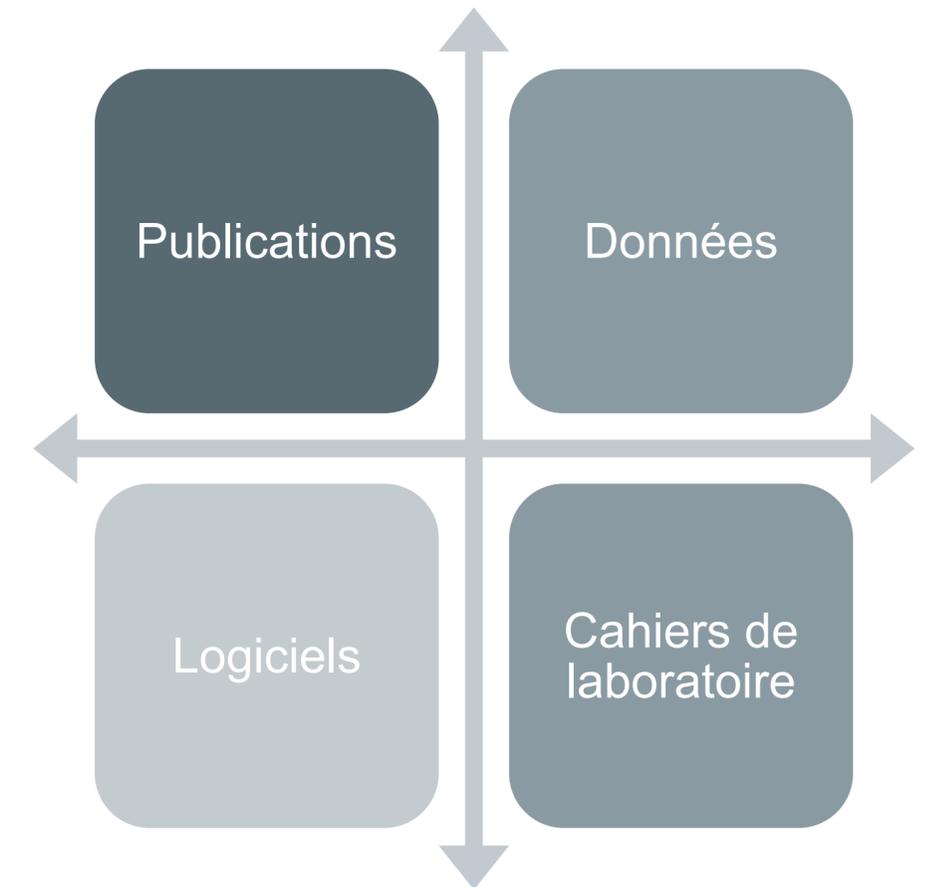


- La crise de reproductibilité
- La crise de déontologie scientifique (P-Hacking)  
P-hacking in clinical trials and how incentives shape the distribution of results across phases.  
Proc Natl Acad Sci U S A. 2020;117: 13386–13392.)
- La crise politique (éthique/démocratique)  
(Publish or Perish & indicateurs / ego metrics, producteurs de savoir biomédical privés, marché de l'édition scientifique)
- La crise de la “bioinformatique”

## → Améliorer la science et l'innovation



- Rendre les **résultats de la recherche scientifique** accessibles à tous (un des aspects de la Science Ouverte)
- En permettant une diffusion et une réutilisation sans entrave des données et résultats de la recherche

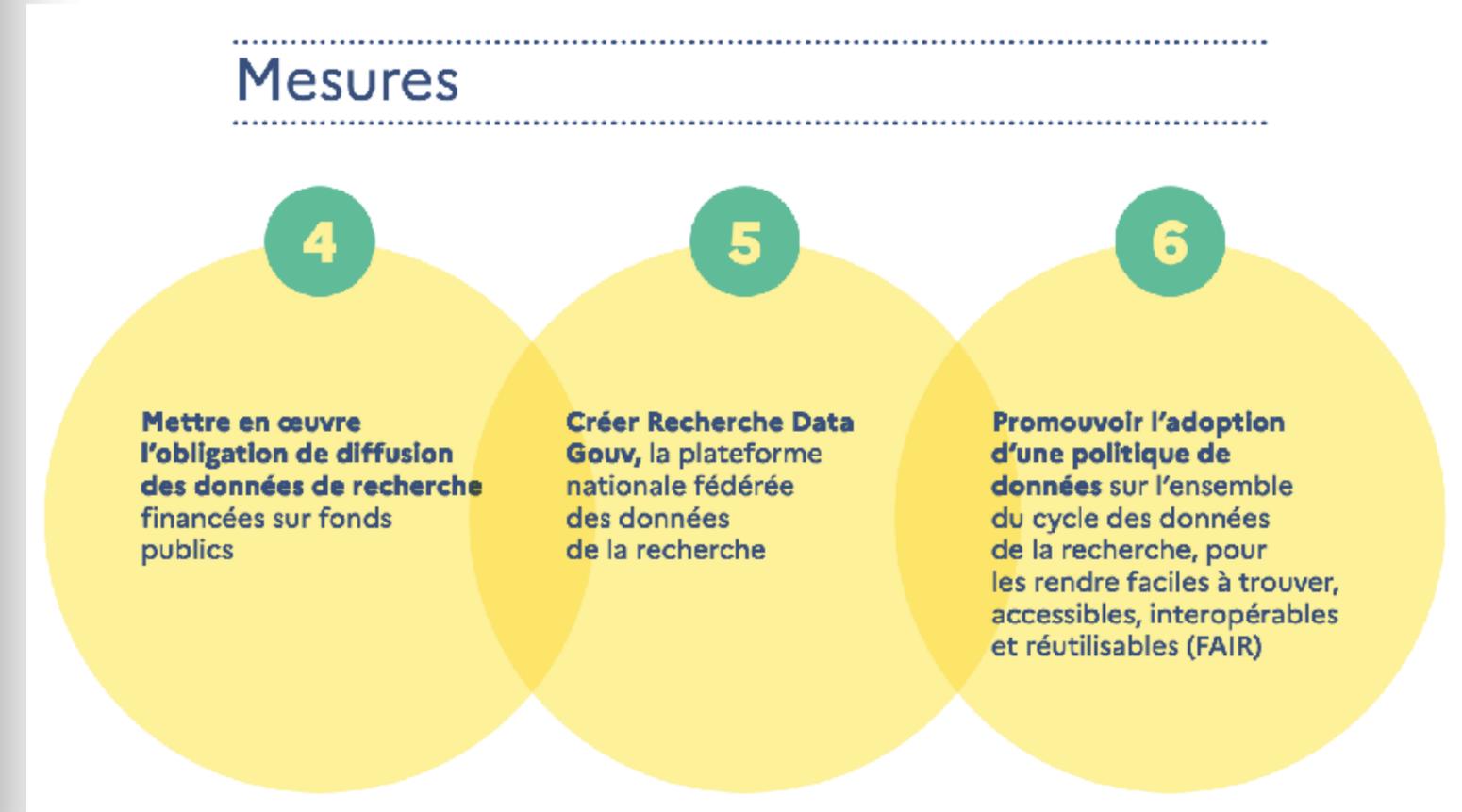




2018



2021





## Citez les types de données de la recherche que vous connaissez

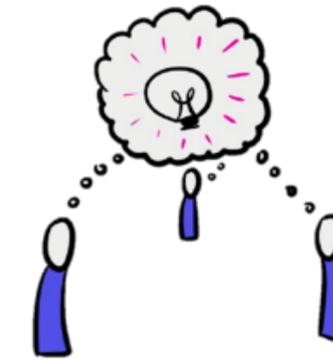
En utilisant le tableau à Post-It  
[http://scrumblr.ca/Data\\_Type](http://scrumblr.ca/Data_Type)



- Les données de recherche sont les **preuves** qui sous-tendent la réponse à la question de recherche et peuvent être utilisées pour **valider** les **résultats**, quelle que soit leur forme (i.e. imprimée, numérique ou physique).
- Il peut s'agir de **renseignements quantitatifs** ou d'énoncés **qualitatifs** recueillis par les chercheurs dans le cadre de leurs travaux par **expérimentation, observation, modélisation, entrevue** ou autres méthodes, ou de renseignements tirés de preuves existantes.
- Les données peuvent être **brutes** ou **primaires** (par exemple, directement issues de mesures ou de collectes) ou **dérivées** de données primaires par analyse ou interprétation (e.g. nettoyées ou extraites d'un ensemble de données plus vaste), ou encore dérivées de sources existantes dont les droits peuvent être détenus par d'autres.



[https://youtu.be/66oNv\\_DJuPc](https://youtu.be/66oNv_DJuPc)



@picto-dico

Notez les points marquants (bon ou mauvais) en gestion des données

**Data Sharing and Management Snafu in 3 Short Acts**  
by Karen Hanson, Alisa Surkis & Karen Yacobucci  
NYU Health Sciences Libraries  
August 3, 2012 (Last Update: December 12, 2012)





## Citez 5 conditions

En utilisant le tableau à Post-It  
<http://scrumblr.ca/solutions>



**F**indable **A**ccessible **I**nteroperable **R**eusable





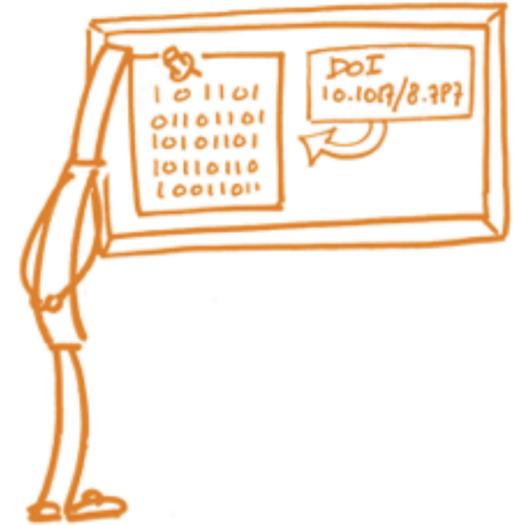
Faciliter la découverte des données (et de leurs métadonnées)  
tant pour les humains que pour les machines



- Les données ont un **PID** (Persistent Identifier, identifiant unique et pérenne)
- Les données sont décrites par des **métadonnées**
- Ces métadonnées incluent le PID des données qu'elles décrivent
- Les données sont déposées dans un **entrepôt de données**



Permettre l'accès aux données et leur téléchargement, ce qui peut inclure l'authentification et l'autorisation.

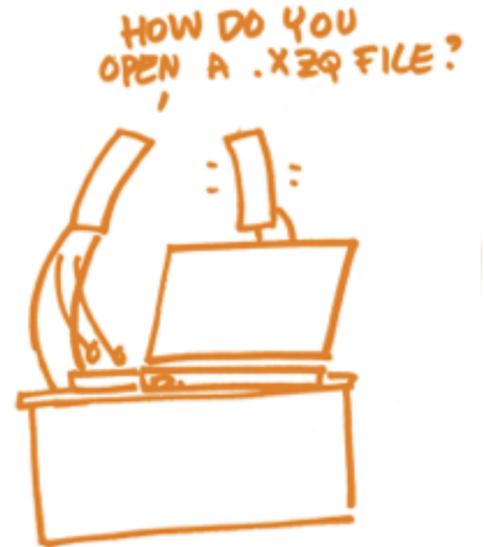


- Les données sont accessibles à travers un **protocole de communication standard**
- Ce protocole est **libre et ouvert**
- Ce protocole permet un accès par **authentification** si besoin
- Les **métadonnées restent accessibles** même si les données ne le sont pas (disparues ou inaccessibles)



Permettre l'exploitation et l'intégration des données quel que soit l'environnement informatique utilisé

- Les données sont **décrites avec un vocabulaire contrôlé**
- Le vocabulaire utilisé **respecte les principes FAIR**
- Les **métadonnées sont contextualisées** avec des liens vers d'autres données

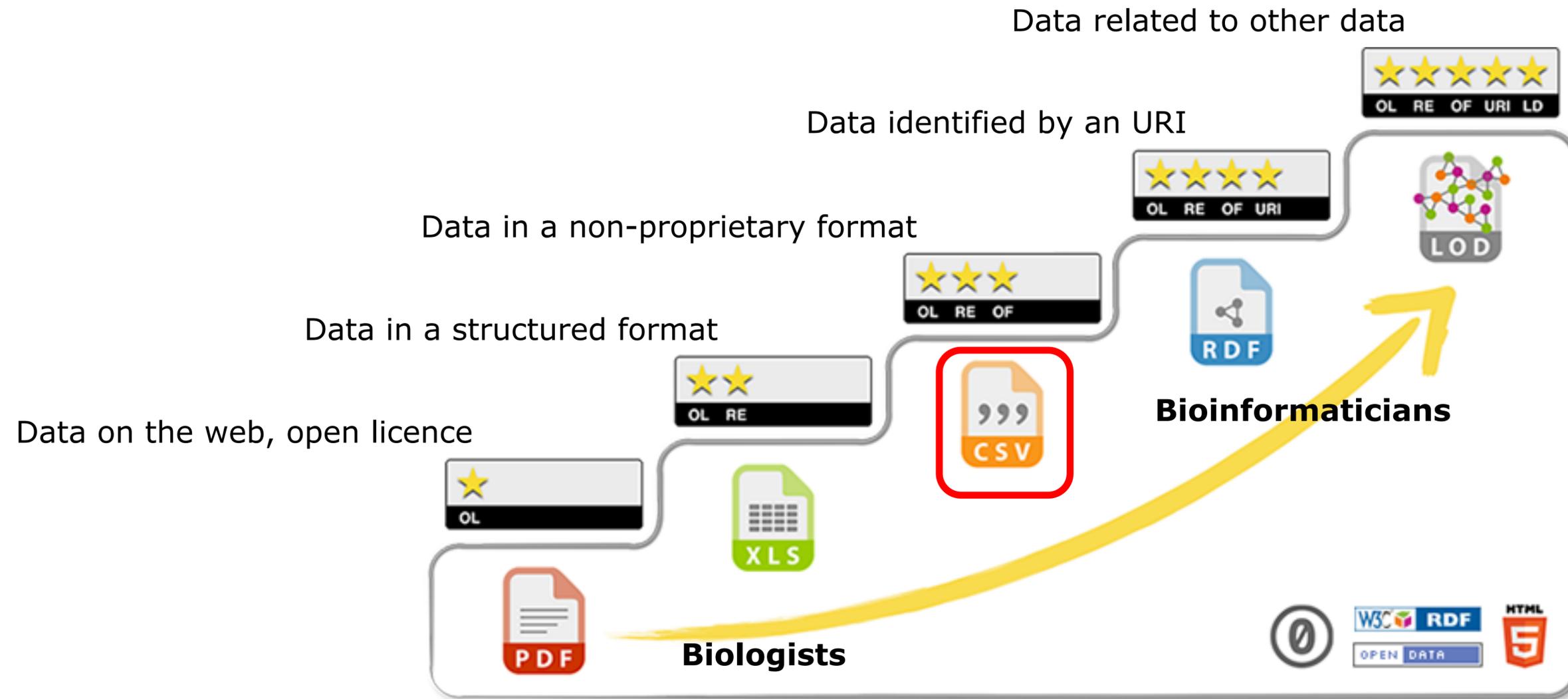




Permettre la réutilisation des données pour de futures recherches



- Les métadonnées contiennent toutes les informations qui peuvent être utiles (**pluralité d'attributs**)
- Une **licence de réutilisation** est attribuée aux données
- La description des données indique leur **provenance**
- Le partage des données suit les **standards de la communauté scientifique**



La progression vers FAIR et l'Open Data nécessite une coopération multidisciplinaire :

- Biologistes
- Bioinformaticiens
- Spécialistes des ontologies/sémantiques

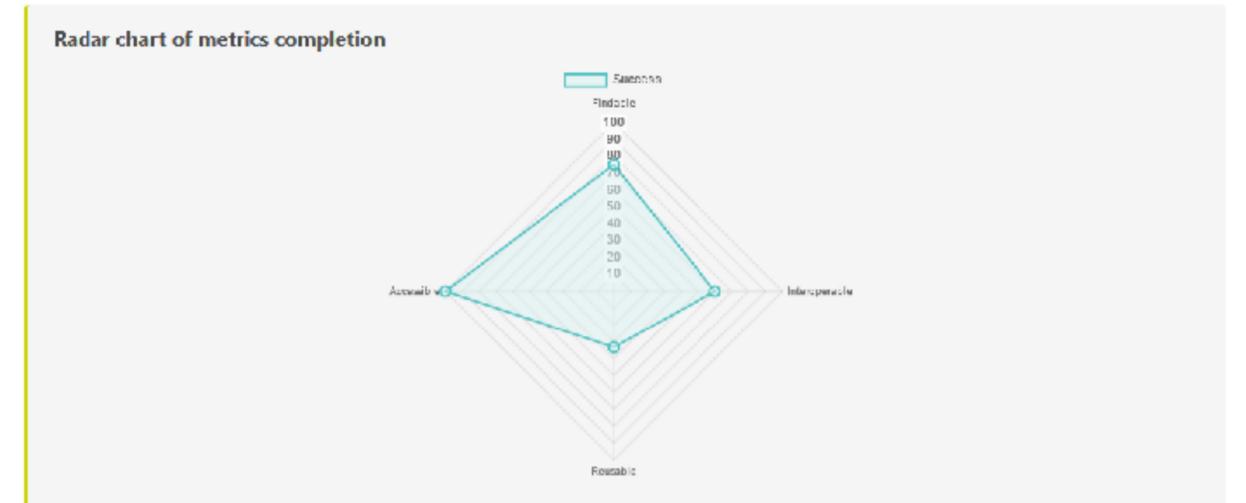
# Evaluation de la FAIRness des données

- Grille SHARC (SHARing Reward & Credit)
  - RDA (Research Data Alliance)
  - [https://zenodo.org/record/2551500#.X4hC\\_-2re70](https://zenodo.org/record/2551500#.X4hC_-2re70)
- FAIR self assessment tool
  - ARDC (Australian Research Data Commons)
  - <https://ardc.edu.au/resources/aboutdata/fair-data/fair-self-assessment-tool/>
- FAIR checker
  - IFB
  - [https://fair-checker.france-bioinformatique.fr/base\\_metrics](https://fair-checker.france-bioinformatique.fr/base_metrics)

Enter resource identifier (URL/DOI)

The URL/DOI is valid The input contains the following DOIs that you can also test: [doi:10.15454/P27LDX](#)

[Examples](#) [Dataset Database](#) [Workflow](#) [Publication Database](#) [Dataset](#) [Tool](#)

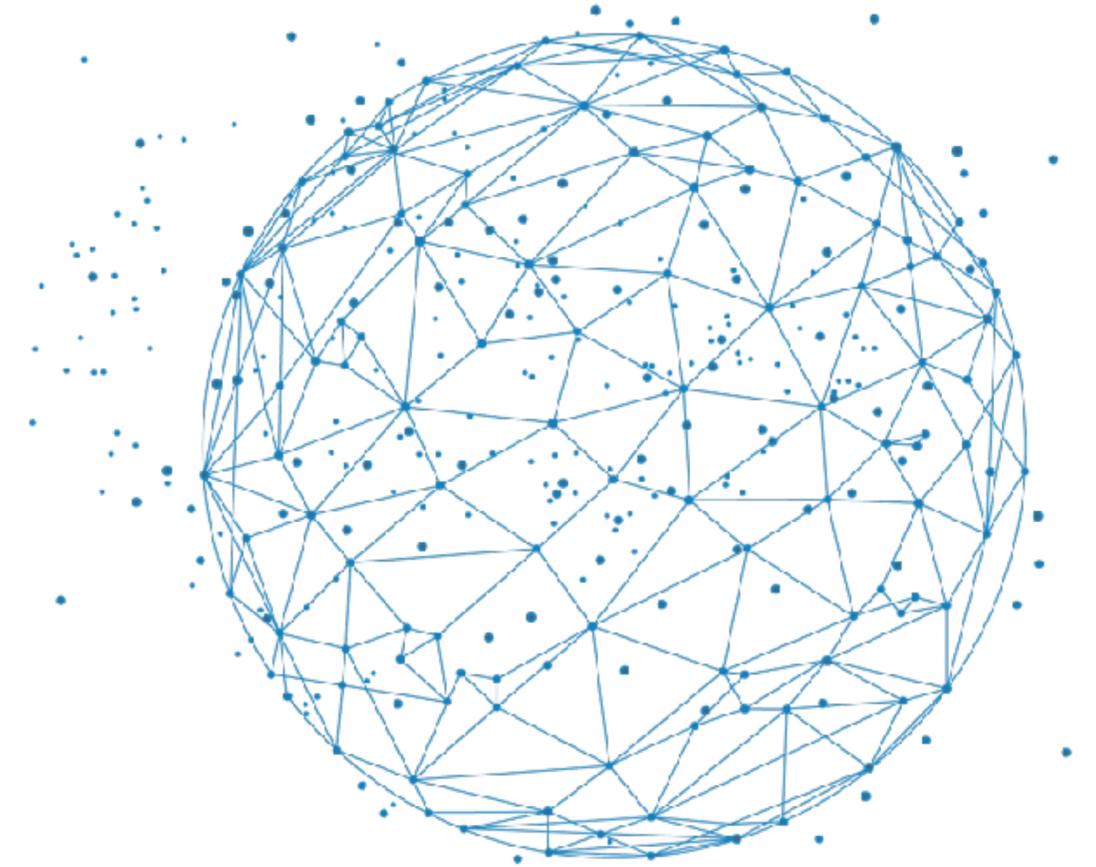


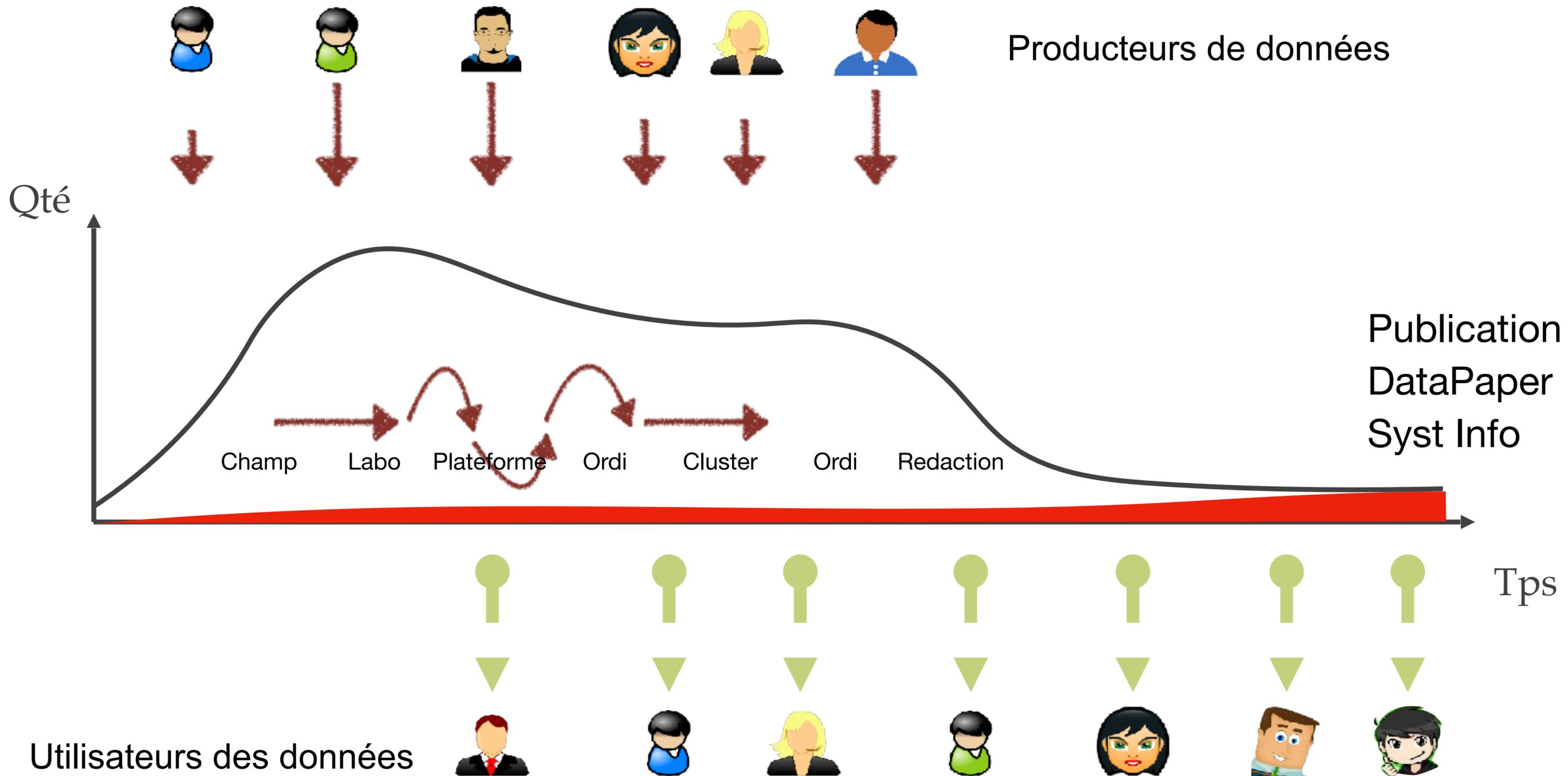
List of metrics with details and results

Principle	Name	Description	Comment	Recommendation	Score	Result	Test	Details
F1A	Unique IDs				2	Success	<a href="#">Check</a>	
F1B	Persistent IDs				0	Failure	<a href="#">Check</a>	
F2A	Structured metadata				2	Success	<a href="#">Check</a>	
F2B	Shared vocabularies for metadata				1	Success	<a href="#">Check</a>	
A1.1	Open resolution protocol				2	Success	<a href="#">Check</a>	
I1A	Any structured information				2	Success	<a href="#">Check</a>	
I1B	Ontological and machine-resolvable formats				2	Success	<a href="#">Check</a>	
I2A	Human-readable vocabularies				0	Failure	<a href="#">Check</a>	
I2B	Machine-readable vocabularies				2	Success	<a href="#">Check</a>	
R1	External links				0	Failure	<a href="#">Check</a>	
R1.1	Metadata includes license				0	Failure	<a href="#">Check</a>	
R1.2	Metadata includes provenance				0	Failure	<a href="#">Check</a>	
R1.3	Community standards				1	Success	<a href="#">Check</a>	

For additional tips and recommendations on how to improve your resource, we recommend you to use the FAIR Cookbook: <https://fairplus.github.io/the-fair-cookbook/content/home.html>

# Cycle de vie des données



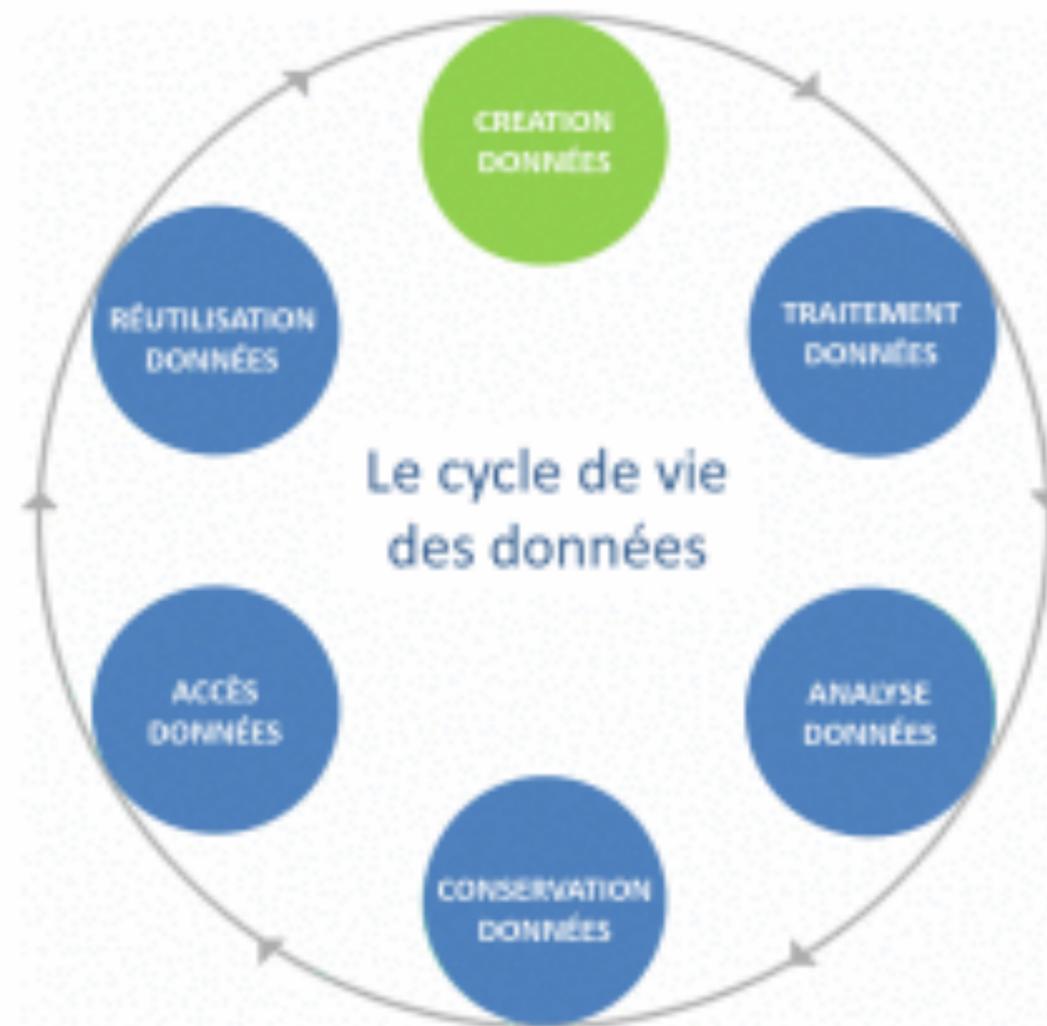




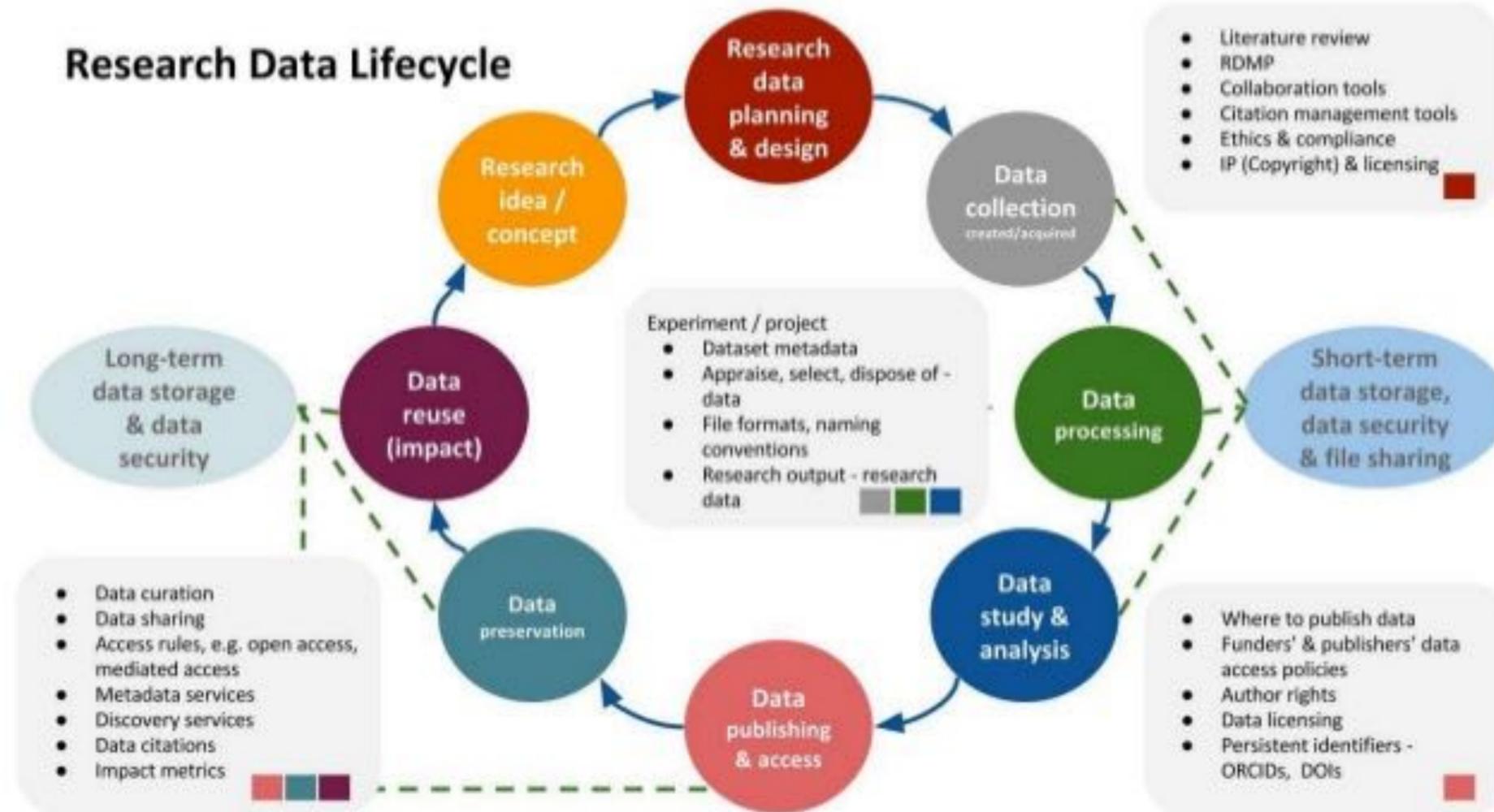
Le modèle de UK Data Archive définit les six étapes suivantes :

- **Création ou collecte** des données (creating data) ;
- **Traitement** des données (processing data) ;
- **Analyse** des données (analysing data) ;
- **Conservation** des données (preserving data) ;
- **Accès** aux données (giving access to data / data discovery) ;
- **Réutilisation** des données (reusing data).

[Une introduction à la gestion et au partage des données de la recherche - Le cycle de vie des données](#)



[Une introduction à la gestion et au partage des données de la recherche - Le cycle de vie des données](#)



# Plan de Gestion de Données



# Avez vous déjà pris l'avion ?





## PGD de projet

**Document** qui définit comment seront gérées les données d'un projet **pendant et après le projet**

- La difficulté : penser à traiter toutes les étapes du cycle de vie des données
- L'avantage : les modèles vous aident à penser à tout, en vous posant une série de questions

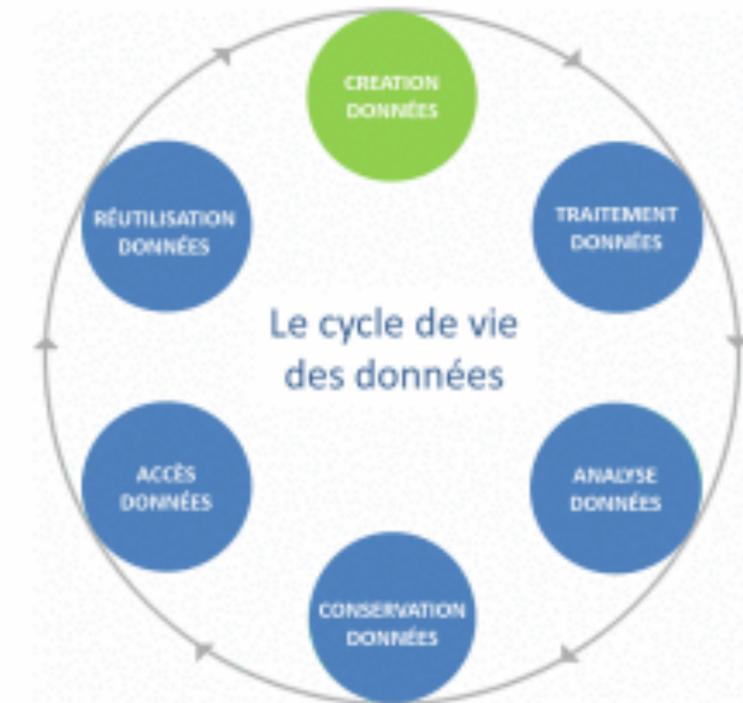
## QUAND

Document évolutif = au moins 3 versions :

- Première version au début du projet de recherche
- Mises à jour régulières au cours du projet : versions intermédiaires
- Version finale à la fin du projet

## QUI

**Rédaction par l'équipe de recherche coordinatrice du projet**





## Modèle de PGD = une liste de questions à remplir pour rédiger un PGD

### Modèles existants :

- Modèles des financeurs : ANR, Commission Européenne...
- Modèle Science Europe - <https://doi.org/10.5281/zenodo.4915862>
- Modèles institutionnels, des centres de calcul/stockage... (INRAe, CEA, Institut Pasteur, IN2P3, Unviversités, Grandes écoles...)

Peu / pas de modèles spécifiques à des types de données



- **Plan** : on planifie (donc on anticipe)
- **Gestion** : on gère, on fait fructifier (on commence déjà par ne plus perdre)
- **Données** : à bien définir au préalable

→ Se poser les bonnes questions, le plus tôt possible

→ Penser à toutes les démarches à effectuer à chaque étape du cycle de vie des données

- Assurer la reproductibilité des expériences
- Décrire comment les données sont obtenues
- Faciliter la réutilisation des données
  - Garantir la compréhension des données
- Respecter le droit et les personnes
  - Clarifier le cadre juridique et éthique
- Éviter les pertes de données
  - Assurer un stockage adapté
- Clarifier les droits de réutilisation
  - Spécifier les modalités de partage
- Établir le rôle de chacun
  - Définir les responsabilités

## **Pour l'équipe de recherche**

Se référer au PGD pour :

- Retrouver les données
- Comprendre les données
- Savoir où sont conservées les données ...

## **Pour la communauté scientifique**

Publier le PGD pour indiquer :

- Quelles données existent
- Où elles sont conservées
- Qui peut y accéder, sous quelles conditions...



## 1. Description des données et collecte ou réutilisation de données existantes

tout développer | tout réduire

**1.1 Description générale du produit de recherche**



**1.2 Est-ce que des données existantes seront réutilisées ?**



**1.3 Comment seront produites/collectées les nouvelles données ?**



→ Objectif : Assurer la reproductibilité des expériences

- Décrire les données



## 2. Documentation et qualité des données

tout développer | tout réduire

**2.1 Quelles métadonnées et quelle documentation (par exemple mode d'organisation des données) accompagneront les données ?**



**2.2 Quelles seront les méthodes utilisées pour assurer la qualité scientifique des données ?**



## 4. Traitement et analyse des données

tout développer | tout réduire

**4.1 Comment et avec quels moyens seront traitées les données ?**

- Objectif : Faciliter la réutilisation des données
- Garantir la compréhension des données



## 3. Exigences légales et éthiques, code de conduite

tout développer | tout réduire

**3.1 Quelles seront les mesures appliquées pour assurer la protection des données à caractère personnel ?**



**3.2 Quelles sont les contraintes juridiques (sensibilité des données autres qu'à caractère personnel, confidentialité, ...) à prendre en compte pour le partage et le stockage des données ?**



**3.3 Quels sont les aspects éthiques à prendre en compte lors de la collecte des données ?**



- Objectif : Respecter le droit et les personnes
- Clarifier le cadre juridique et éthique



## 5. Stockage et sauvegarde des données pendant le processus de recherche

tout développer | tout réduire

**5.1 Comment les données seront-elles stockées et sauvegardées tout au long du projet ?**



- Objectif : Éviter les pertes de données
- Assurer un stockage adapté



## 6. Partage des données et conservation à long terme

tout développer | tout réduire

**6.1 Comment les données seront-elles partagées ?**



**6.2 Comment les données seront-elles conservées à long terme ?**



- Objectif : Clarifier les droits de réutilisation
- Spécifier les modalités de partage

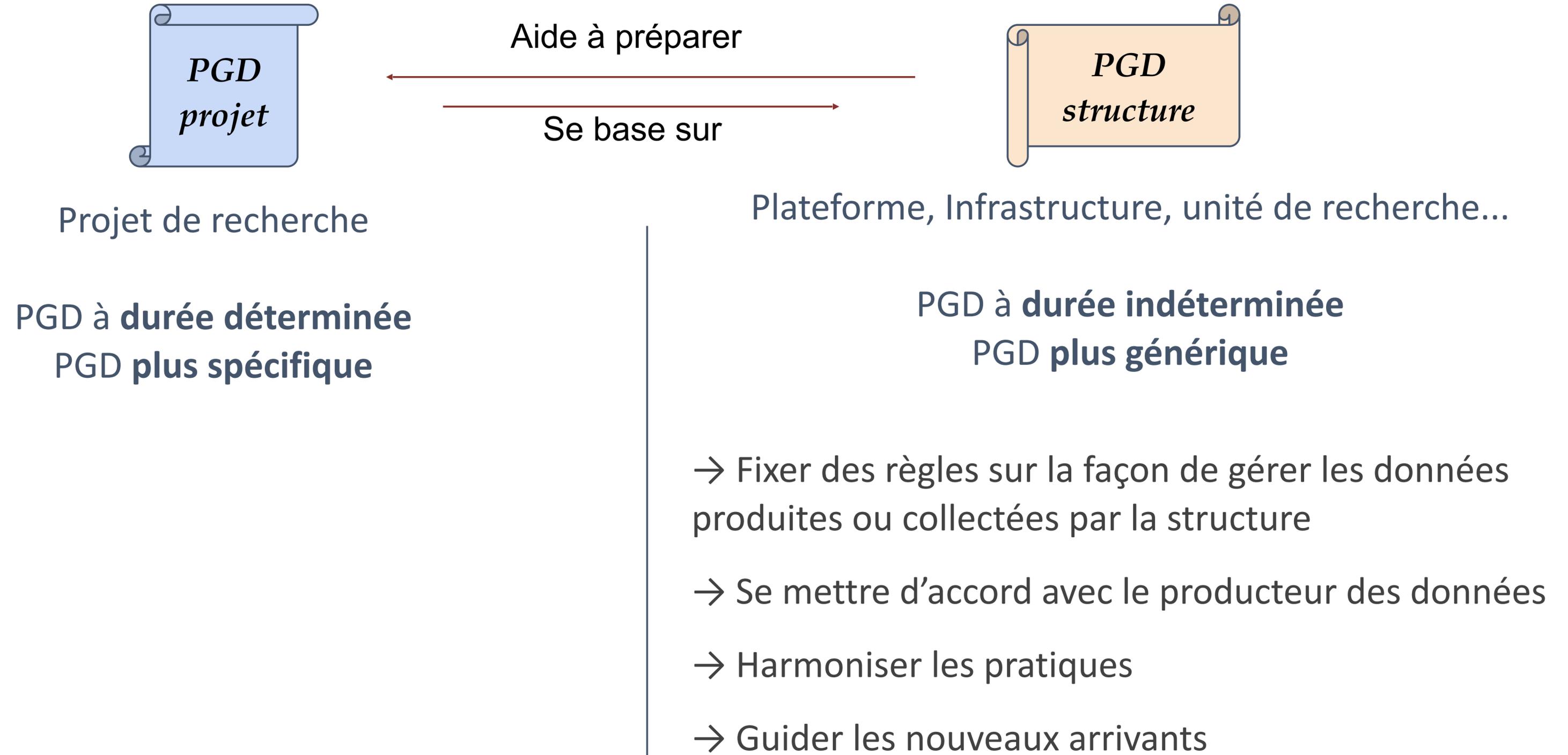


Liste des personnes contribuant à la gestion des produits de recherche au cours d'un projet et leurs rôles. L'attribution d'un rôle à une personne s'effectue dans l'onglet "Rédiger".

Nom	Affiliation	Rôles attribués (Produits de recherche associés)	
Anne-Caroline Delétoille	Institut Pasteur	<ul style="list-style-type: none"><li>• Personne contact pour les données (JD2 - PCR)</li></ul>	 
Fanny SEBIRE	Institut Pasteur	<ul style="list-style-type: none"><li>• Coordinateur du projet</li><li>• Personne contact pour les données (JD1 - Images)</li><li>• Responsable du plan</li></ul>	 

→ Objectif : Établir le rôle de chacun

- Définir les responsabilités





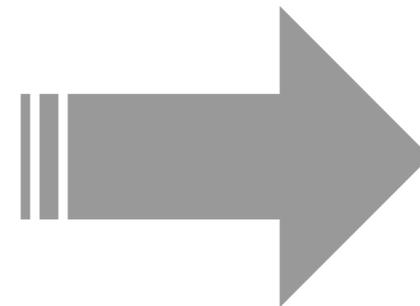
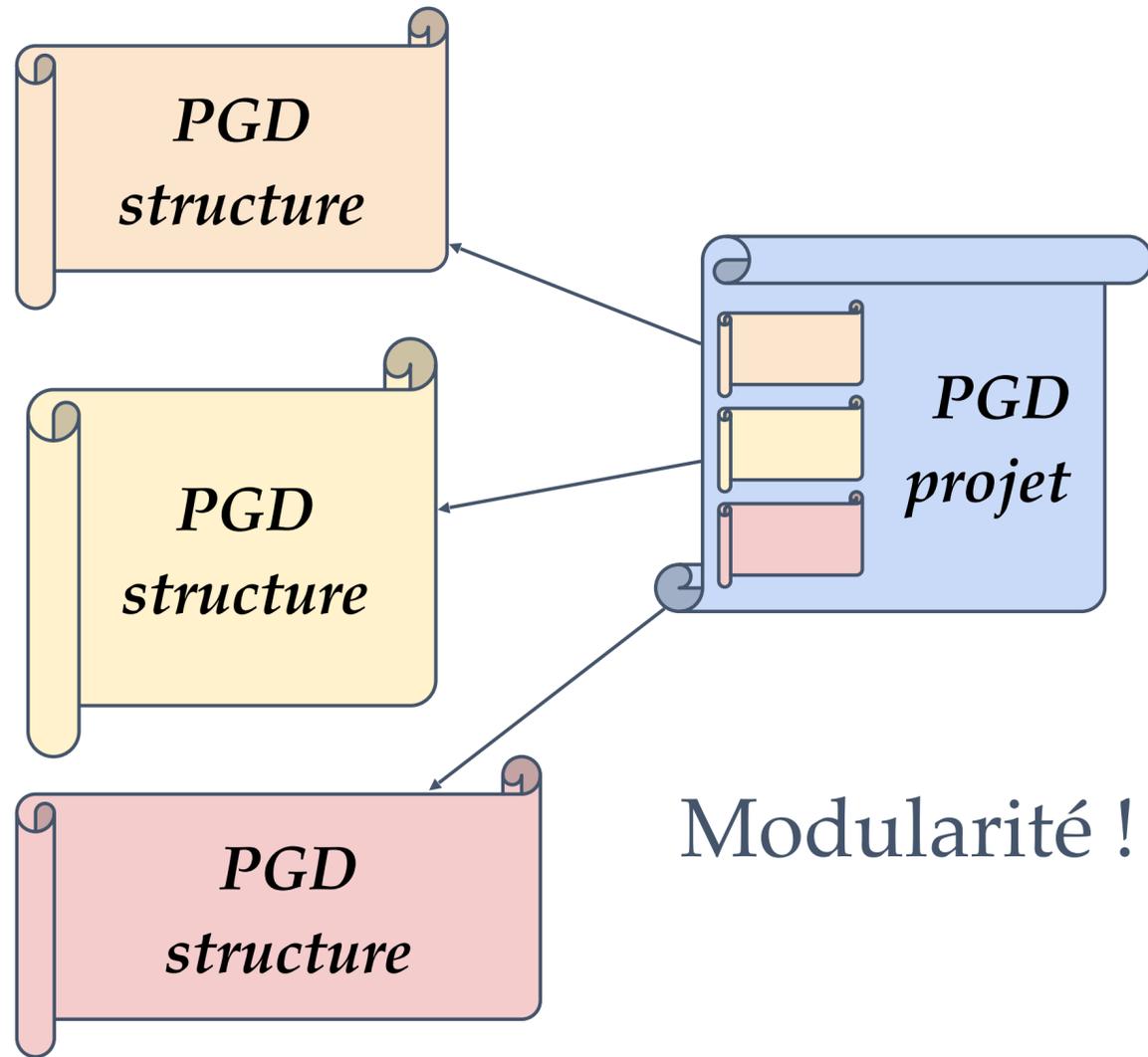
PGD publics sur DMP OPIDoR (projet et structure) : [https://dmp.opidor.fr/public\\_plans](https://dmp.opidor.fr/public_plans)

Exemples de PGD de projet :

- [PGD du projet INFRAVEC2](#)
- [Collection de 841 PGD H2020](#)
- [11 PGD primés](#)

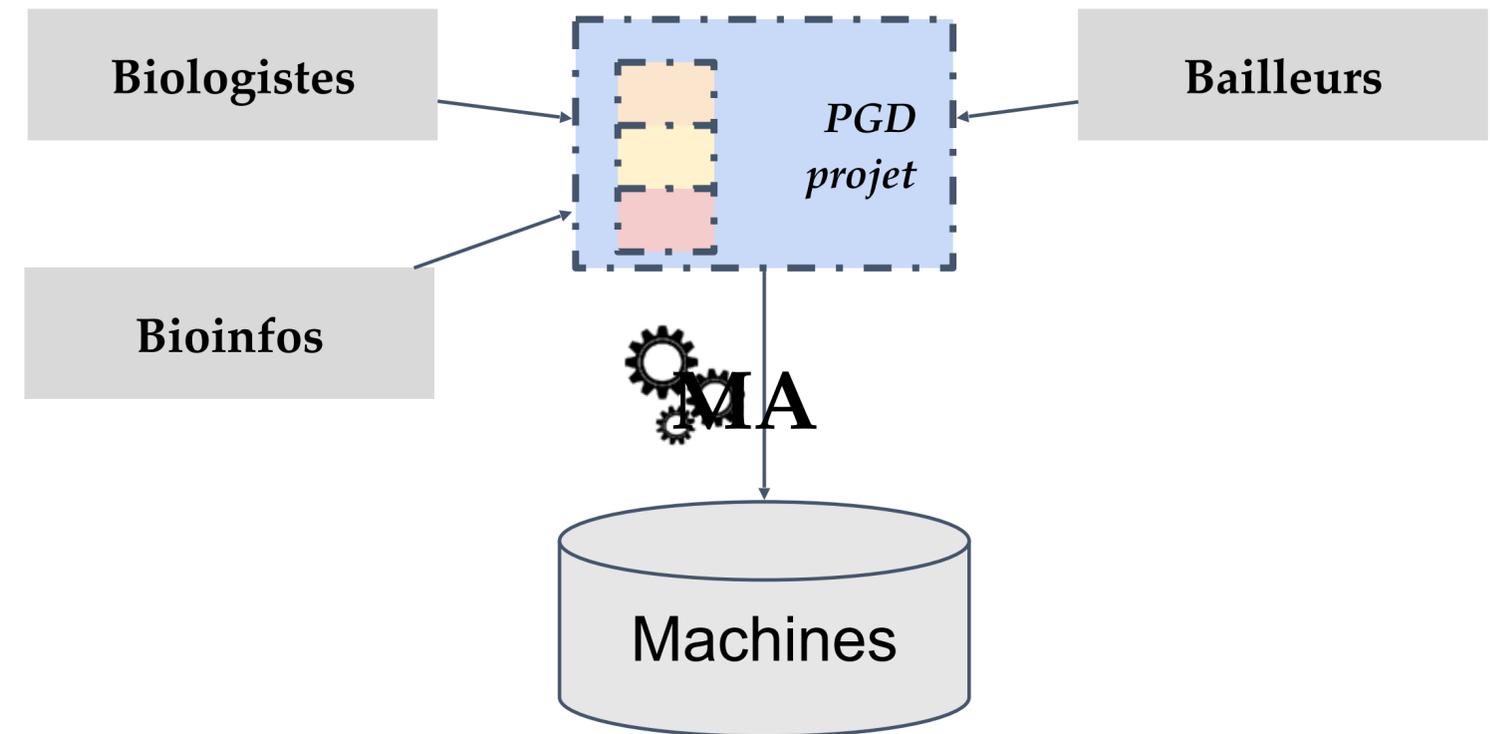
Exemples de PGD de structure :

- [Plan de Gestion des Données du CIRM-BIA](#) :
- [Data management plan of the Plant Bioinformatics Facility](#)



**Machine actionable DMP** : un plan de gestion lisible par les machines

**Objectif** : faire du Plan de gestion des Données un outil de configuration des environnements des infrastructures



# Différents outils pour rédiger des PGD

- [DMP OPIDoR](#) - solution nationale
- [DSW - Data Stewardship Wizard](#) - solution européenne (ELIXIR)
- [ARGOS](#) - solution de la Commission Européenne



DS Wizard

20210224\_exemple\_canevas\_IFB\_bioimage

Questionnaire | KOLs | Metrics | Preview | Documents | Settings

Chapters

- I. Préface ✓
- II. Introduction 1
- III. Informations générales 24
- IV. Données de la recherche 91

**1.a.2.a.2.a.3 Quelles mesures de contrôle de la qualité sont prises pour ce jeu de données ?**

Ceci concerne l'acquisition des images. Voir le premier white paper de <https://quarep.org/> <https://arxiv.org/pdf/2101.09153.pdf>. La liste d'options présentée ici est en anglais : ne devrait-on pas traduire cela pour être consistant/cohérent avec le reste de ce document ? (quoique les modes d'acquisition sont bien donnés en anglais...)

- a. Illumination power
- b. Detection system performance
- c. Field of view uniformity/flatness
- d. Chromatic aberrations
- e. Lateral and axial resolution
- f. Image quality
- g. Commentaires ☰

**1.a.2.a.2.a.4 Des versions différentes du jeu de données sont-elles créées ?**

- Get familiar with the **DMP Lifecycle** and excel in **Open and FAIR RDM planning**
- **Co-create DMPs** and manage workload
- **Publish and cite** DMPs as living documents
- **Configure** DMPs to tailored community needs
- **Link** DMPs to research outputs, EO SC services and the **OpenAIRE Research Graph**



- Liste exhaustive de vos données
- Réfléchir à
  - Où mettre les données pendant le projet
    - Volume, besoins d'accès
  - Où mettre les données pour les partager
    - Comment les documenter
- Comment gérer « le reste »
  - CR, annexes, gestion ...
- Stratégie de nommage de tous les fichiers

- Un grand merci à :



H. Chiapello



T. Denecker



J.-F. Dufayard



P. Lieby



L. Maurel



G. Sarah



J. Seiler



Christophe Antoniewski



Naïra Naouar



Anne-Caroline Delétoille



Fanny Sébire



Valentin  
Loux



Cédric Midoux



Célia  
Michotey



Cyril Pommier



*Et merci à toutes celles et ceux qui nous ont stimulé par leurs questions et encouragé par leur enthousiasme*