

Retour d'expérience sur la mise en place d'un inventaire et d'une politique de données

Monica Heintz & Hélène Gautier

C@fés Renatis, 7/04/2022



Présentation du Lesc

- Laboratoire d'anthropologie non spécialisé dans une aire géographique
- Environ 75 chercheurs et enseignants-chercheurs et 60 doctorants
- Pôle documentaire regroupant les ingénieurs en charge des ressources documentaires:
 - 2 ingénieurs de la Bibliothèque Eric de Dampierre (archives chercheurs + thesaurus GeoEthno)
 - Responsable du service audiovisuel (vidéos)
 - Chargée des ressources documentaires du CREM (Archives sonores du CNRS-Musée de l'Homme)
 - Chargée d'édition de corpus numériques

Des politiques en faveur de la science ouverte à plusieurs échelles

Variété des politiques et des niveaux:

- organisations internationales et nationales
- agences de financement
- universités et institutions de recherche
- revues

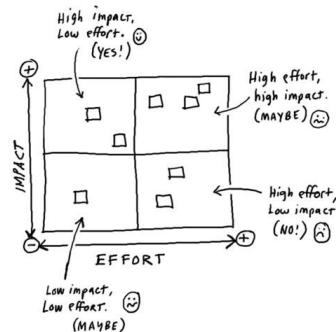


Une expérience dans le cadre d'ateliers de l'InSHS

- Cycle d'ateliers animé par Lionel Maurel et Emmanuelle Morlock
- Objectif: rédiger un document résumant les choix du laboratoire en matière de science ouverte, de gestion et de partage des données ainsi que les responsabilités des différents acteurs impliqués
- Plusieurs unités de recherche dont le Lesc (UMR 7186, Nanterre)
- Binôme chercheur-ingénieur
- Travail entre les ateliers et en collaboration avec d'autres collègues (notamment ingénieurs)

Une première étape: l'inventaire des données

- Recensement (non exhaustif) des jeux de données
 - Avoir une idée des types de données, des volumétries, de leur localisation, etc.
 - Inclure tout ce sur quoi le laboratoire investit des ressources (matérielles et/ou humaines) ou qui est pris en charge par un membre du laboratoire
- Outil pour la FAIRisation
 - Matrice impact/effort
- Aide à la rédaction de la politique
 - Identifier les priorités
 - Définir le périmètre des “données” pour le laboratoire



Identification du jeu				Caractérisation des données produites				Informations complémentaires				Modes de diffusion / archivage				Modalités de partage				
Programme de recherche	Désignation du jeu	Année de collecte ou de production	Nom du chercheur principal ayant collecté ou produit les données	Personne ou service en charge de la gestion des données	Caractérisation des données produites	Etat de finalisation (en cours / publié)	Format(s)	Taille (nombre de fichiers, espace disque)	Données brutes ou primaires (oui / non)	Type	Observations particulières (origine des données, réalisations potentielles, etc.)	Plan de gestion de données (oui/non/en cours)	Zones géographiques couvertes	Serveur et modalités de stockage / sauvegarde	Diffusion via le projet de publication numérique	Partage / archivage comme jeu de données réutilisables	Références/Indicateurs par des méta-moteurs, des portails, etc.	Périmètre du partage potentiel	Licence attribuée	
Projet Neurons	Neurons - neurosciences TD	2020		JL pda ingelavia	Transcriptomiques (transcriptome) et métabolomes descriptifs des cellules support combinés de séquençage génomique haut débit.	en cours	CSV	1000 fichiers, moins de 1 Mo	non	données brutes	Le dataset TD des données brutes et des métabolomes est un jeu de données de séquençage de génomique dans un format de données de séquençage de génomique et est destiné à être utilisé pour l'analyse.	cf PSC (en ligne)		Home-num (2ans)	Home-num (10 ans)	Zenodo/ Home-num (10 ans) (Zenodo)		Tout le monde	CC BY 4.0	
	Neurons images	2019		JL pda ingelavia	colfection d'images	en cours	TIFF, JPEG	1000 fichiers de fichiers image			cf PSC (en ligne)		Home-num (10 ans)	Home-num (10 ans)	Zenodo et Zenodo en continu		Tout le monde	cf Licence MIT pour les images de cellules CC BY pour les images de cellules et les données de cellules.		
	Neurons système de publication	2020		JL pda ingelavia	application personnalisée	en cours	JSON, CSV, XML, PDF	moins de 1 Mo	non	données informatiques, programmes		cf PSC (en ligne)		Home-num (10 ans)	Home-num (10 ans)	Zenodo		Tout le monde	cf Licence MIT pour les images de cellules et les données de cellules.	
Plateforme PathMed	PathMed et PathMed	2017		le développement d'applications, de jeux de données, de jeux de données	PathMed et PathMed	publié	JSON	1000 fichiers		jeux de données	les fichiers et métabolomes pub liés au PathMed sont pour une partie de génomique.		Home-num (10 ans)	Home-num (10 ans)	Home-num (10 ans)		Home-num (10 ans)	Home-num (10 ans)	Tout le monde	CC-BY-NC et CC-BY-NC pour les images, CC-BY pour les métabolomes.

<https://bit.ly/3KdR05I>

L'inventaire de données au Lesc en pratique

- Adaptation du modèle proposé avec ajout de champs
- Sollicitation de l'ensemble des membres du laboratoire
 - Remplir l'inventaire ne constitue pas un engagement à ouvrir ses données
 - Informations de nature documentaire et technique
 - Aucune colonne obligatoire, donc aussi possible de compléter si toutes les informations ne sont pas disponibles
 - Possibilité d'avoir plusieurs lignes pour un même projet/jeu de données (selon le type de données, le périmètre de partage, l'état de finalisation, le lieu de stockage, etc.)
- Question de la granularité pertinente

L'inventaire de données, bilan

- Données surtout pensées en termes de projet (ANR) ou de terrain, non numériques et “passives”
- Certaines colonnes peu remplies: périmètre de partage, licence
- Côté ingénieurs, prise de connaissance de jeux de données
- Côté chercheurs, prise de conscience de l'anticipation nécessaire des questions de gestion et de partage des données

Politique de données vs DMP

Politique de données

- Niveau: laboratoire
- Porte sur ce qui relève de la science ouverte (les données mais pas que!)
- Rubriques axées sur la science ouverte
- Communique une vision stratégique
- Révisable mais vocation à durer dans le temps

Data Management Plan (DMP)

- Niveau: projet
- Porte sur les données
- Rubriques axées sur le cycle de vie des données
- Outil de planification/gestion
- Document évolutif (a minima 3 versions)



Un modèle de politique à adapter

- Proposition du [modèle OpenAIRE](#) (2018) pour les infrastructures de recherche traduit en français et adapté par les animateurs de l'atelier: <https://bit.ly/3x3TXIf>
- Sections
 - Préambule
 - Périmètre d'application de la politique
 - Droits, obligations et responsabilités (laboratoire et membres)
 - Accès ouvert aux publications
 - Accès ouvert aux données de la recherche
 - Science ouverte
 - Infrastructure et services d'appui
 - Evaluation des travaux de recherche et des chercheurs
 - Formation
 - Application de la politique
 - Annexes

La politique de données en pratique

- Des questions
 - sur les tutelles (ont-elles des politiques? quels sont leurs engagements? sous quelle forme?...)
 - sur l'environnement (y a-t-il des services dédiés dans l'environnement du laboratoire (MSH, université, etc.)?)
 - sur l'unité (les questions de science ouverte font-elles partie du plan quinquennal? comment sont gérées les données en interne?)
 - sur les pratiques (départ à la retraite, changement d'unité: que deviennent les données? Qui prend les décisions relatives aux données? Quand?)
- Peu de modifications mais ajout d'informations concrètes (personne ou service ressource sur telle question)
- Choix de ne pas inclure de mesures contraignantes

Bilan

- Outil le plus utile actuellement (et le plus simple à mettre en oeuvre):
l'inventaire de données
 - Réflexions sur sa poursuite et sa mise à jour régulière
- Un point d'entrée particulièrement pertinent pour les anthropologues: la mission