# Code & données de la recherche (1/2)

#### Formateurs: Code & données de la recherche

- Linda Angulo, Chargée de Mission Données de la Recherche
- Matteo Camier, Responsable HPC @pmcs2i.ec-lyon

#### **Formateurs: Publications**

- Nicolas Jardin, Directeur Adjoint Biblio. Michel Serres
- Stéphanie Lamaison, Bibliothécaire Biblio. Michel Serres



36, avenue Guy de Collongue 69130 Écully - France +33 (0)4 72 18 60 00

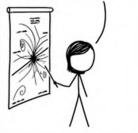
www.ec-lyon.fr





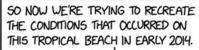
# Code & données de la recherche (1/2)

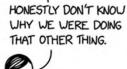
OUR LAB WAS TRYING TO RECREATE THE CONDITIONS THAT OCCURRED SECONDS AFTER THE BIG BANG.

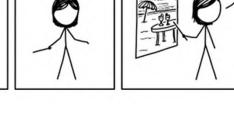










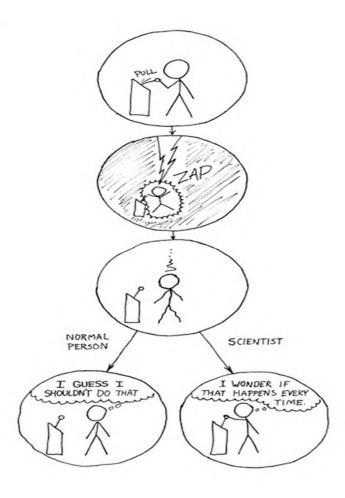






36, avenue Guy de Collongue 69130 Écully - France +33 (0)4 72 18 60 00

www.ec-lyon.fr



#### **Budapest Open Access Initiative, 22eme anniversaire**



- Héberger la recherche OA sur des infrastructures ouvertes
- Réformer l'évaluation de la recherche et les récompenses pour améliorer les incitations
- Favoriser des canaux de publication et de distribution inclusifs
- Dépenser l'argent pour publier la recherche OA en gardant en mémoire les objectifs de l'OA

**#BOAI20** 



#### **2eme Plan Nationale Pour La Science Ouverte 2021-2024**

#### Les 4 axes du 2e Plan national pour la science ouverte

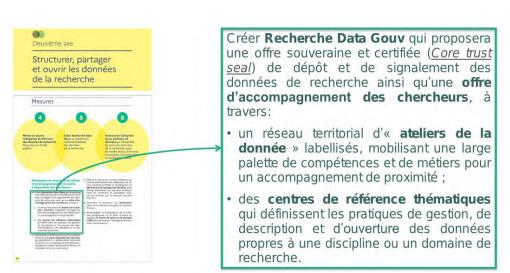
- Généraliser l'accès ouvert aux publications
- Structurer, partager et ouvrir les données de la recherche
- Ouvrir et promouvoir les codes sources produits par la recherche
- Transformer les pratiques pour faire de la science ouverte le principe par défaut

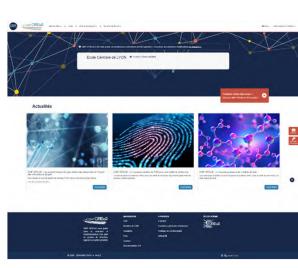


objectif de 100% de publications en accès ouvert en 2030 fixé par la loi de programmation de la recherche.

#### 2e Feuille de Route Science Ouverts, à l'ECL

- Ouvrir le processus de recherche et les résultats à tous les acteurs de la société
- Faciliter les collaborations scientifiques et tirer profit de la science des données pour nos enjeux de recherche
- Saisir les opportunités et maîtriser les risques potentiels liés à l'ouverture
- Chargée des données de Recherche

















#### Acteurs d'Ouverture de Code et Données de la Recherche

**Clarisse MARANDIN Direction Bibliothèque**  **Christophe CORRE** Direction de la Recherche

**Bénédicte MARTIN Resp Centrale Innov** Resp Aff. EU DPRV

**Elisabeth DALVERNY Responsable DRPV** 

**Camille ZAMI-PIERRE Juriste ECL** 

Nicolas JARDIN Resp. services recherche

**Matteo CAMIER** Resp. tech Pôle calcul / gestion des données de la recherche

**Laurine MAIRE** Chargée d'Aff. **Centrale Innovation** 

**Véronique MERAT** Ingénierie de projet **DPRV** 

**Guillaume EMPTAZ FSD** 

**Stéphanie LAMAISON** HAL/open accès

**Anne CADIOU IR CNRS LMFA** Pôle calcul

**Marine PICO** Chargée projets EU **Centrale Innovation** 

**Angélique** CATEUX **Suivi contrats** Recherche

Christophe **FESSART DPO** 

**Linda ANGULO LOPEZ Chargée mission** données

**Benoît PIER Dir. Rech CNRS LMFA Membre GT COSO** 



**AMI données GT** juridique



**Robin MATEJICEK** Suivi projets **Recherche ENISE** 

**Groupe de Travaille Science Ouvert de Centrale Lyon** 

GT Traitement des données DataLyste Fork



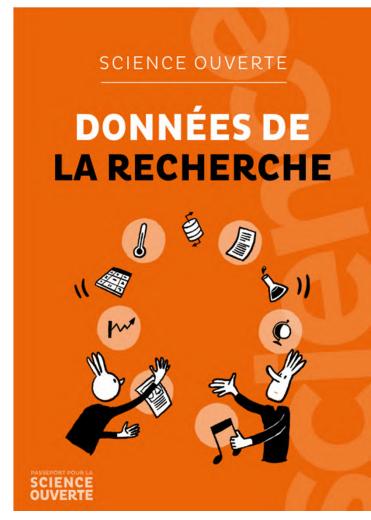




#### **Recommandations aux doctorantes**











### Politiques de sécurité et de confidentialité

La Loi République numérique

#### **Données Publiques**

# **Exceptions**

- Droits d'auteur
- Droit sui generis
  - Brevets
- Obtentions végétales
  - Essais cliniques
    - Biodiversité

# **RGPD**

# Données à caractère personnel avec dérogations :

- durée de conservation;
- recherches exploratoires;
- droit à l'information des personnes;
  - droit à l'opposition

#### Aspects juridiques des données de la recherche

#### **Droits d'auteur et Propriété intellectuelle**

- Chercheurs conservent leurs droits d'auteur sur les œuvres créées dans leurs fonctions et peuvent les céder.
- Protège les investissements dans la création de bases de données et s'oppose à des extractions substantielles de contenu. Valable 15 ans, renouvelable.



# Loi République numérique (2016) et Open Data

- Données de recherche, assimilées à des données publiques, doivent être accessibles en ligne et réutilisables gratuitement.
- Exception Text and Data Mining (TDM), permise pour la recherche publique sous conditions (non-lucratif, accès licite, sécurité des systèmes et des fichiers).

#### **Exception & Politiques de Sécurité et de Confidentialité**

La Loi République numérique

#### **Données Publiques**

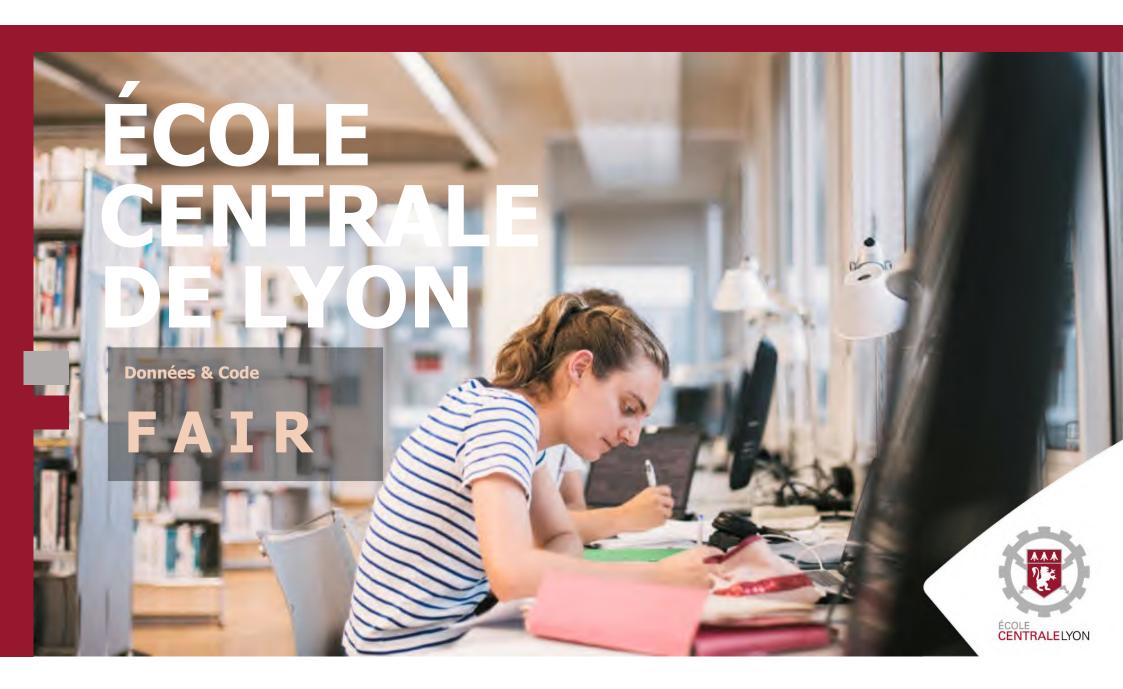
# **Exceptions**

- Obtention Droits d'auteur
  - Droit sui generis
    - Brevets
    - Végétales
  - Essais cliniques
    - Biodiversité

# **RGPD**

# Données à caractère personnel avec dérogations :

- durée de conservation;
- recherches exploratoires;
- droit à l'information des personnes;
  - droit à l'opposition



#### FAIR (Facile à trouver, Accessible, Interopérable, Réutilisable)

Facile à trouver

DOI ; métadonnées indiquant l'identifiant ; entrepôt permettant la recherche ; Décrie les métadonnées riches et utilisant des vocabulaires

**Accessible** 

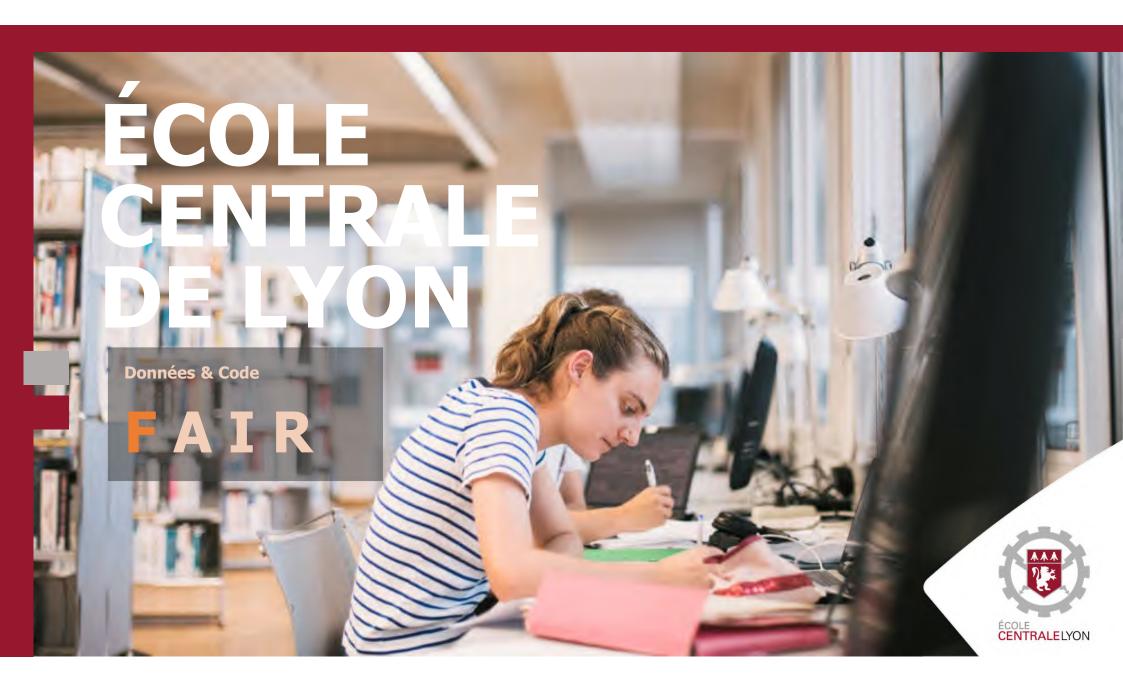
Données accessibles avec le DOI; métadonnées toujours disponibles; Déposer des données ouverte, éventuellement sous conditions

Interopérable

Métadonnées standardisées et machine-readable ; Format de fichiers ouverts et machine-readable ; vocabulaires avec URI et URL

Réutilisable

Licences disponibles ; Décrie provenance des données ; standards correspondant à la communauté





#### Plan de gestion des données

- L'équilibre entre ouverture et sécurité
- Le PGD évolue tout au long du projet
- Edité collectivement par exemple sur DMP OPIDoR
- Agences de financement de la recherche demandent désormais la fourniture d'un PGD



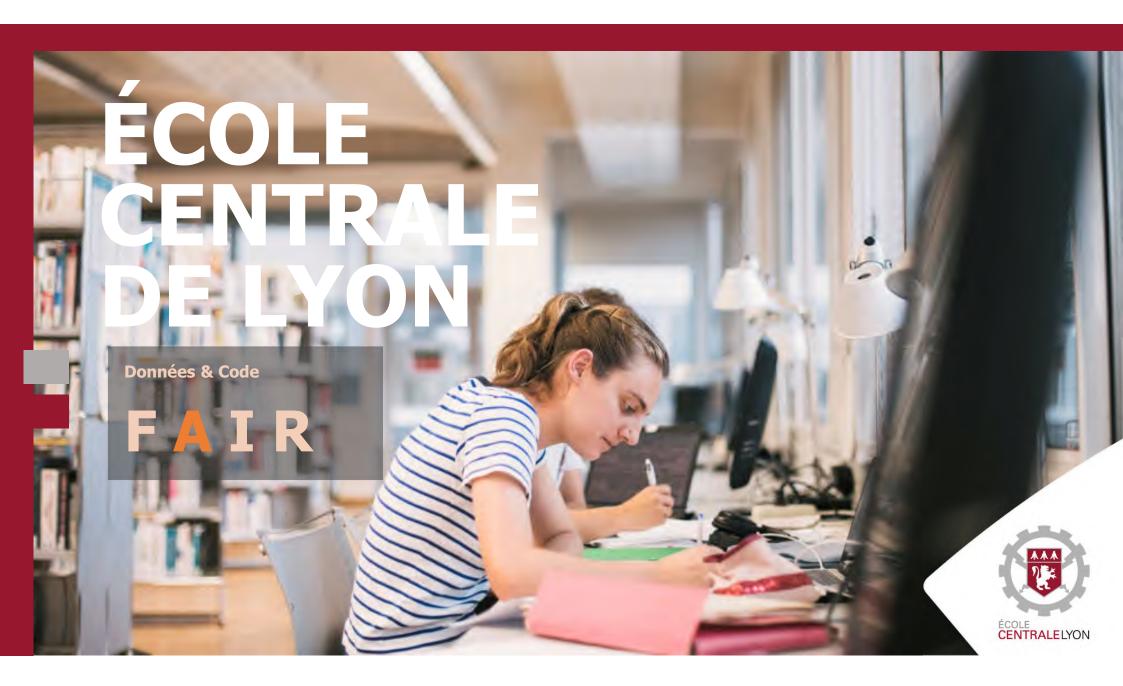


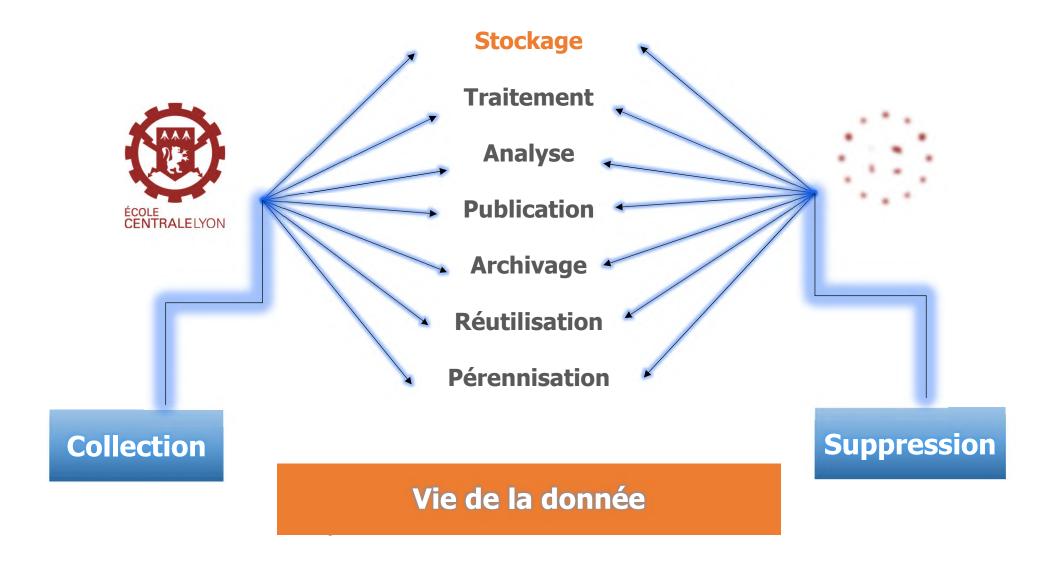






- Activités principales d'une équipe de recherche
- Garantie de la confidentialité de la donnée doit être intégré dès cette première étape
- Recueil doit se faire de façon justifiée, sécurisée et cloisonnée
- Consentement des personnes, active ou passive
- Limitation facilite le stockage et surtout la protection des données



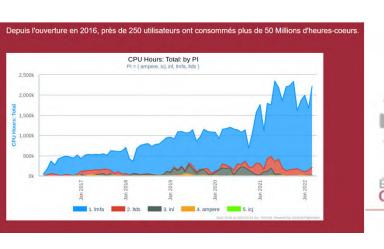




- DMP incitent à mieux préparer cette étape
- Fichiers nommés, retrouver facilement & distinguer versions
- La sauvegarde, 3 copies sur 2 supports différents, dont 1 copie à distance.
- Stocker ses données en cours de projet
- Données confidentialité, nécessite un espace adapté et s'interroger sur la durée de conservation
- Stockage Centraliser

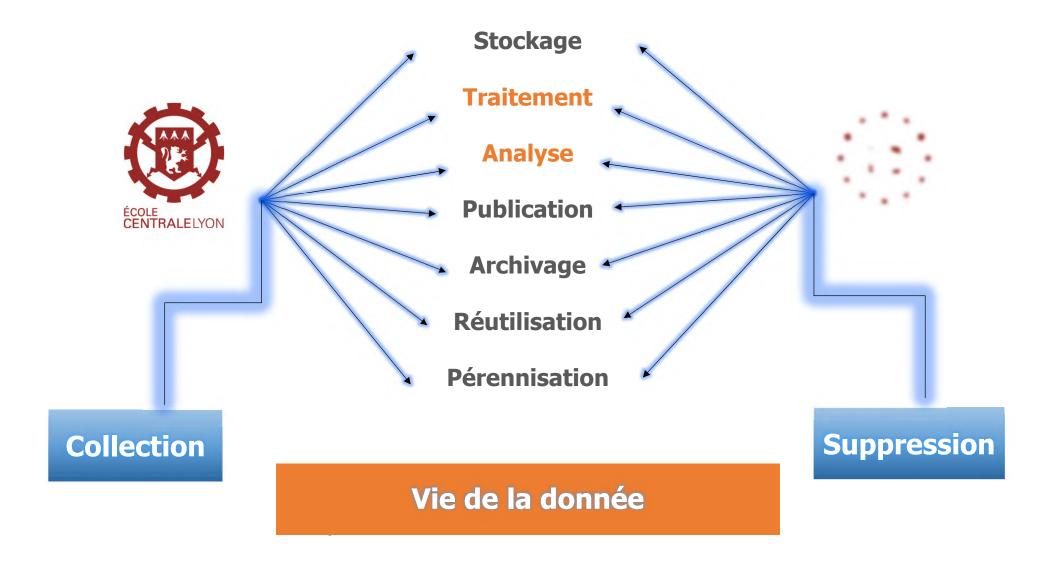
#### Le Pôle de Calcul de l'École Centrale

- Hébergement dans le datacentre ECL, au CRI22
- Un cluster de calcul parallèle, Newton (+3500 cœurs CPU, 8 GPU)
- Espaces de stockage partagé (+2 Po)
- Visualisation et Post-traitement des Données
- Ateliers et Séminaires, Informatiques et Calcul Scientifique







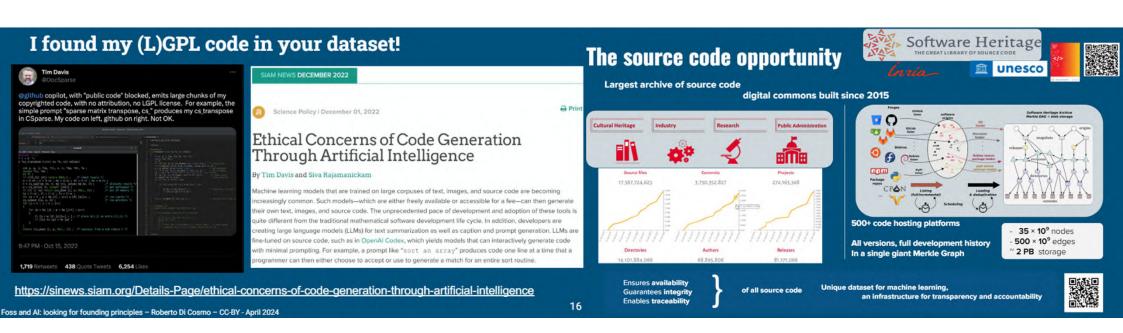




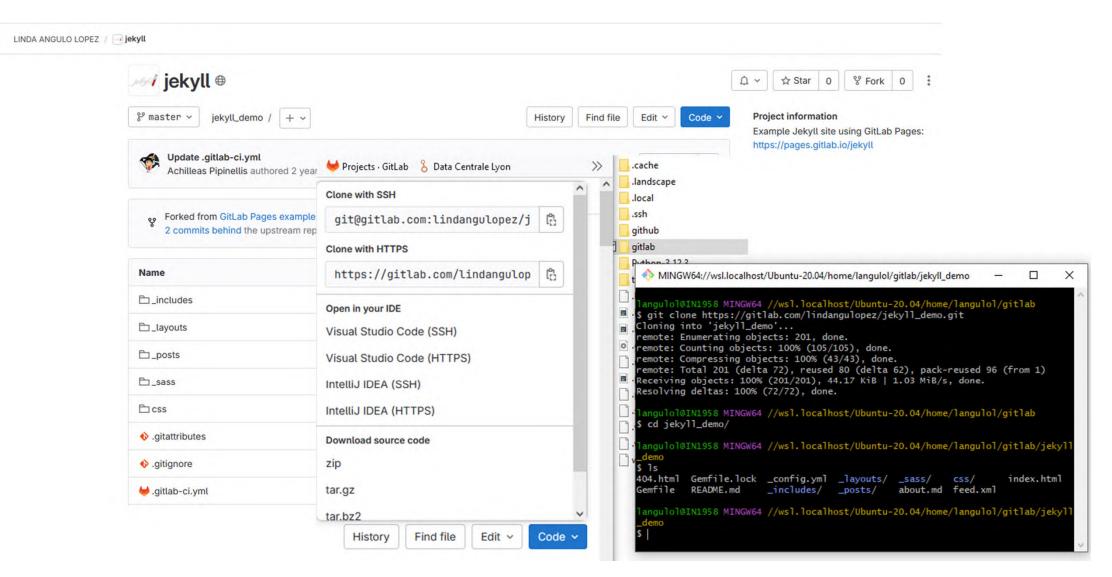
- Transcrite, traduite, vérifiée, validée & nettoyée
- Données personnelles identifiantes, anonymisées
- Code soient également accessibles
- L'analyse permet de décrire les données
- L'interprétation, donner sens à ces résultats
- Métadonnées persistantes

### **Duplication et Création de Code**

- La génération des données
- L'analyse et traitement des données
- Duplication et modification des projets

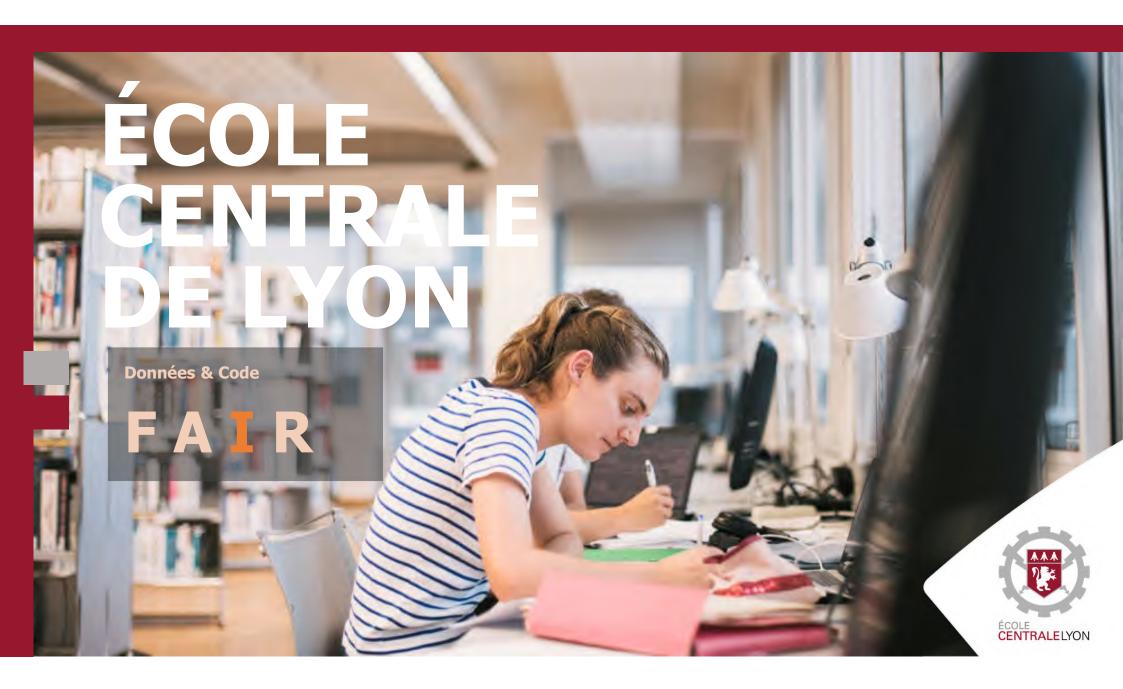


### Forkez et contribuez à un projet Git





- Création d'une branche thématique à partir de la branche master,
- validation de quelques améliorations (commit),
- poussée de la branche thématique sur votre projet Git (push),
- ouverture d'une requête de tirage sur Git (Pull request),
- discussion et éventuellement possibilité de nouvelles validations (commit).
- Le propriétaire du projet fusionne (merge) ou ferme (close) la requête de tirage.
- Synchronisation de la branche master mise à jour avec celle de votre propre dépôt.



#### Préparer les fichiers ouvert, auto documenté

#### **NetCDF**

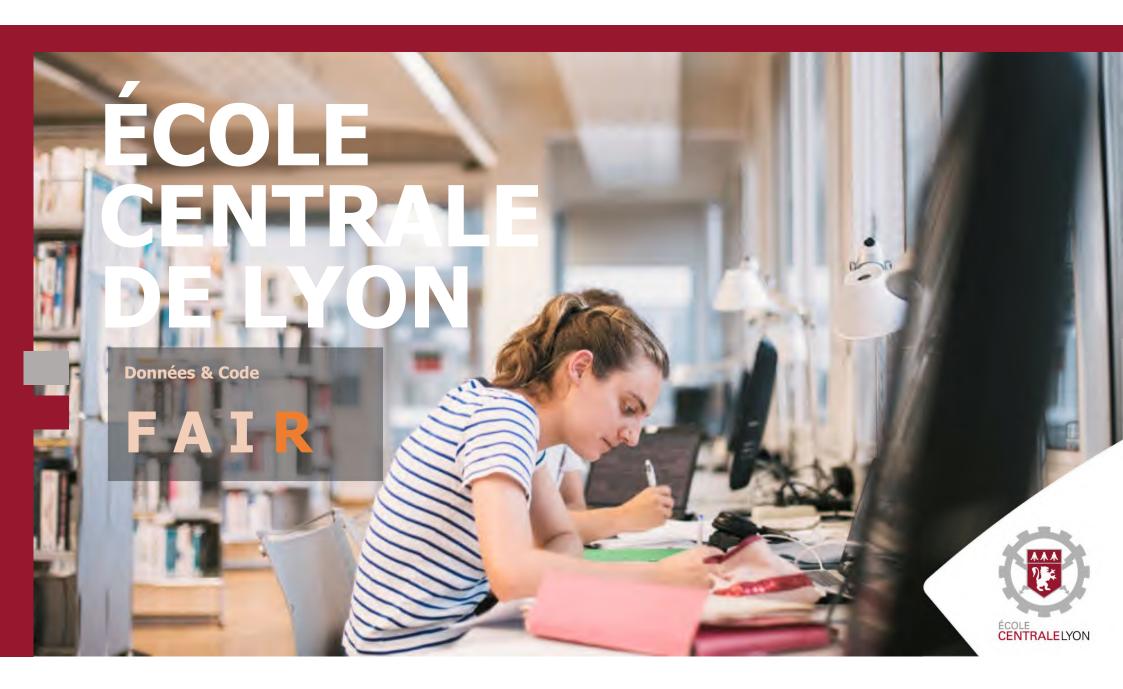
- Représenter et formater des données dimensionnées sous forme de tableaux
- intégrant les métadonnées
  directement dans l'entête du fichier



#### HDF5

- Type conteneur
- Structure de fichier hiérarchique
- Simuler des données grâce au calcul intensif
- compression et d'écriture/lecture parallèles







#### Valoriser, promouvoir et partager

#### **Publication d'articles scientifiques**

- une petite partie des données ou métadonnées est accessible

#### **Bases pluridisciplinaires, Recherche Data Gouv**

- Recommandations du financeur
- La possibilité de pérenniser l'accès
- Entrepôts certifiée, critères CoreTrustSeal
- Niveau de sécurité nécessaire

#### **Data-papers**

- la description détaillée des jeux, métadonnées

### **Crédits image et sources principales**

#### **Formations**

- Ecole Centrale de Lyon, A.S.P.I.C.S.
- FUN MOOCS
  - Sorbonne Université, La science ouverte
  - INRIA, Recherche reproductible : principes méthodologiques pour une science transparente
  - INRIA, Recherche reproductible II: Practices and tools for managing computations and data

#### Ressources en ligne

- <u>https://mi-gt-donnees.pages.math.unistra.fr/site/guide.html</u>
- https://www.recherche-reproductible.fr/
- https://www.ouvrirlascience.fr/science-ouverte-donnees-de-la-recherche
- https://recherche.data.gouv.fr/fr/page/classes-virtuelles
- https://dmp.opidor.fr
- https://cat.opidor.fr
- https://printempsdeladonnee.fr/calendrier/
- https://doranum.fr/
- https://www.canal-u.tv/chaines/ad/journee-gitlab/software-heritage-l-archive-universelle-a-la-croisee-desforges
- https://xkcd.com

# **GLOSSAIRE DONNEES DE LA RECHERCHE (1/2)**

	Entrepôt de données	Espace numérique dans lequel on peut chercher ou déposer des données de recherche. Un entrepôt ne recueille que des jeux de données achevés, à la différence d'une plateforme de stockage qui doit être utilisée tant que l'on travaille encore sur ses données.
	FAIR	Findable Accessible Interoperable Reusable, ensemble de bonnes pratiques à suivre pour gérer et partager ses données au mieux.
	Horizon Europe	Programme européen pour la recherche et le développement pour 2021-2027. Le programme précédent s'intitulait « Horizon 2020 »
	Inist	Institut national de l'information scientifique et technique, unité du CNRS chargée de faciliter l'accès à l'information scientifique et aux données de la recherche ainsi que d'accompagner les chercheurs dans ces démarches.
	Loi CADA	Loi française de 1978 créant la Commission d'accès aux documents administratifs. Elle permet aux citoyens de demander l'accès aux documents administratifs rendus publics.
	Loi pour une République numérique	Loi française de 2016 mettant en place, entre autres, l'ouverture par défaut des données publiques, la publication en accès ouvert des articles de recherche financés au moins à 50% par des fonds publics (avec un embargo), l'ouverture des données de la recherche sous conditions.
	Métadonnées	Ensemble d'informations permettant de décrire une donnée. Par exemple, les métadonnées de données d'enquêtes, vont être le nom des personnes interrogées, leur âge, la date de l'entretien Et l'enregistrement audio de l'entretien sera la donnée en elle-même.



