La gestion des données



Cette fiche générale présente les six étapes principales du cycle de vie des données de recherche, chacune détaillée dans une fiche dédiée. En complément, deux fiches transversales, sur les principes FAIR et les aspects juridiques, éthiques et d'intégrité scientifique doivent être prises en compte à chaque étape.

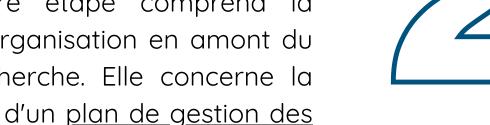
Avant le projet



Planification initiale

Cette première étape comprend la réflexion et l'organisation en amont du projet de recherche. Elle concerne la mise en place d'un <u>plan de gestion des</u> données (PGD).

TESTING



Au cours du projet



Description des données collectées/créées

Cette phase concerne l'origine des données, à savoir la méthodologie de collecte utilisée et la description des données en privilégiant les systèmes utilisant standards des de métadonnées.



Au cours du projet



Validation, Traitement et Stockage

Étape clé du projet de recherche, elle comprend la vérification, le nettoyage, scientifique validation traitement des données. Elle inclut également leur gestion via un stockage sécurisé pour assurer leur intégrité et accessibilité



Après le projet



Accès et Partage

Cette étape consiste à publier des articles de descriptions des données (data papers); à appliquer une licence d'utilisation appropriée, permettant l'exploitation des données par des tiers; et à choisir un entrepôt adapté pour le dépôt, l'accès et la réutilisation des données.



Après le projet



Préservation

Cette démarche consiste à archiver les données à long terme pour conserver les documents, les rendre accessibles et en préserver l'intelligibilité, avec une perspective allant de 10 ans à 30 ans.



Réutilisation

Cette étape permet la création de nouvelles recherches, le réexamen des résultats et des données, et leur pour l'enseignement et utilisation l'apprentissage, tout en respectant les termes de la licence d'utilisation.











Le plan de gestion des données (PGD) est un outil clé pour accompagner la gestion des données de recherche et tracer les différentes étapes de leur cycle de vie. Requis par de nombreuses agences de financement, telles que l'ANR ou Horizon Europe, il permet de répondre aux exigences en matière de gestion et de partage des données. Structuré en rubriques, il doit être mis à jour tout au long du projet. Il permet de décrire et de suivre l'évolution des jeux de données, tout en préparant leur partage, leur réutilisation et leur pérennisation. Cette fiche présente les principales rubriques d'un PGD.

INFORMATIONS ADMINISTRATIVES

Les informations administratives comprennent le nom, l'identifiant et la description du projet, les agences de financement, le responsable principal, le contact pour les données, les dates de première version et de dernière mise à jour, ainsi que les politiques associées.



DESCRIPTION DES DONNÉES

La description des données inclut le type, le format et le volume des données, les jeux de données utilisés, les méthodes de collecte et de création, le système d'organisation et de gestion des fichiers, et les processus d'assurance qualité.



DOCUMENTATION ET MÉTADONNÉES

La documentation et les métadonnées comprennent les informations nécessaires pour lire et interpréter les données, l'organisation de la collecte et de la documentation, ainsi que les standards de métadonnées adoptés.



STOCKAGE, SAUVEGARDE, SÉCURITÉ

Cette étape inclue le lieu de stockage des données, le plan de sauvegarde, les responsables de la sauvegarde, les procédures de récupération, les risques et leur gestion, les dispositifs d'accès, et les mesures pour le transfert sécurisé des données collectées sur le terrain.



Des mise précis peu

CALENDRIER

OBLIGATION LÉGALE

En France, le décret n° 2021-1572 du 3

décembre 2021 (article 6) relatif au

respect des exigences de l'intégrité

scientifique impose la rédaction d'un plan de gestion des données à tous les chercheurs des établissements publics

LA RÉALISATION DU

PGD EN LIGNE

<u>DMP OPIDOR</u> est un **outil de** rédaction d'un PGD en ligne mis à la disposition de la communauté de

l'Enseignement Supérieur et de la Recherche.

Le PGD est un **document évolutif**. Des mises à jour et des livrables précis peuvent être définis selon les financeurs et/ou projets.

SÉLECTION ET CONSERVATION

La sélection et la conservation précisent les détails sur les données retenues, partagées ou conservées, leurs utilisations prévues pour la recherche, la durée de conservation au-delà du projet, le système d'archivage électronique pour l'archivage à long terme, ainsi que les responsabilités, la préparation des données.



DIFFUSION DES DONNÉES

La diffusion des données concerne les étapes pour faciliter leur découverte, les conditions de restriction du partage des données via les licences de réutilisation, le mécanisme de partage, le délai de publication, ainsi que la procédure pour obtenir un identifiant persistant pour les données.



ASPECTS JURIDIQUES ET ÉTHIQUES

L'éthique et le cadre légal précisent l'accord de conservation et de partage des donnée personnelles, la protection de l'identité des participants, la sécurité des données sensibles, les droits de propriété intellectuelle, les propriétaires des données, les licences de réutilisation et les restrictions d'utilisation.



RESPONSABILITÉS ET MOYENS

Cette rubrique inclue le nom de la personne chargée de mettre en œuvre le PGD, les responsables de chaque activité de gestion des données, les équipements et logiciels nécessaires, les besoins en expertise ou formation supplémentaires.



LE SAVIEZ-VOUS ?

Un PGD peut être établi
aussi bien dans une optique
de partage des données que
pour des données en accès
restreint ou fermé, total
ou partiel. Le PGD
mentionnera dans ce cas les

raisons de non partage.

ORIENTÉ LIBRE ACCÈS

Le PGD est très lié au principe du libre accès aux données de recherche. En fonction de votre choix et de vos contraintes en matière de partage, des critères sont à définir.







Dans votre Plan de Gestion des Données (PGD), il est essentiel de préciser l'origine des données, qu'elles soient collectées, créées ou réutilisées. Pour chaque type de données, il est important de détailler la méthodologie de collecte et de les décrire précisément, afin de garantir une gestion des données rigoureuse et conforme aux principes FAIR, facilitant ainsi leur compréhension, leur partage et leur réutilisation par d'autres chercheurs à plus long terme. Cette fiche présente les éléments essentiels à renseigner lors de l'étape de la description des données dans le cadre d'un PGD.



DESCRIPTION DES DONNÉES

Chaque PGD doit comporter une description des données de recherche du projet. Il convient de décrire le plus précisément possible ces données pour qu'elles soient comprises et exploitées.

Provenance

S'agit-il de donnée d'observation, expérimentales, computationnelles, dérivées ou compilées ?

Type

S'agit-il de données textuelles, numériques, audiovisuelles, données spécifiques ?

Stabilité et criticité

Les données sont elles fixes, croissantes, révisables ? Les données sont elles sensibles ?

PROVENANCE

Données administratives

Collectées dans le cadre d'activités institutionnelles courantes (ex: enregistrement d'état civil)

Données d'observation

Capturées en temps réel (ex : données d'enquête)

Données expérimentales

Obtenues en laboratoire, reproductibles mais coûteuses (ex: session expérimentale)

Données dérivées /compilés

Résultant du traitement de données brutes (ex. : indicateurs, modèles économétriques).

TYPE

Données textuelles

Note de terrain, réponses d'enquêtes

Données numériques

Tableaux, mesures ...

Données audiovisuelles

Images, vidéos

Données spécifiques

Liées à une discipline ou un instrument.

STABILITÉ ET CRITICITÉ

Fixe

Données immuables après leur collecte

Croissantes

Nouvelles données ajoutées sans modification des anciennes

Révisables

Nouvelles données ajoutées avec possibilité de modification des anciennes.

Criticité

Les données sont-elles sensibles, confidentielles ? (c.f. fiche "Aspect juridiques, éthiques et intégrité scientifique")



CONTEXTE DE PRODUCTION

Collecte

S'agit-il de collecter des données existantes ? Dans ce cas, il faudra renseigner le chemin d'accès à ces données, c'est-à-dire citer la source et les modalités d'accès.

Création

Pour des données nouvelles, créées ou générées dans le cadre du projet, on détaillera le processus de création ou le mode opératoire mis en place.

Réutilisation

Allez-vous réutiliser des données déjà existantes ?



FORMAT DES DONNÉES

Pour anticiper l'étape de l'archivage des données, le choix du format ouvert est crucial afin que la lecture des données reste accessibles dans le temps.

Formats ouverts

Fichiers non-propriétaires, ils sont encodés de façon transparente, font partie du domaine public. Garantit l'accessibilité et la pérennité des données. (ex : CSV, TXT, PDF/A)

Formats fermés

Fichiers propriétaires, les formats fermés n'appartiennent pas au domaine public. Ils requièrent l'utilisation de logiciels adéquats pour leur lecture et modification. (ex : DOCX, XLSX)



UTILISATION DES STANDARDS DE MÉTADONNÉES

Les métadonnées sont des informations qui décrivent les caractéristiques essentielles d'un jeu de données, telles que son contenu, sa provenance, son format, son producteur, etc.

Cette méthode est une approche spécifique relative à la description des données, reconnue, normalisée, et largement utilisée.

Exemple de Standards

Dublin Core : Utilisé pour décrire divers types de ressources numériques.

METS: Utilisé pour l'encodage de métadonnées dans des documents XML.



UTILISATION D'UN README

Un <u>README</u> (LISEZMOI) peut s'avérer nécessaire pour décrire plus en détails le contexte de production et/ou les données dans les fichiers.

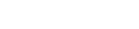
Il vient en complément des métadonnées saisies lors du dépôt du jeu de données, du dictionnaire de données (qu'il peut aussi contenir) et/ou d'autres supports de

Un README est généralement diffusé dans un format ouvert, largement utilisé tel que le texte plein (txt) ou en markdown (md).

documentation accessibles.

CONTACT PROJET

SOURCE



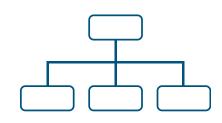


Validation, Traitement et Stockage des données



Cette fiche aborde les étapes concernant la validation, le traitement, et le stockage des données dans un Plan de Gestion des Données (PGD). La validation garantit l'exactitude et la fiabilité des données, tandis que le traitement les transforme en informations exploitables. Enfin, le stockage sécurisé des données garantit leur accessibilité et leur sécurité, en respectant les bonnes pratiques et les normes en vigueur.









Nommage des fichiers

- Donner un nom bref et explicite
- Ne pas mettre d'espace ni de caractères spéciaux
- Indiquer les dates selon la norme ISO 8601 : AAAA-MM-JJ
- Placer l'élément clé en premier pour une identification rapide
- Indiquer les versions des fichiers (ex VP: version provisoire, VD: version définitive ...)

Arborescence des fichiers

- Adopter une organisation thématique (ou sujets)
- Hiérarchiser les répertoires
- Limiter le niveau hiérarchique à 4 ou 5
- Éviter les dossiers "Divers"
- Utiliser des intitulés de dossiers clairs et intelligibles

Prétraitement des données

Vérifier les erreurs, dans les données présentes comme les erreurs de frappe, les valeurs incohérentes, les informations manquantes, les doublons, les imprécisions, etc. On parle de "preprocessing" ou prétraitement des données.







Stockage des données

Au LEM, le stockage institutionnel recommandé pour les données est <u>s-Drive</u>, service officiel du CNRS.

- Accessible à tous les membres du LEM via leur compte Janus (chaque agent CNRS/UMR en dispose automatiquement).
- Synchronisation fluide via un navigateur web ou via NextCloud (type Dropbox/Google Drive) .
- Espace personnel de 100 Go, non extensible.
- Corbeille intégrée : fichiers supprimés conservés jusqu'à 30 jours.
- Hébergement sécurisé en France, avec analyse antivirus et sauvegardes hebdomadaires.
- Possibilité de créer des espaces partagés (workspaces) pour les projets collectifs, avec quota dédié.

Traitement et Analyse des données

L'analyse peut ensuite être réalisée à l'aide d'outils comme la visualisation ou le machine learning. Il est important de conserver une traçabilité des modifications, en particulier pour les données sensibles. Une forge (espace collaboratif) permet de suivre les différentes versions tout en facilitant la collaboration. Il est conseillé d'utiliser une forge sauveraine. Par exemple l'université de Lille a une forge institutionnelle <u>GitLab</u> accessible via l'ent.

Disponible sur : https://doranum.fr/stockage-archivage/stockage-donnees_10_13143_z0ge-nc29/







Cette fiche présente les moyens essentiels pour garantir un accès et une diffusion efficaces des données de recherche dans un Plan de Gestion des Données (PGD). Cette étape repose sur trois leviers principaux : publier des articles décrivant les données (data papers) pour en accroître la visibilité et l'impact scientifique, appliquer une licence d'utilisation adaptée pour encadrer leur exploitation par des tiers, et choisir un entrepôt de données fiable pour assurer leur stockage, leur accessibilité et leur réutilisation dans une démarche de science ouverte.

POURQUOI PUBLIER UN DATA PAPER?

- Pour communiquer sur l'existence des données et permettre de les trouver.
- Pour créditer les auteurs (reconnaissance, référence citable) et valoriser les données.
- Pour faciliter la réutilisation des données et le travail (en les rendant intelligibles).

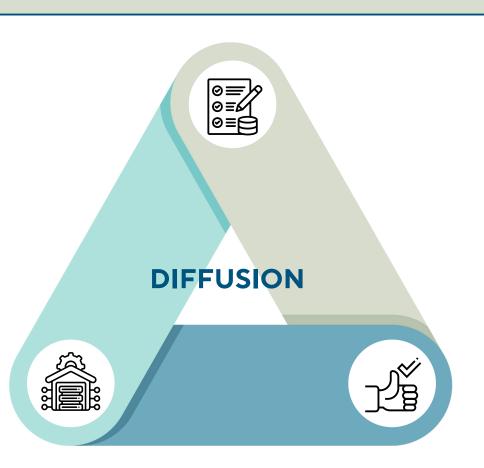
Il est possible de générer une ébauche d'un datapaper à partir d'un jeu de donnée déposé dans l'entrepôt Recherche Data Gouv. Cliquez ici pour plus <u>d'informations</u>

QUEL ENTREPÔT?

Les données de recherche doivent être déposées sur la plateforme recherche data.gouv.fr, dans l'<u>espace</u> dédié à la collection au laboratoire Lille Économie Management (LEM) de l'entrepôt de données de l'Université de Lille, Lillodata.

DATA PAPERS

Publication qui décrit un jeu de données scientifiques, notamment à l'aide d'informations structurées, appelées métadonnées. Le <u>data paper</u> fournit une voie formalisée au partage des données contrairement à l'article de recherche classique qui teste des hypothèses ou présente de nouvelles analyses. Il suit également un processus de publication avec peer review.



OÙ PUBLIER UN DATA PAPER?

- Dans un data journal, revue dédiée à ce type de publication.
- Dans une revue classique qui publie des data papers en plus des articles traditionnels.

QUE CONTIENT UN **DATA PAPER**

Les data papers ont une structure très variable selon les revues. Ils ont cependant des composantes communes.

ENTREPÔT

Pour déposer les données de recherche scientifiquement validées (sur le court et moyen terme - 5 à 10 ans), et les rendre accessibles et réutilisables. Il en existe plusieurs types : disciplinaires, multidisciplinaires, propres à un éditeur, institutionnelles, spécifiques d'un projet de recherche.

AVANT DE DÉPOSER SES DONNÉES

- Les jeux de données à partager ont été sélectionnés
- Les principes éthiques sont respectés
- Les droits de diffusion sont vérifiés
- Les modalités d'accès sont définies
- Les fichiers sont organisés et nommés de manière explicite
- Les fichiers sont dans des formats pérennes et ouverts
- Les données sont décrites et documentées
- Une licence est attribuée aux données

LICENCE D'UTILISATION

Contrat par lequel le ou les producteur(s) autorise(nt) l'exploitation par un tiers d'utiliser les données et précise(nt) dans quels buts elles peuvent être utilisées, à l'aide de conditions d'accès et de réutilisation. Pour favoriser la réutilisation des données, il faut privilégier les licences ouvertes.

QUELLES LICENCES?

Pour encadrer la réutilisation des données de recherche, vous pouvez choisir une licence parmi celles proposées sur <u>data.gouv.fr</u>, comme la Licence Ouverte V2 Etalab (permettant une réutilisation libre et gratuite des données, à condition de mentionner la source).









Différente du stockage des données, l'étape de préservation, ou archivage pérenne, a pour objectifs principaux de conserver les données sur le long terme et d'en garantir l'intelligibilité dans le futur. Contrairement au stockage, qui se concentre sur la continuité de l'exploitation durant le projet, l'archivage pérenne prévoit une conservation durable (plus de 30 ans) et s'accompagne d'enjeux financiers, environnementaux et scientifiques. Bien qu'importante, cette étape non obligatoire est payante, et nécessite une réflexion approfondie pour évaluer sa pertinence en fonction des objectifs du projet.

PRÉPARATION DES DONNÉES Assurer l'accès et la lisibilité des données à long terme

ARCHIVAGE DES DONNÉES

QUELS PRÉREQUIS?

- Organiser les données : hiérarchiser, nommer, versionner et choisir des formats de fichiers appropriés.
- Conserver de manière pérenne ou être détruit à l'issue de sa durée d'utilité administrative.
- Documenter les données en utilisant des standards de métadonnées pour assurer leur traçabilité et leur réutilisation.
- Estimer la volumétrie des données pour anticiper les besoins en stockage.
- Planifier des sauvegardes régulières pour éviter les pertes
- Prévoir le financement nécessaire pour la gestion et le stockage des données en fonction de leur volumétrie.

CLÉS DE L'ARCHIVAGE PÉRENNE

L'archivage pérenne repose sur trois axes principaux pour assurer la conservation des données à long terme :

- L'intégrité des supports de stockage : garantir la durabilité et la sécurité des supports à long terme.
- La lisibilité des formats de fichiers : préserver la compatibilité des formats afin d'éviter les risques de dépréciation technologique.
- L'intelligibilité des données : maintenir la compréhension et l'utilisation des données, grâce à une documentation appropriée (métadonnées, guides d'utilisation, etc.).

COMMENT ARCHIVER SES DONNÉES

L'archivage des données de recherche repose sur des services d'archivage électronique dédiés (SAE). En France, l'acteur principal est le CINES, dont la mission est d'archiver les données et documents numériques produits par la communauté académique et de recherche via cette plateforme. Il propose des solutions d'archivage numérique payantes, à moyen et long terme, tout en offrant une expertise en informatique et en archivistique. La sécurité et l'intégrité des données sont assurées grâce à diverses procédures, telles que l'attribution de métadonnées, le choix de formats de fichiers pérennes, la réplication des données et la mise à disposition d'un environnement informatique sécurisé.



Les principes FAIR



Dans votre Plan de Gestion des Données (PGD), il est essentiel de garder à l'esprit les principes FAIR, car ceux-ci interviennent à chaque étape du cycle de vie des données. Les principes FAIR (Findable, Accessible, Interoperable, Reusable) sont des recommandations internationales visant à améliorer la gestion et le partage des données de recherche. Ils visent à rendre les données plus faciles à trouver, accessibles, interopérables entre systèmes, et réutilisables par d'autres chercheurs, contribuant ainsi à la transparence et à la collaboration scientifique.

FACILITER LA DÉCOUVERTE DES DONNÉES

- Les données ont un identifiant pérenne ou Persistent IDentifier (PID) en anglais (par exemple le Digital Object Identifier ou DOI) afin de disposer d'un accès stable à la ressource.
- Les données sont décrites par des métadonnées scientifiques et des métadonnées documentaires.
- Les données, ou au moins leurs métadonnées, sont indexées ou enregistrées dans un outil de recherche, i.e. un dépôt de donnée dans un entrepôt ou via un catalogue de données.

PERMETTRE L'ACCÈS AUX DONNÉES ET AUX MÉTADONNÉES

- Sur Internet, à travers un protocole standard,
 libre et ouvert. (par exemple https)
- Sur **authentification** pour les données en a**ccès restreint**.
- Les métadonnées doivent rester accessibles même si les données sont temporairement inaccessibles ou si l'accès aux données est restreint.

LES PRINCIPES FAIR

RENDRE SES DONNÉES INTEROPÉRABLES

- Les données sont **décrites dès le début** du cycle de vie à l'aide de vocabulaires contrôlés.
- Les métadonnées doivent autant que possible faire référence aux autres données qui peuvent être mises en relation et ainsi permettre des liens entre elles.
- Les formats de fichiers utilisés sont **ouverts et documentés** pour permettre l'exploitation et une pérennisation des données par différents outils.

PERMETTRE LA **RÉUTILISATION**DES DONNÉES

- Les métadonnées ont une pluralité d'attributs utiles pour la compréhension et la réutilisation des données.
- Une **licence de réutilisation** est attribuée au données.
- La description des données indique leur provenance.
- La structure des données suit les standards de la communauté scientifique pour faciliter leur analyse.

Aspects juridiques, éthiques, intégrité scientifique



Dans votre Plan de Gestion des Données (PGD), cette fiche sur les aspects juridiques et éthiques doit être prise en compte à chaque étape du cycle de vie des données de recherche. Elle accompagne tout le processus, depuis la collecte jusqu'au partage, en intégrant des principes tels que la protection des données sensibles, le partage des données et les droits et devoirs du chercheur. Suivre ces recommandations garantit une gestion rigoureuse et conforme, essentielle pour une recherche éthique et responsable.

LES LICENCES À METTRE EN OEUVRE

La licence choisie par l'auteur définit ce que le réutilisateur est autorisé ou non à faire avec les données. Il faudra a minima qu'il respecte l'intégrité des données et qu'il mentionne la paternité de l'information (sa source et la date de dernière mise à jour)

L'IÉSEG étant un établissement privé, les chercheurs affiliés à cette tutelle ne sont pas soumis aux mêmes obligations d'ouverture des données que les chercheurs des établissements publics.

Cliquez ici pour plus d'informations

RÉGLE JURIDIQUE POUR LES DONNÉES

Les données relèvent d'un régime lié au droit des bases de données. Dans ce cas, le droit de propriété appartient légalement à l'établissement de tutelle des chercheurs (CNRS, Université de Lille, Université ULCO, IESEG ou Université d'Artois). Il est considéré comme le titulaire effectif du droit de propriété.

L'ÉTHIQUE ET L'INTÉGRITÉ SCIENTIFIQUE

Le respect de la vie privée, la propriété intellectuelle, la qualité et l'intégrité des données sont des dimensions éthiques de la gestion des données. Le code de conduite européen pour l'intégrité en recherche identifie quatre principes fondamentaux: fiabilité, honnêteté, respect et responsabilité

UNE OUVERTURE PAR DÉFAUT

Depuis la loi pour une République numérique de 2016, les données de recherche « achevées », financées au moins pour moitié par des fonds publics, sont assimilées à des données administratives et font donc l'objet d'un « principe d'ouverture par défaut ». Par conséquent, elles sont censées être publiées et rendues accessibles sur Internet, sauf exceptions (données personnelles, sensibles, secret industriel, secret défense, protection du potentiel scientifique et technique (PPST), zones à régime restrictif (ZRR), etc.)

Diffusion des données

Droit et devoir du chercheur

Données à caractères juridiques particulier

QUELLES DONNÉES?

D'après la CNIL, tous les types de données sont concernées, à savoir les données textuelles, numériques, audiovisuelles, spécifiques ...

Si vous êtes amené à collecter et traiter des données personnelles dans le cadre du RGPD,

pensez à consulter le DPO de votre tutelle de rattachement, à savoir:

- <u>DPO de l'Université de Lille</u>
- <u>DPO de l'IESEG</u>
- <u>DPO du CNRS</u>
- <u>DPO de l'Université ULCO</u>
- <u>DPO de l'Université d'Artois</u>

LES DONNÉES PERSONNELLES

Les données personnelles sont des informations permettant de vous identifier directement ou indirectement. Il s'agit de votre prénom, nom, âge, genre, numéro de téléphone, adresse IP, vos adresses postale et électronique, vos voix, images, empreintes digitales, votre signature ...

LES DONNÉES SENSIBLES

Elles font partie des données personnelles. Elles révèlent la prétendue origine raciale ou ethnique, les opinions politiques, les convictions religieuses ou philosophiques ΟU l'appartenance syndicale, le traitement des données génétiques, des données biométriques aux fins d'identifier une personne physique de manière unique, les données concernant la santé, des données concernant la vie sexuelle ou l'orientation sexuelle d'une personne phyjque

CONTACT PROJET

