

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 Écullly - France
+33 (0)4 72 18 60 00

www.ec-lyon.fr

ÉCOLE CENTRALE DE LYON

Formation doctorale : Bibliothèque Michel Serres



ÉCOLE
CENTRALE LYON

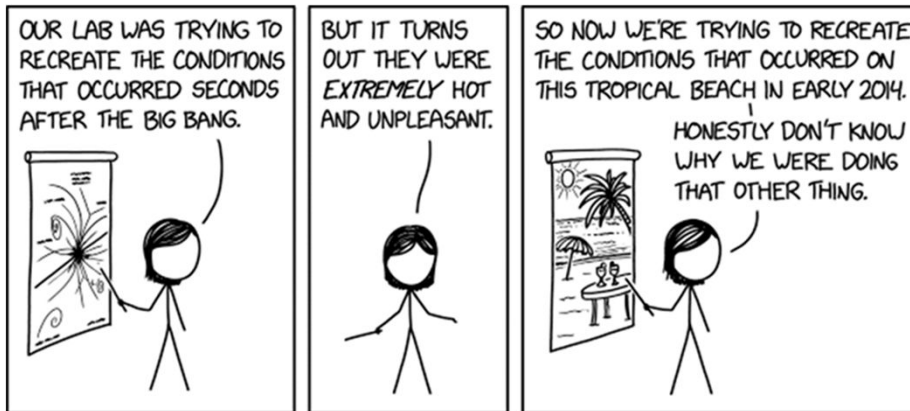
ÉCOLE CENTRALE DE LYON

Gestion des données de recherche : optimisation de l'impact de vos travaux(1/2).



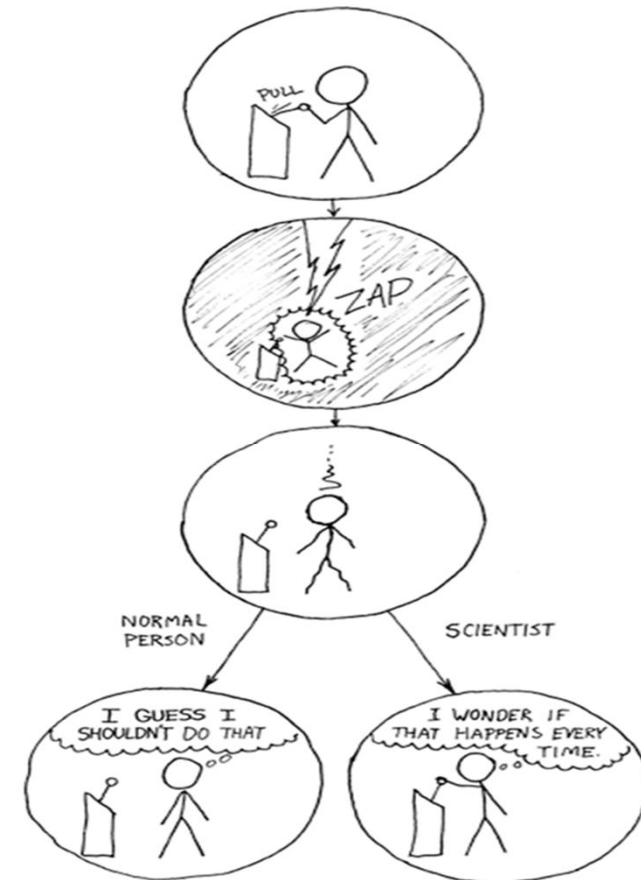
ÉCOLE
CENTRALE LYON

Code & données de la recherche (1/2)



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

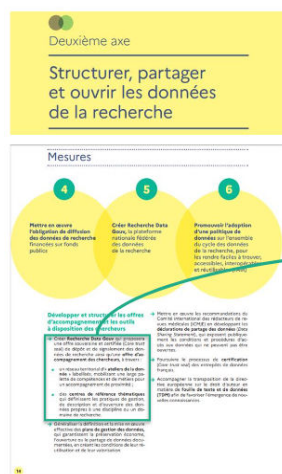


Deuxième
Plan national
pour la science
ouverte

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

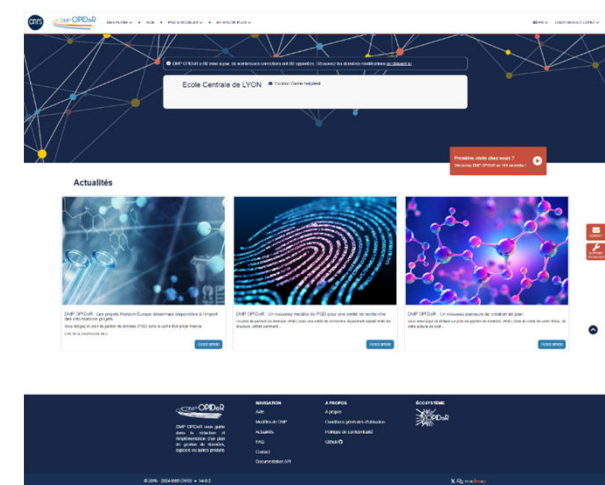
2e Feuille de Route Science Ouverts, à l'ECL

- 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



Créer **Recherche Data Gov** qui proposera une offre souveraine et certifiée (*Core trust seal*) de dépôt et de signalement des données de recherche ainsi qu'une **offre d'accompagnement des chercheurs**, à travers :

- un réseau territorial d'« **ateliers de la donnée** » labellisés, mobilisant une large palette de compétences et de métiers pour un accompagnement de proximité ;
- des **centres de référence thématiques** qui définissent les pratiques de gestion, de description et d'ouverture des données propres à une discipline ou un domaine 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

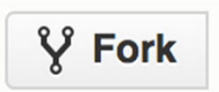

Benoît PIER
 Dir. Rech CNRS LMFA
 Membre GT COSO



AMI données
GT juridique


Robin MATEJICEK
 Suivi projets
 Recherche ENISE

Groupe de Travail Science Ouvert de Centrale Lyon
GT Traitement des données DataLyste



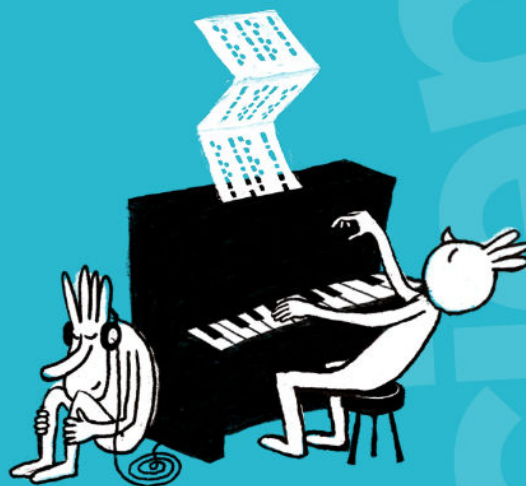
Recommandations aux doctorantes

PASSEPORT POUR LA SCIENCE OUVERTE



GUIDE PRATIQUE
À L'USAGE DES
DOCTORANTES ET
DES DOCTORANTS

SCIENCE OUVERTE CODES ET LOGICIELS



PASSEPORT POUR LA
SCIENCE
OUVERTE

SCIENCE OUVERTE

DONNÉES DE LA RECHERCHE



PASSEPORT POUR LA
SCIENCE
OUVERTE

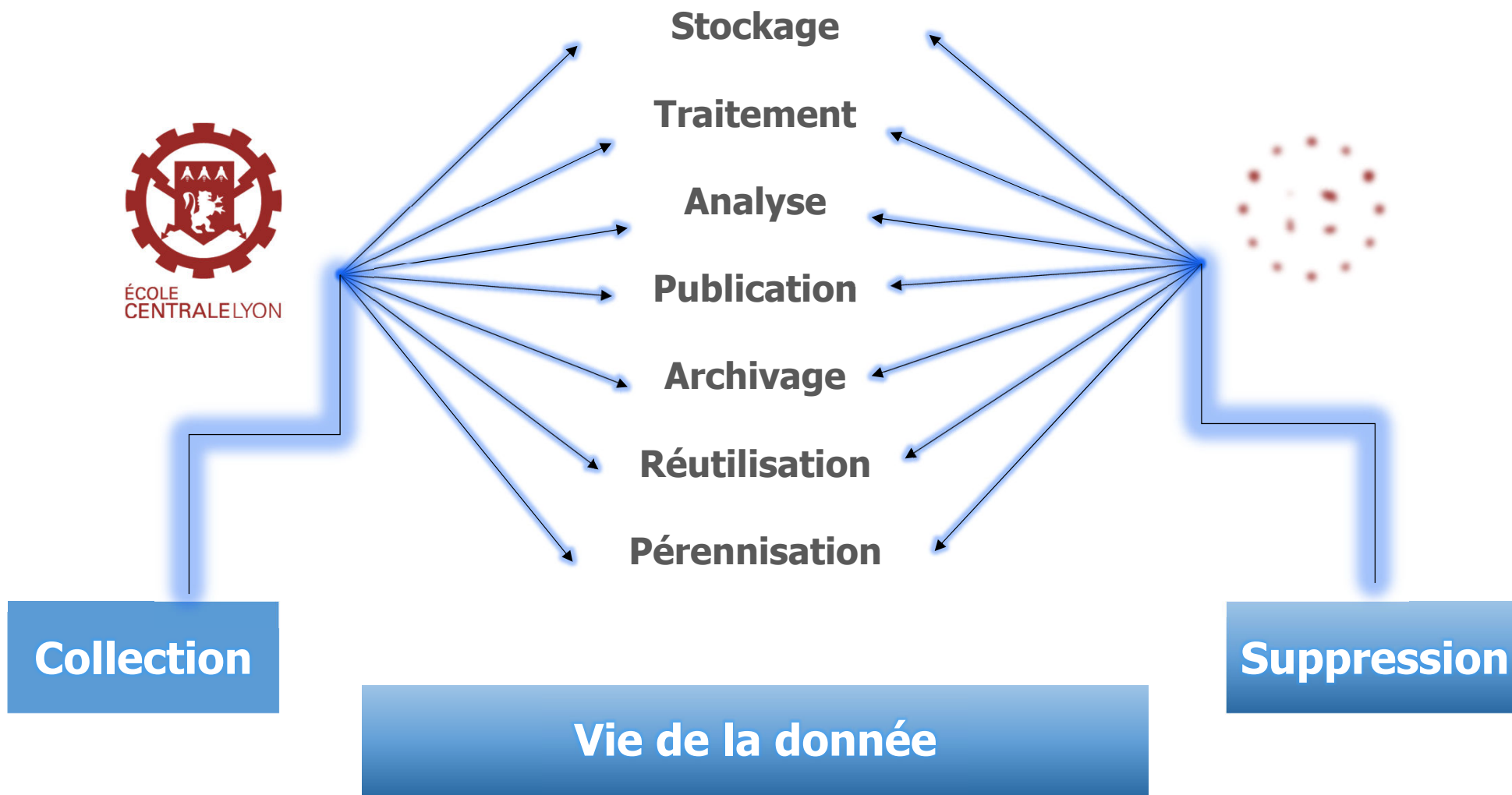
ÉCOLE CENTRALE DE LYON

Gestion des données de recherche : optimisation de l'impact de vos travaux(1/2).



ÉCOLE
CENTRALE LYON

Cycle de vie des données



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

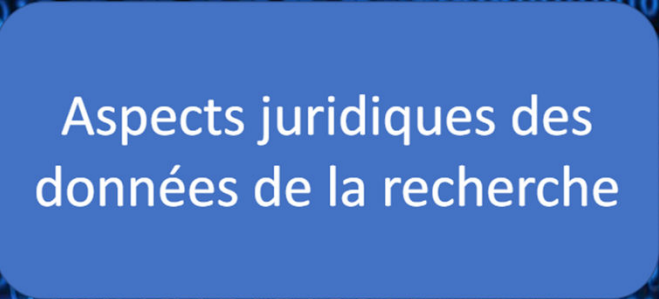
Données à caractère personnel avec dérogations :

- 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.



Aspects juridiques des données de la recherche

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

ÉCOLE CENTRALE DE LYON

Données & Code

FAIR



ÉCOLE
CENTRALE LYON



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

Facile à trouver

DOI ; métadonnées indiquant l'identifiant ; entrepôt permettant la recherche ; Décrit 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 ; Décrit provenance des données ; standards correspondant à la communauté

ÉCOLE CENTRALE DE LYON

Données & Code

F A I R



ÉCOLE
CENTRALE LYON

ÉCOLE CENTRALE DE LYON

Plan de Gestion des Données



ÉCOLE
CENTRALE LYON

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



Cycle de vie des données



Collection

Stockage
Traitement
Analyse
Publication
Archivage
Réutilisation
Pérennisation

Suppression

Vie de la donnée





La collecte de données, RGPD

- **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**

ÉCOLE CENTRALE DE LYON

Données & Code

FAIR



ÉCOLE
CENTRALE LYON

Cycle de vie des données





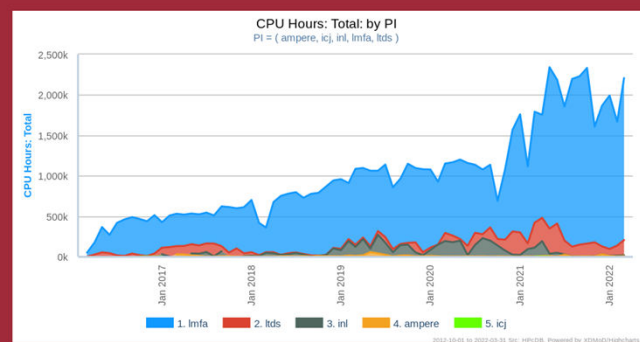
Le stockage des données

- **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

Depuis l'ouverture en 2016, près de 250 utilisateurs ont consommés plus de 50 Millions d'heures-coeurs.



Bienvenue au PMCS21

English Français
Laboratoire de Mécanique des Fluides et d'Acoustique - UMR 5509

LE LABORATOIRE ACTUALITÉS ÉQUIPES SUPPORT À LA RECHERCHE PUBLICATIONS FORMATIONS EMPLOI

Accueil > Support à la Recherche > Systèmes, Infrastructures et Calcul > Ateliers et Formations

NOS TUTELLES

- cnrs
- UNIVERSITÉ DE LYON
- INSA

NOS PARTENAIRES

- UNIVERSITÉ DE LYON
- INSTITUT CARNOT

Ateliers et Formations

- Architectures - partie 1
- Architectures - partie 2
- Prefabrique au calcul parallèle
- Introduction à l'informatique scientifique
- Précision numérique
- Atelier Précision numérique (intranet)
- Les outils du développeur de codes de calcul
- Atelier Unix (intranet) fiche solutions

ANNUAIRE

- Bash
- Scd
- Ask
- Introduction à ssh et aux clés ssh
- Choix des langages de programmation
- Comment structurer son code ?
- Introduction à git (2016) (cheatsheet)
- Atelier Git et Github (intranet)
- Travailler avec des branches avec git et gitlab
- Introduction à OpenMP

CONTACT ET ACCÈS

- Introduction à MPI
- Introduction au calcul sur GPU

INTRANET

- Matériau de la plateforme PMCS21

Cycle de vie des données





Avant de pouvoir être utilisée

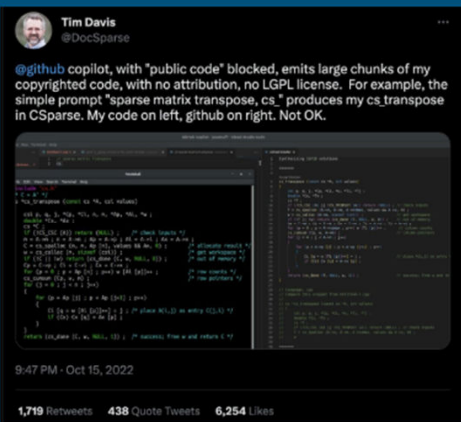
- **Transcrite, traduite, vérifiée, validée & nettoyée**
- **Données personnelles identifiantes, anonymisées**
- **Code soient également accessibles**
- **L'analyse permet de décrire les données**
- **L'interprétation, donner sens à ces résultats**
- **Métadonnées persistantes**

Duplication et Création de Code



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

I found my (L)GPL code in your dataset!



SIAM NEWS DECEMBER 2022

Science Policy | December 01, 2022

Ethical Concerns of Code Generation Through Artificial Intelligence

By Tim Davis and Siva Rajamanickam

Machine learning models that are trained on large corpuses of text, images, and source code are becoming increasingly common. Such models—which are either freely available or accessible for a fee—can then generate their own text, images, and source code. The unprecedented pace of development and adoption of these tools is quite different from the traditional mathematical software development life cycle. In addition, developers are creating large language models (LLMs) for text summarization as well as caption and prompt generation. LLMs are fine-tuned on source code, such as in *OpenAI Codex*, which yields models that can interactively generate code with minimal prompting. For example, a prompt like “sort an array” produces code one line at a time that a programmer can then either choose to accept or use to generate a match for an entire sort routine.

<https://sinews.siam.org/Details-Page/ethical-concerns-of-code-generation-through-artificial-intelligence>

The source code opportunity

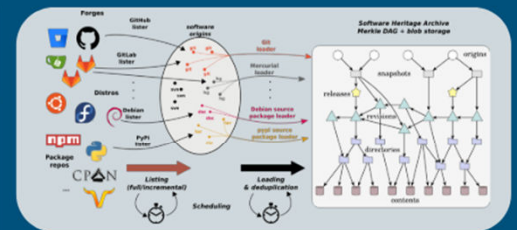
Largest archive of source code digital commons built since 2015



Ensures availability
Guarantees integrity
Enables traceability

of all source code

Unique dataset for machine learning, an infrastructure for transparency and accountability



500+ code hosting platforms

All versions, full development history
In a single giant Merkle Graph

- 35 × 10⁹ nodes
- 500 × 10⁹ edges
~ 2 PB storage



Forkez et contribuez à un projet Git

LINDA ANGULO LOPEZ / jekyll



🔔 Star 0 🍴 Fork 0

🔗 master jekyll_demo / +

History Find file Edit Code

Project information
Example Jekyll site using GitLab Pages:
<https://pages.gitlab.io/jekyll>

Update .gitlab-ci.yml
Achilleas Pipinellis authored 2 year

Projects · GitLab Data Centrale Lyon

Forked from [GitLab Pages example](#)
2 commits behind the upstream rep

Name
_includes
_layouts
_posts
_sass
css
.gitattributes
.gitignore
.gitlab-ci.yml

Clone with SSH
git@gitlab.com:lindangulopez/j

Clone with HTTPS
https://gitlab.com/lindangulop

Open in your IDE

- Visual Studio Code (SSH)
- Visual Studio Code (HTTPS)
- IntelliJ IDEA (SSH)
- IntelliJ IDEA (HTTPS)

Download source code

- zip
- tar.gz
- tar.bz2

History Find file Edit Code

- .cache
- .landscape
- .local
- .ssh
- github
- gitlab

```
MINGW64://wsl.localhost/Ubuntu-20.04/home/langulo1/gitlab/jekyll_demo
langu1o1@IN1958 MINGW64 //wsl.localhost/Ubuntu-20.04/home/langulo1/gitlab
$ git clone https://gitlab.com/lindangulopez/jekyll_demo.git
Cloning into 'jekyll_demo'...
remote: Enumerating objects: 201, done.
remote: Counting objects: 100% (105/105), done.
remote: Compressing objects: 100% (43/43), done.
remote: Total 201 (delta 72), reused 80 (delta 62), pack-reused 96 (from 1)
Receiving objects: 100% (201/201), 44.17 KiB | 1.03 MiB/s, done.
Resolving deltas: 100% (72/72), done.
langu1o1@IN1958 MINGW64 //wsl.localhost/Ubuntu-20.04/home/langulo1/gitlab
$ cd jekyll_demo/
langu1o1@IN1958 MINGW64 //wsl.localhost/Ubuntu-20.04/home/langulo1/gitlab/jekyll
_demo
$ ls
404.html  Gemfile.lock  _config.yml  _layouts/  _sass/      css/         index.html
Gemfile   README.md    _includes/  _posts/    about.md    feed.xml
```




Duplication du projet, par <<Fork>>

- **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.**

ÉCOLE CENTRALE DE LYON

Données & Code

FAIR



ÉCOLE
CENTRALE LYON

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



ÉCOLE CENTRALE DE LYON

Données & Code

FAIR



ÉCOLE
CENTRALE LYON

Cycle de vie des données





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 certifiées, 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

- <https://mi-gt-donnees.pages.math.unistra.fr/site/guide.html>
- <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-des-forges>
- <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.

ÉCOLE CENTRALE DE LYON

Je vous remercie pour votre attention.



ÉCOLE
CENTRALE LYON

ÉCOLE CENTRALE DE LYON

Gestion des données de recherche : optimisation de l'impact de vos travaux(2/2).



ÉCOLE
CENTRALE LYON