



**Géographie-cités**  
UMR 8504

**Equipe PARIS**



# **Des questions et des enjeux contemporains associés aux données de la recherche en SHS**

Illustration par des travaux menés au sein de l'équipe

Jordi Calabuig Serra, Julie Gravier, Marianne Guérois, Thibault Le Corre,  
Thomas Louail, Hugues Pécout, Pierre Pistre, Olivier Telle,  
Céline Vacchiani-Marcuzzo et Julie Vallée

séminaire d'équipe PARIS // 25 novembre 2022

# Contexte

Global : Numérisation des sociétés

- > augmentation du volume de données de recherche, notamment des données personnelles
- > + grande diversité de producteurs, données massives collectées par entreprises utiles à la recherche

Européen : RGPD depuis 2018, renforcement des droits des personnes, implications pour la recherche

Promulgation de la science ouverte par les institutions scientifiques et mécanismes incitatifs/obligatoires pour adopter une démarche FAIR (*Findable, Accessible, Interoperable, Reusable*)

- > Critère d'évaluation pour financement ; obligation de rédaction d'un PGD dans les projets financés

National : Loi pour une république numérique ("loi Lemaire) depuis 2016, "loi Valter" depuis 2015

- > 2 IR\* : Huma-num et PROGEDO

UMR : Recherches à l'échelle des individus plus nombreuses

- > thèses et projets qui collectent des données personnelles

"Commande" des directions de l'UMR et de l'équipe d'une séance avec retours d'expériences et discussion collective

Fonction(s)

- intro à des sujets techniques et pointeurs pour auto-formation
- évaluation de la demande pour des formations spécifiques
- état des lieux/pratiques utile pour rédaction du bilan et projet

# Plan de la présentation

RGPD et recherches en SHS, démarches, 2 exemples  
(Thomas et Jordi)

IR\* PROGEDO et données de la statistique publique (Pierre)

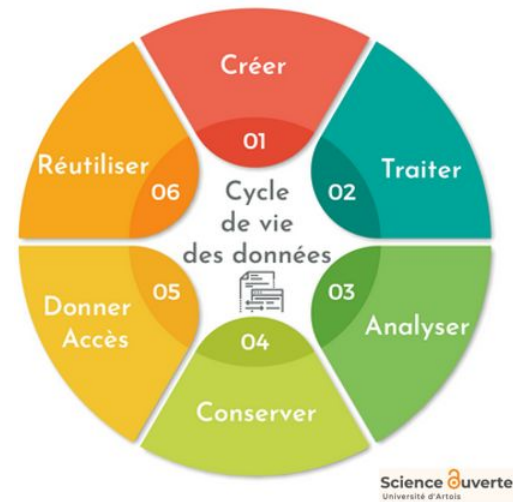
Ouverture de données (Julie V. et Thibault)

Cas de la ré-utilisation de données collectées par des entreprises privées  
(Marianne et Olivier)

Inter-opérabilité des données et ontologies (Julie G.)

Les outils mis à disposition par la IR\* Huma-num (Hugues)

Questions & Discussion



# Traiter des données personnelles/“à caractère personnel”

**Les données individuelles anonymisées de manière irréversible**

– sans possibilité de ré-identification des personnes –

**ne sont pas considérées comme des données personnelles**

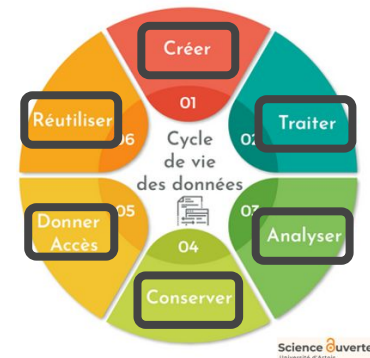
→ pas soumises au RGPD

Données personnelles :

- directement identifiantes (nom, prénom, photo/video/voix, numéro SS, adresse, etc.)
- indirectement identifiantes (recoupement d'informations)

données pseudonymisées : qui ne peuvent plus être directement attribuées à la personne concernée, mais le croisement avec une table d'informations extérieures peut permettre la ré-identification

données sensibles : orientation sexuelle, religieuse, appartenance syndicale, convictions politiques ou philosophiques, données de populations vulnérables, ...



# Notions et rôles définis dans le RGPD

**Traitement** de données personnelle : toute opération sur des données personnelles, quelque soit le support utilisé, informatisé ou non

La notion légale de traitement plus large que l'analyse et l'exploitation des données  
→ inclut aussi la collecte, le stockage, l'archivage, la réutilisation de données

**Délégué.e à la protection des données (DPD, DPO en anglais)** : tout organisme public en a un.e, toute unité de recherche aussi ; pour l'UMR, la DPD du CNRS

**Responsable de traitement** : dans les unités CNRS c'est le/la DU

Le/la **responsable du projet scientifique** s'assure avec la/le DPD de la conformité du traitement de données personnelles

Le **registre** des traitements de l'unité enregistre tous les traitements de données personnelles effectuées dans des recherches à l'UMR

# Droits des personnes et obligations des chercheur.e.s

## Droits des personnes :

- d'information précise sur le traitement
- d'accès à ses données
- de rectification
- d'opposition
- d'être informé d'une violation des données
- à la portabilité
- à une utilisation restreinte de ses données

## Le traitement doit être (6 principes) :

- licite (avoir un fondement légal) :  
consentement des personnes, mission de service public, "intérêt légitime"
- avoir une finalité bien définie
- pertinent et proportionné
- sécurisé
- conservation limitée dans le temps
- transparent

# Le cas des données personnelles collectées dans les thèses réalisées à l'unité

- Collecte de données potentiellement sensibles : orientation sexuelle, identité de genre, populations vulnérables (migrants)
- **La/le doctorant.e avec sa/son directrice de recherche s'assurent de la conformité du traitement avec le RGPD, en dialogue avec la cellule DPD du CNRS**
- Le traitement doit faire l'objet d'un enregistrement dans le registre de l'UMR
- Le fondement légal peut être le consentement, l'intérêt légitime, la mission de service public → à déterminer avec le DPD
- La finalité est le sujet de la thèse
- Les autres principes peuvent être garantis notamment par l'usage des services huma-num (stockage sécurisé, archivage pérenne après anonymisation, cf. présentation d'Hugues)
- Obligation d'information précise des personnes, rédaction d'une mention d'information
- Exemples à suivre :
  - La collecte de données réalisée par Jordi Calabuig Serra dans sa thèse ;
  - L'information aux personnes participant à l'enquête par questionnaire du projet RECORDS

**Responsable du traitement :** Les informations recueillies vous concernant vont faire l'objet d'un traitement de données à caractère personnel par le responsable du traitement, le géographe Jordi Calabuig Serra, sur la population LGBT+ résidente à Barcelone. Une approche géographique est adoptée, cotutelle à l'Université Paris 1 Panthéon Sorbonne (France) et à l'Université de Girona (Espagne). Géographie-cités : <https://geographie-cites.cnrs.fr/membres/jordi-calabuig-serra/>

**Finalités du projet :** Le traitement a pour objet de cerner les espaces et les formes de socialisation des touristes LGBT+ à Barcelone (Espagne) en se focalisant sur les expériences vécues et les perceptions des individus.

**Modalités de collecte :** Il y a dans ce projet trois modalités de collecte de l'information : des entretiens individuels avec des touristes et des résidents, et des observations participantes :

- Un questionnaire en ligne adressé aux touristes LGBT+ qui ont visité Barcelone entre 2018 et 2022.
- Des entretiens individuels avec des touristes et des résidents. Ils seront enregistrés au format audio et durent entre 45 et 90 minutes, la durée moyenne étant d'une heure. Ils se dérouleront de préférence en face à face ou, à défaut, en ligne.
- Les observations participantes ne concernant pas la collecte de données permettant d'identifier les personnes.

**Nature des données collectées :** Seules les données strictement nécessaires à la réalisation de notre recherche seront collectées et traitées. Au cours des entretiens et dans le questionnaire, des sujets tels que l'orientation sexuelle, l'identité de genre, les lieux de loisirs et de socialisation, le lieu de résidence approximatif et la mobilité (et les modalités de socialisation) seront abordés. Pour les personnes prêtes à être interviewées, un moyen de contact sera communiqué.

**Base légale du traitement :** La base légale du traitement repose sur la mission d'intérêt public de la recherche. Si vous êtes français·e, vous devez être âgé·e de quinze ans ou plus. Si vous êtes d'un autre pays, vous devez répondre à des questionnaires sans le consentement de votre représentant de l'autorité parentale.

**Participation libre :** Votre participation au projet de thèse « Effets du tourisme LGBT+ sur la population LGBT+ résidente à Barcelone. Une approche géographique des interactions » est entièrement libre et volontaire.

**Retrait du consentement :** Vous êtes libre de vous retirer ou de cesser votre participation à tout moment sans conséquence pour vous.

**Pseudonymisation et confidentialité :** Le responsable du projet de thèse « Effets du tourisme LGBT+ sur la population LGBT+ résidente à Barcelone. Une approche géographique des interactions » prend les engagements suivants :

- Si vous participez à un entretien, votre identité sera dissimulée à l'aide d'un faux prénom dans les publications académiques mais aussi les comptes rendus de communications académiques, etc. Aucune autre information ne sera conservée qui puisse révéler votre identité : les notes d'entretien, comptes rendus d'entretiens, notes d'observation, notes d'analyses et publications seront complètement anonymisées.
- Si vous répondez au questionnaire en ligne, votre identité sera dissimulée à l'aide d'un numéro aléatoire pour tous les renseignements collectés.
- Seulement le responsable de projet détient la table de correspondance qui permet de faire le lien entre votre identité et le faux prénom ou le numéro aléatoire attribué dans les différents fichiers.

**Destinataires des données personnelles :** Le destinataire de ces données est le responsable du traitement, NADINE CATTAN et LUIS PRATS PLANAGUMA, ayant accès à l'information complète mais de manière anonyme.

**Transferts de données :** Toutes les données seront gardées en France.

**Durée de conservation :** Vos données personnelles sont conservées en base active jusqu'au 31/12/2023. Après cette date, elles seront définitivement archivées de manière anonyme.

**Mesure de sécurité :** Afin de garantir la confidentialité de vos données et éviter leur divulgation, le responsable du projet est autorisé à accéder aux données.

- Les mesures de sécurité, tant physiques que logiques, sont prises et conformes à la politique de confidentialité de l'Université de Girona.

**Diffusion :** Les résultats de cette recherche seront diffusés sans que l'identification des personnes enquêtées soit possible. Dans un document public, les résultats seront diffusés dans des colloques professionnels et scientifiques, des revues professionnelles et académiques et dans des médias destinés au grand public.

**Droits des personnes :**

- Vous pouvez poser des questions au sujet de ce projet à tout moment en communiquant avec le responsable du traitement à l'adresse [Serra@etu.univ-paris.fr](mailto:Serra@etu.univ-paris.fr).
- Vous pouvez accéder et obtenir copie des données vous concernant, vous opposer au traitement de vos données, ou demander la suppression de vos données. Vous disposez également d'un droit à la limitation du traitement de vos données. Vous pouvez exercer ces droits en écrivant à [Serra@etu.univ-paris.fr](mailto:Serra@etu.univ-paris.fr).
- Vous pouvez contacter la Déléguée à la Protection des Données du CNRS : CNRS, Service de la Protection des Données, 17 Avenue de la République, 91190 Vandœuvre-lès-Nancy, mail : [dpd.demandes@cnrs.fr](mailto:dpd.demandes@cnrs.fr).
- Si après avoir contacté le responsable du projet la DPD du CNRS vous estimez que vos données ne sont pas traitées conformément à la loi, vous avez la possibilité d'introduire une réclamation en ligne auprès de la Commission nationale de l'Informatique et des Libertés (CNIL), 3 Place de Fontenay, TSA 80716, 75334 Paris Cedex 07, <https://www.cnil.fr/>.

## Un exemple : Retour d'expérience RGPD et enquête

1/3

### Entretiens – Observations sur le terrain – Questionnaire en ligne

**Pseudonymisation et confidentialité :** Si vous participez à un entretien, votre identité sera dissimulée à l'aide d'un faux prénom [...]

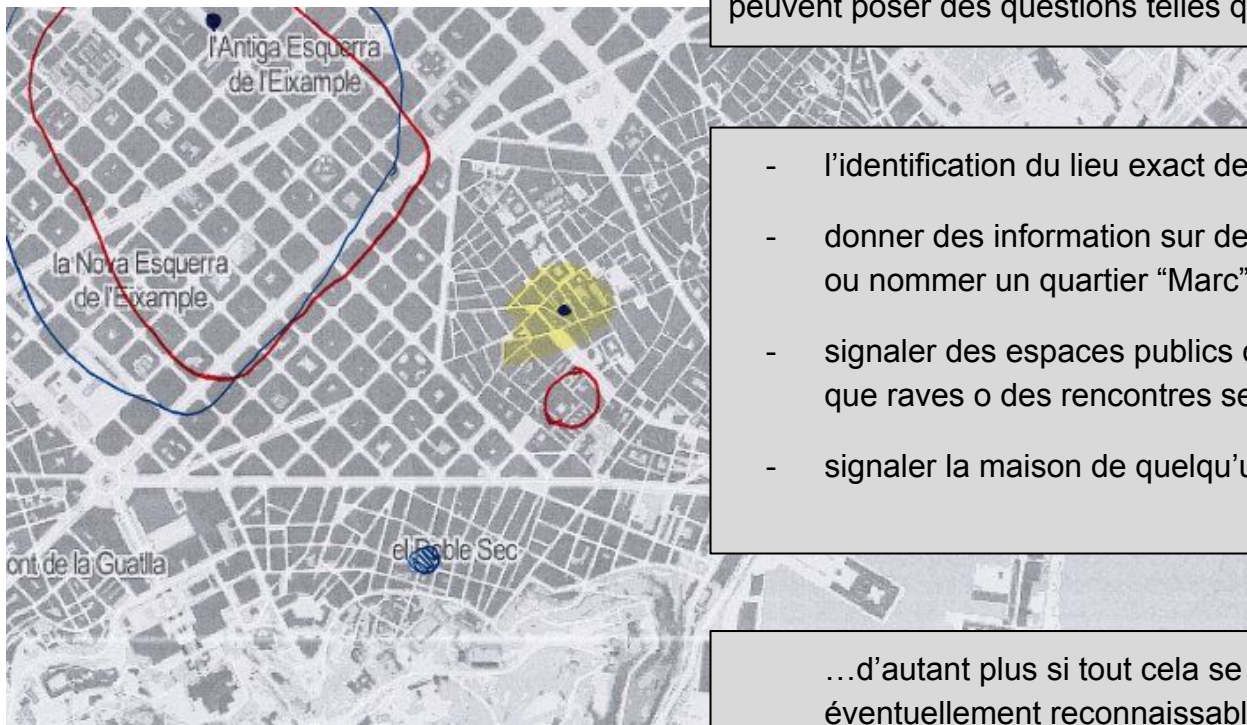
**Diffusion :** Les résultats de cette recherche seront diffusés sans que l'identification des personnes enquêtées soit possible.

**Base légale du traitement :** Si vous êtes français·e, vous devez être âgé·e de quinze ans ou plus. Si vous êtes d'une autre nationalité, vous devez avoir l'âge minimal requis pour répondre à des questionnaires sans le consentement de votre représentant de l'autorité parentale.

**Durée de conservation :** Vos données personnelles sont conservées en base active jusqu'à un an après la soutenance de la thèse de Jordi Calabuig Serra, prévue en 2023. Après cette date, elles seront définitivement archivées de manière anonymisée.



**Modalités de collecte** : Il y a dans ce projet trois modalités de collecte de l'information [...] parmi lesquelles des entretiens avec une carte de la ville qui peuvent poser des questions telles que :



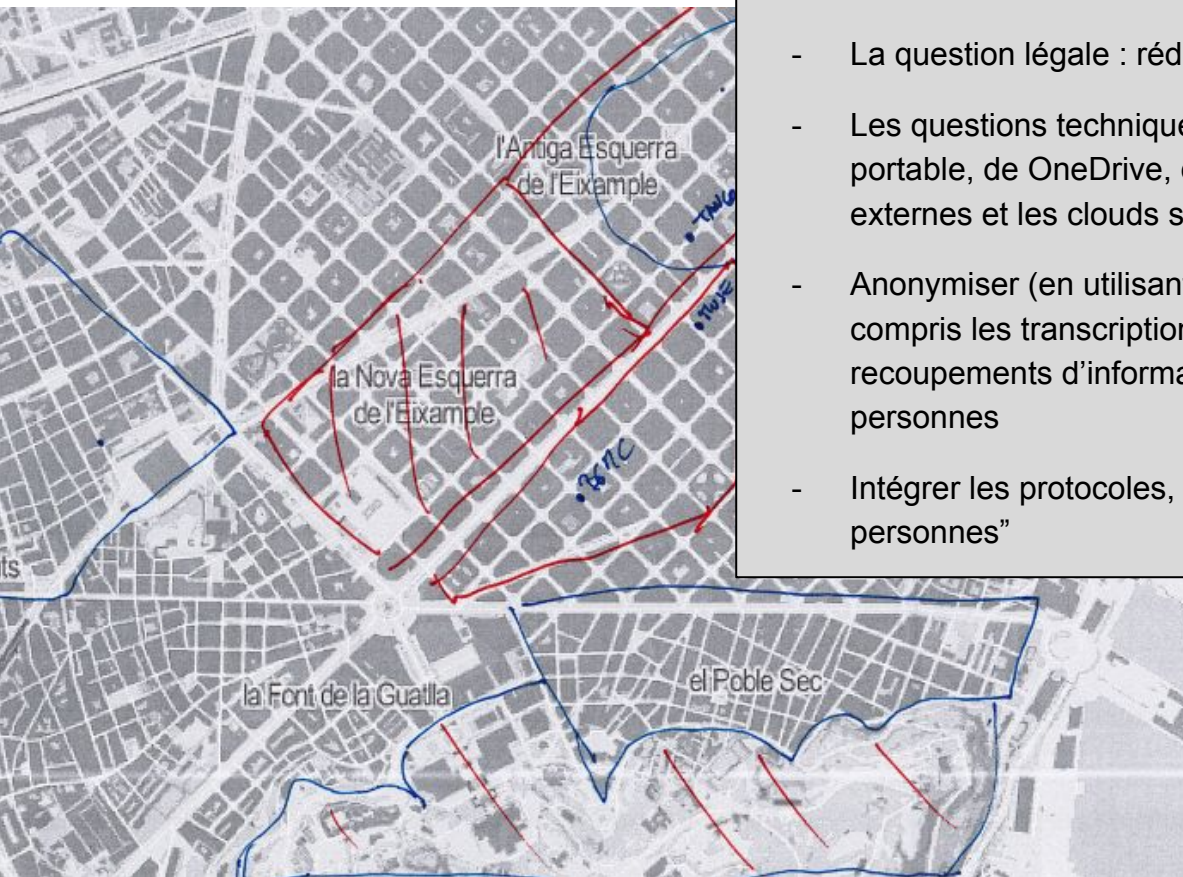
- l'identification du lieu exact de travail ou de résidence
- donner des information sur des tiers (p. ex. identifier "la maison de Júlia" ou nommer un quartier "Marc" d'après le prénom de son petit ami)
- signaler des espaces publics où se déroulent des activités illégales, telles que raves o des rencontres sexuelles en plein air
- signaler la maison de quelqu'un où l'on pratique le chemsex

...d'autant plus si tout cela se fait par écrit, donc, avec une écriture éventuellement reconnaissable

**Résumé des points principaux :**

3/3

- La question légale : rédaction et publication de la notice d'information
- Les questions techniques : suppression de fichiers de mon téléphone portable, de OneDrive, de Gmail et similaires. Privilégier les sauvegardes externes et les clouds sûrs
- Anonymiser (en utilisant les noms d'emprunt dans tous les fichiers, y compris les transcriptions d'entretiens) et être vigilant aux éventuels recoupements d'informations qui peuvent rendre identifiables les personnes
- Intégrer les protocoles, notamment en ce qui concerne les "Droits des personnes"



# Autre exemple : projet collaboratif mené à l'UMR

## Mention d'information d'une enquête par questionnaire

Thomas L.

### Responsable(s) du traitement

Les informations recueillies dans le questionnaire en ligne font l'objet d'un traitement dans le cadre du programme de recherches *RECORDS* (pratiques d'Es publics des plateformes de Streaming). Le coordinateur scientifique de ce projet est M. Thomas Louail, chargé de recherches au CNRS et membre de l'unité Géographie-cités, dont les locaux sont situés sur le Campus Condorcet, Bâtiment de recherche Sud, 93322 Aubervilliers Cedex. Les autres responsables scientifiques du programme RECORDS sont : (i) pour Sciences Po, M. Philippe Coulangeon, directeur de recherches au CNRS et membre de l'unité Observatoire Sociologique du Changement (ii) pour le Centre Marc Bloch, M. Camille Roth, chargé de recherches au CNRS (iii) pour Orange, Mme Valérie Peugeot, chercheuse à Orange Labs (iv) pour Deezer, M. Manuel Moussallam, chercheur au sein du département R&D de Deezer.

Vous pouvez les contacter par mail à l'adresse [records@parisgeo.cnrs.fr](mailto:records@parisgeo.cnrs.fr)

Extraits de la page :

<https://records.huma-num.fr/rqpd/>

### Finalités du projet

L'objet du traitement des données personnelles collectées est la réalisation d'études scientifiques :

- sur les consommations de contenus musicaux (« streams »), les goûts et les pratiques culturelles des utilisatrices et utilisateurs de la plateforme de streaming musical Deezer
- sur leurs usages de cette plateforme, les contextes dans lesquels elles/ils l'utilisent, et comment elles/ils l'associent à d'autres supports d'écoute

Nous attendons de vous que vous participiez à une enquête par questionnaire dans laquelle nous vous poserons des questions sur votre usage de Deezer, vos goûts musicaux, vos habitudes de sortie et de divertissement, ainsi que sur votre situation sociale et économique.

10 à 15 minutes sont nécessaires en moyenne pour répondre intégralement au questionnaire.

### Pseudonymisation et confidentialité

Votre identité sera dissimulée à l'aide d'un numéro aléatoire pour tous les types d'informations collectées.

Seul.e.s les responsables scientifiques de RECORDS détiennent la table de correspondance qui permet de faire le lien entre votre identifiant sur Deezer (c'est-à-dire l'adresse mail avec laquelle vous vous connectez à Deezer) et le numéro aléatoire qui vous sera attribué dans les différents fichiers manipulés par les chercheur.e.s et ingénieur.e.s participant au programme RECORDS. Cette table de correspondance ne sera jamais transmise à des partenaires hors du programme RECORDS.

### Destinataires des données personnelles

Après pseudonymisation des données personnelles collectées dans le questionnaire RECORDS, les destinataires de ces données sont exclusivement les chercheur.e.s et ingénieur.e.s membres des établissements partenaires du programme de recherche RECORDS, dont la liste est disponible sur la page : <https://records.huma-num.fr/participant-e-s/>

### Mesure de sécurité

Afin de garantir la confidentialité de vos données et éviter leur divulgation, toutes les mesures sont prises, notamment :

- seuls les chercheur.e.s impliqué.e.s dans le projet sont autorisés à accéder aux données collectées ;
- les mesures de sécurité, tant physique que logique, sont prises ;
- les données sont hébergées sur un serveur sécurisé du CNRS localisé

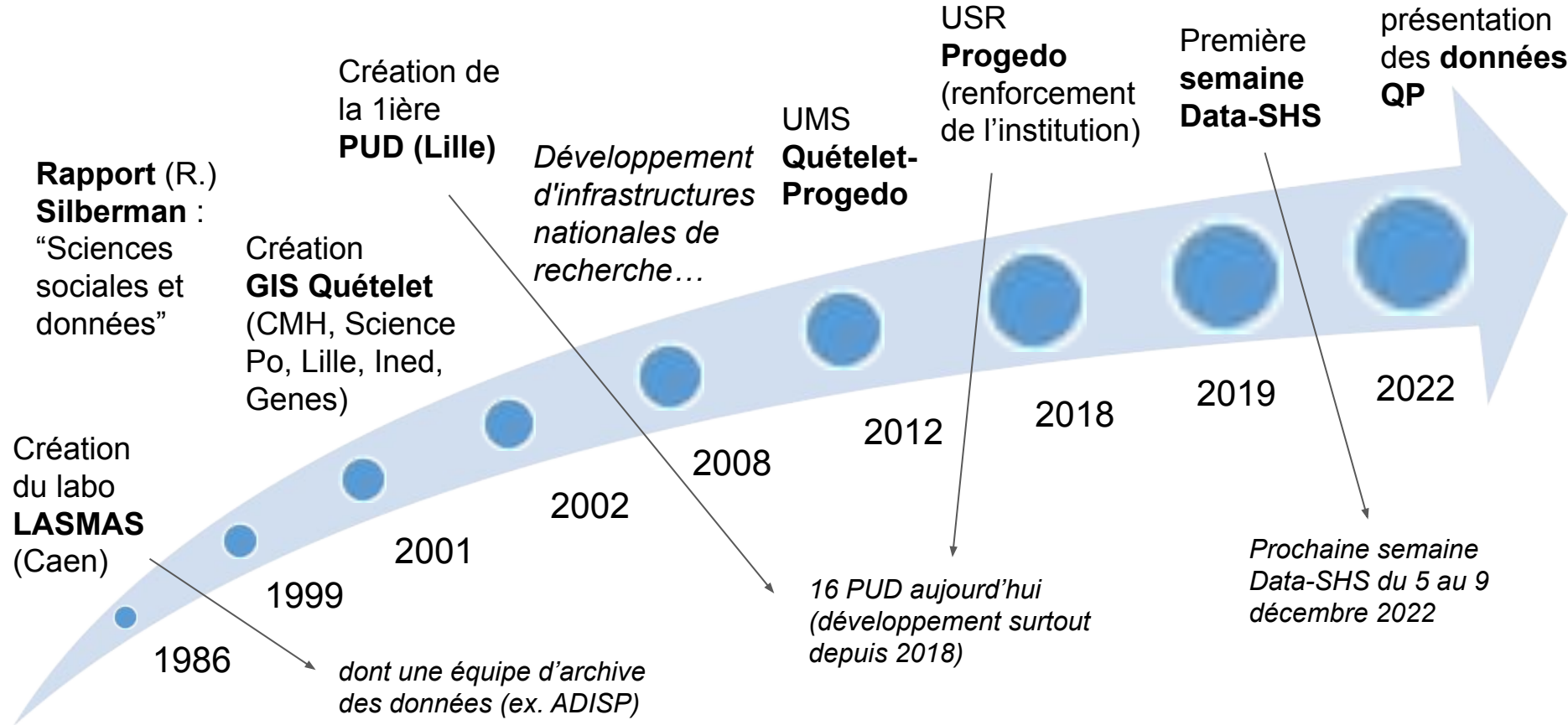




IR\* PROGEDO,  
une grande infrastructure  
de recherche en SHS  
pour traiter, analyser,  
conserver, donner accès  
aux données de la  
statistique publique



# PROGEDO, un acteur "historique" de l'accès aux données de la statistique publique pour les SHS



# PROGEDO, un riche catalogue de données de la statistique publique pour les recherches en SHS



Producteurs : surtout **INSEE, services ministériels** (DARES, DREES...), **INED, Cerema...**

<u>Données</u> : des <b>séries</b> de sources	=	<b>97</b>
(ex. recensement de la population, enquête emploi, enquêtes Mobilités et déplacements...)		
↳ Des <b>jeux de données</b> (ex. enquête emploi 2021)	=	<b>1 540</b>
↳ Des <b>variables</b> (ex. AAH dans l'enquête emploi 2021)	=	<b>395 000</b>

Conditions d'accès et types de commande :

- Fichiers standards anonymes (cf. enjeux de conservation), hors RGPD = service gratuit
- **Fichiers de Production et de Recherche pseudonymisés (FPR)**, RGPD = service gratuit
- *PSM (Produit sur Mesure) pour certaines sources INSEE*, RGPD = service gratuit
- **Données sensibles / fichiers confidentiels** : accès à distance sécurisé (cf. **CASD**), RGPD = service payant

# PROGEDO, un **nouveau portail** de présentation des données

(<https://data.progedo.fr/>) ...

(1/3)

- Une **interface plus interactive** (moteur de recherche, filtres, rubriques de metadonnées harmonisées : Résumé, Détails, Variables, Données, Documentation, Accès...)
- Une **commande facilitée** (cf. onglet “Panier” pour chaque jeu de données)
- Un **meilleur référencement des enquêtes** (cf. un DOI unique pour chaque enquête pour faciliter les citations)

Logo: CNRS ined | a member of CESSDA | Pierre P.

CATALOGUE DE DONNÉES EN SCIENCES HUMAINES ET SOCIALES POUR LA RECHERCHE

1 543 Jeux de données | 396 514 Variables | 13 133 Utilisateurs

niveau de vie...

Séries (93) | Jeux de données (1543) | Variables (396514)

**Capacité des communes en hébergement touristique - 2022**

Thème(s) : Données localisées  
Version 1 Date de la dernière version : 2023-02-28  
Producteur(s) : INSEE Diffuseur : Progedo-Adison

**Table d'appartenance géographique des communes - 2022**

Thème(s) : Données localisées  
Version 1 Date de la dernière version : 2023-02-18  
Producteur(s) : INSEE Diffuseur : Progedo-Adison

**Table de passage annuelle - 2022**

Thème(s) : Données localisées  
Version 1 Date de la dernière version : 2023-02-18  
Producteur(s) : INSEE Diffuseur : Progedo-Adison

**Baromètre d'opinion de la DREES - 2021**

Thème(s) : Opinions — Santé et protection sociale  
Version 1 Date de la dernière version : 2022-07-19  
Producteur(s) : INSEE Diffuseur : Progedo-Adison

# PROGEDO, un nouveau portail de présentation des données et de nouvelles opportunités (2/3)

## Variables :

- Des informations détaillées
- Des explorations facilitées

←

résidence secondaire

1950 2022

**Années**

- 2020 9
- 2019 9
- 2018 9
- 2017 54
- 2016 8

— Voir plus —

**Thèmes**

- Conditions de vie et société 349
- Données localisées 131
- Opinions 61
- Salaire et revenus 27
- Santé et protection sociale 15

— Voir plus —

**Mots-clés**

- Logement 245


Séries (24) Jeux de données (312) Variables (481)

**XESN** Nombre de logements situés à l'étranger, dans un DOM ou un TOM utilisés comme résidence secondaire ou comme autre résidence principale  
[Enquête Logement \(version FPR\) - 2001-2002](#)

**DISRS** Disposition d'une résidence secondaire  
[Enquête quadrimestrielle de conjoncture auprès des ménages - 1982](#)

**NSEC99** Nombre de résidences secondaires et logements occasionnels au RP99  
[Capacité des communes en hébergement touristique - 2004](#)

**MRAPPORTRS...** Gains liés à la location de la résidence secondaire depuis un an - définitif  
[Budget des familles - 2011](#)

 JEU DE DONNÉES

Enquête Emploi en continu (version FPR) - 2021 Enlever du panier

Version 1 date : 2023-06-15

Résumé Détails **Variables** Données Documentation Accès

Identifiants :  
iii-1525b  
doi:10.13144/iii-1525b

Thème :  
Travail et emploi

Séries : Enquête Emploi / Enquête Emploi en continu (EE / EEC)

Couverture géographique :  
France métropolitaine  
Guadeloupe  
Martinique  
Guyane  
La Réunion

Producteur :  
[INSEE](#)

Diffuseur :  
[Progedo-Adifp](#)

**AAC** Exercice d'une activité professionnelle régulière antérieure, pour les inactifs, chômeurs et personnes ayant une activité

**AACGR\_PAR1** Parent 1 a déjà travaillé

**AACGR\_PAR2** Parent 2 a déjà travaillé

**AAH** Perception de l'AAH (allocation aux adultes handicapés)

Fichier : indiv21  
Type de variable : caracter

Statistiques :

Nombre de réponses	-	338023
Nombre total d'observations	-	343304

Modalités de réponses :

Modalité	Code	Fréquence
	1	5566
	2	331829
	9	628



# PROGEDO, un nouveau portail de présentation des données et de **nouvelles opportunités** (3/3)

## Enquêtes :

un **dépôt désormais ouvert pour des enquêtes réalisées par des chercheurs** qui respectent les standards de qualité des enquêtes de la statistique publique (cf. une vingtaine d'enquête en cours d'évaluation)

*Exemple de l'enquête de l'ANR VICO sur la vie pendant le premier confinement*



The screenshot shows the PROGEDO website interface. At the top, there is a banner for 'La vie en confinement' with the subtitle 'Le carnet Hypothèses du projet et de l'ANR VICO'. Below the banner is a navigation bar with links: 'Études & Résultats', 'Actualités', 'Données', 'On en parle', and 'À propos'. The main content area is divided into two columns. The left column contains navigation links: 'ARTICLE SUIVANT' (Enquêter sur un événement historique exceptionnel : objectifs et premiers résultats), 'ARTICLE PRÉCÉDENT' (La vie en confinement), and a 'CATÉGORIES' list including 'Actualités', 'Actualités du programme Vico', 'Appels à contributions', 'Publications', 'Recrutement, bourses, prix...', 'Études & Résultats', and 'On en parle'. The right column displays the article title 'Enquêter sur un événement historique exceptionnel : synthèse des objectifs et des premiers résultats', the publication date 'PUBLIÉ 12/05/2020 - MIS À JOUR 12/12/2020', the author 'Par l'équipe de l'enquête VICO', and a download button 'Télécharger'. A note at the bottom of the article states: 'Note : Pour plus d'informations sur le mode de collecte et d'administration de l'enquête, ainsi que sur les caractéristiques de l'échantillon et sur les premiers résultats, vous pouvez consulter [ici une version plus longue de cette présentation.](#)' Below the note, there is a paragraph: 'Nous traversons une situation de crise sanitaire et sociale mondiale, marquée par une obligation de confinement qui a interdit ou très fortement limité les déplacements, les activités (de travail comme de loisirs) et les contacts qui remplissaient nos vies ordinaires. Dans cette situation'.



# Agréger et libérer des données personnelles initialement non-libres

Agréger des données “publiques” personnelles (obtenues sous convention de recherche, par ex. via PROGEDO)

- Possible si respect du RGPD
- En lien avec la loi pour une République numérique (dite "loi Lemaire")
  - Le premier volet de la loi vise à favoriser la “circulation des données et du savoir” à travers **l’ouverture des données publiques et d’intérêt général**, la création d’un service public de la donnée et le libre accès aux écrits de la recherche publique.
  - Avec la loi, l’ouverture des données publiques ou *open data* **devient la règle et non plus l’exception**.

*« Les données doivent être aussi ouvertes que possibles, et fermées autant que nécessaire »*

Associer à ces données agrégées

- le code source ayant servi à les créer
- le dictionnaire des données

=> pour le R des principes FAIR - principe de **Réutilisabilité**

# Exemple 1 -



## Le Mobiliscope

### **Demande d'accès à des fichiers de production et recherche (FPR) n°23525**

(...)

Je soussigné(e) :

m'engage à respecter les conditions d'utilisation des fichiers diffusés par Quetelet-Progedo Diffusion :

- 
- ne pas céder ces données, sous quelque forme que ce soit, à une tierce personne, que ce soit à titre gratuit ou onéreux ;
  - traiter ces données conformément aux règles de l'art et du secret statistique ;
  - mentionner la source des données dans mes communications, publications... conformément au modèle de citation pour l'utilisation des données (cf. ci-après) ;
  - informer le diffuseur de mes communications, publications... et lui en faire parvenir un exemplaire et les références ;
  - informer le diffuseur des constats relatifs à la qualité des données ou à leur difficulté d'utilisation ;
  - informer le diffuseur de toute réutilisation des données pour une autre recherche que celle spécifiée ci-dessus ;
  - détruire les fichiers à l'issue du travail de recherche ;
  - respecter la réglementation en matière de protection des données personnelles.
- 

Je déclare avoir pris connaissance de ce que toute infraction aux engagements mentionnés ci-dessus m'expose à des poursuites d'ordre pénal et / ou à des poursuites en responsabilité civile, avec toutes les conséquences pécuniaires que cela comporte au titre des dommages causés.

Comment faire pour lever l'ambiguïté de l'expression “*s’engage à ne pas diffuser ces données sous quelque forme que ce soit*” ?

Début 2021 : prise de contact et échanges avec le **producteur** de la donnée (Cerema) et le **diffuseur** de la donnée (Progedo) pour obtenir leur accord sur la diffusion en open-data des données agrégées respectant le RGPD.

Résultats :

- Accord (écrit) obtenu du producteur et du diffuseur en avril 2021
- PROGEDO nous a signalé vouloir modifier le modèle des futures conventions pour lever l'ambiguïté

Feu VERT pour diffuser les données agrégées (et pondérées) de présence par heure et par secteur ventilées par groupe socio-démo

Ce qui a été fait :

- directement via le Mobiliscope : <https://mobiliscope.cnrs.fr/>
- et via le dépôt Zenodo : <https://doi.org/10.5281/zenodo.4670766>

## DOI:

DOI 10.5281/zenodo.4670766



## Keyword(s):

social inequalities; segregation; daily mobility; place effects; daycourse of place; cities; geovisualisation

## Related identifiers:

Compiles  
<https://mobiliscope.cnrs.fr>

Derived from  
<https://github.com/Geographie-cites/mobiliscope>

Supplementary material  
<http://mappemonde.mgm.fr/123geov3/> (Journal article)  
 10.1016/j.jtrangeo.2017.02.003  (Journal article)  
 10.1016/j.socscimed.2017.09.033  (Journal article)

## Communities:

Géographie-cités

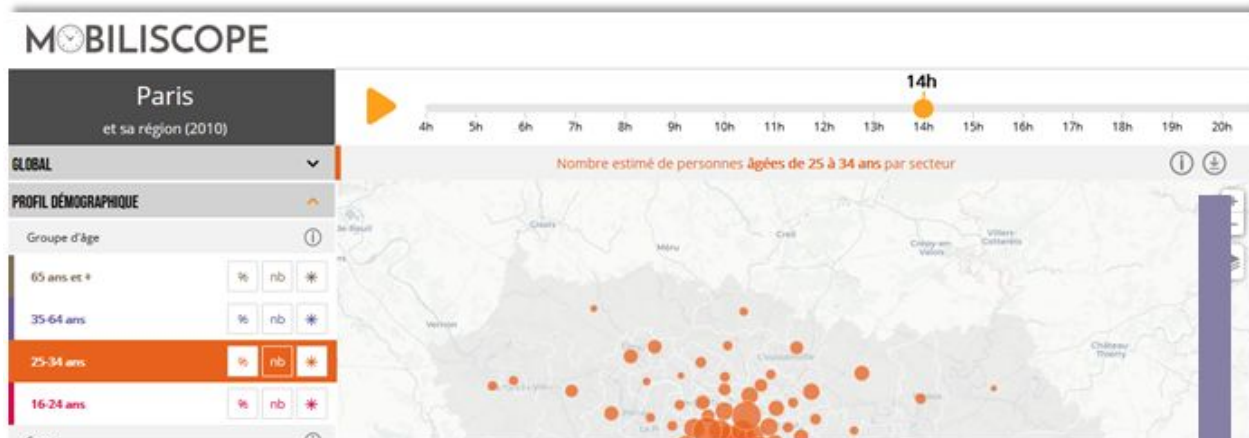

## License (for files):





 [Afferro General Public License v1.0](https://creativecommons.org/licenses/by/4.0/) or later

## Versions

Version v4.0_Extract of French Data	May 28, 2021
10.5281/zenodo.4900655	
Version v4.0	May 28, 2021
10.5281/zenodo.4670766	
Version v3.3	Apr 29, 2020
10.5281/zenodo.4549739	
Version v3.3_full	Apr 29, 2020
10.5281/zenodo.3775577	
Version v1_v2_v3.1	Mar 19, 2019
10.5281/zenodo.3628714	

**Cite all versions?** You can cite all versions by using the DOI [10.5281/zenodo.3628713](https://doi.org/10.5281/zenodo.3628713). This DOI represents all versions, and will always resolve to the latest one. [Read more.](#)


 data\_IDF\_age\_nb

 age\_nb  
 dictionnaire  
 IDF\_secteurs.geojson  
 LicenceODbL\_Mobiliscope\_IDF

# Exemple 2 - CASSMIR (Projet WIsDHoM)

The screenshot shows the Cybergegeo website interface. At the top, the logo 'cybergegeo' is displayed in a yellow box, with the text 'european journal of geography' and 'revue européenne de géographie' below it. A search bar labeled 'Recherche' is on the left. The main content area features a blue bar with 'Data papers' and '2021'. A prominent blue box contains the number '992'. The title of the dataset is 'Une base de données pour étudier vingt années de dynamiques du marché immobilier résidentiel en Île-de-France'. Below the title, there are two lines of descriptive text in English and Spanish. The author's name 'Thibault Le Corre' and a DOI link are also visible.

**cybergegeo**  
european journal of geography  
revue européenne de géographie

Recherche

Data papers

2021

**992**

**Une base de données pour étudier vingt années de dynamiques du marché immobilier résidentiel en Île-de-France**

*A database to analyze twenty years of housing market dynamics in Paris metropolitan region*  
*Creación de una base de datos para estudiar veinte años de dinámicas del mercado inmobiliario residencial en Île-de-France*

**Thibault Le Corre**

<https://doi.org/10.4000/cybergegeo.37430>

## Dataset Description

Name: CASSMIR  
 Repository type: Zenodo [dataset].  
 Dataset available: YES  
 Permanent URL: <http://doi.org/10.5281/zenodo.4497219>  
 DOI: 10.5281/zenodo.4497219  
 Supplemental materials: URL: <https://tlecorre.gitpages.huma-num.fr/cassmir/>

Ce jeu de données est déposé sous la licence Creative Commons paternité - usage non commercial - partage à l'identique 4.0 International (CC BY-NC-SA 4.0) <sup>16</sup>

Objectif : Deux sources (BIEN et PTZ) de données individuelles sensibles qui sont agrégées pour être libérées et diffusées à la communauté scientifique



# Une base accessible avec des données spatiales et populationnelles stockées sur l'entrepôt Zenodo



September 15, 2020

## CASSMIR

 Thibault Le Corre

Project leader(s)

 Thibault Le Corre

Project member(s)

Ronan Ysebaert; Pierre Le Brun; Timothée Giraud; Jean-Baptiste Durand

New version 2.0.0 with majors change








For free and complete informations concerning CASSMIR datasets, please visit our [website](#)

DOI : <http://doi.org/10.5281/zenodo.4030698>

Data Paper

[Preview](#)

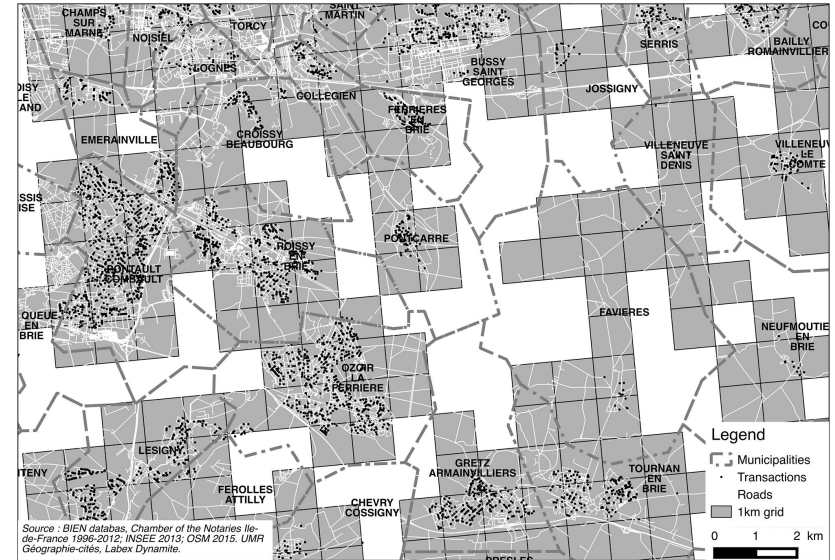
Files (886.8 MB)

Name	Size	
BIENSsampleForTest.csv	9.1 MB	<a href="#">Preview</a> <a href="#">Download</a>
md5:a56ea2640ace7c84a1080e32f7a15e92 		
CASSMIR_GroupesPopDataBase.csv	1.0 MB	<a href="#">Preview</a> <a href="#">Download</a>
md5:a8f57c0e76c8c95a4c4933dbfbd794f1 		
CASSMIR_PopMetadata.csv	65.1 kB	<a href="#">Preview</a> <a href="#">Download</a>
md5:e3436f3a8cb7f73397bf6d09e3fdbeb7 		
CASSMIR_SpatialDataBase.gpkg	873.0 MB	<a href="#">Download</a>
md5:53eb19d31b04d78d99608bf09cb2a5a0 		
CASSMIR_SpatialMetadata.csv	41.2 kB	<a href="#">Preview</a> <a href="#">Download</a>
md5:1b31bbc3db98a50ab275172fb932a280 		
LEXIQUE.csv	3.1 kB	<a href="#">Preview</a> <a href="#">Download</a>
md5:765412c38b090f19e38551733f70ca42 		
PTZSampleForTest.txt	3.5 MB	<a href="#">Download</a>
md5:4db806c5a0cc66fc18521ebfe2e0cf97 		

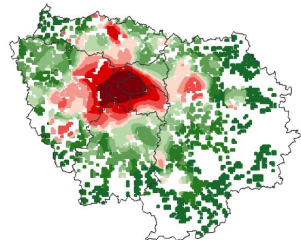


# Des données d'origine à CASSMIR : un focus sur les transactions immobilières (BIEN)

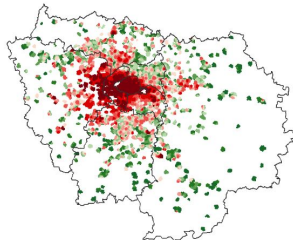
- Une donnée propriétaire sous contrainte
- Gestion et amélioration de la qualité de l'information géographique et statistique
- Détection de communautés par la densité du voisinage
- Interpolation spatiale par la méthode des potentiels qui privilégie hypothèse d'interaction



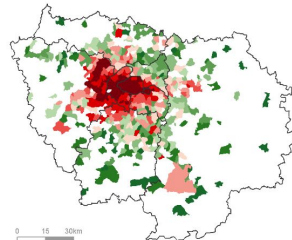
Carroyage 1km lissé



Carroyage 200m lissé

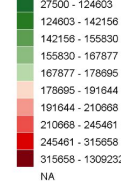


Communes



Prix moyen en 2012

(méthode déciles)





D'une donnée ponctuelle...

...à une donnée agrégée et interpolée à plusieurs échelles

# Un site internet pour comprendre, trouver et manipuler la donnée

GitLab Search GitLab /


Thibault LE CORRE > CASSMIR




**CASSMIR**  Project ID: 708 

94 Commits 1 Branch 0 Tags 13.7 MB Project Storage

Contribution à l'Analyse Spatiale et Sociale des Marchés Immobiliers Résidentiels

master cassmir

 **ModifSite**  
Thibault LE CORRE authored 1 year ago

 README  CI/CD configuration  No license. All rights reserved

Name	Last commit
BIEN_cache/html	ModifSite

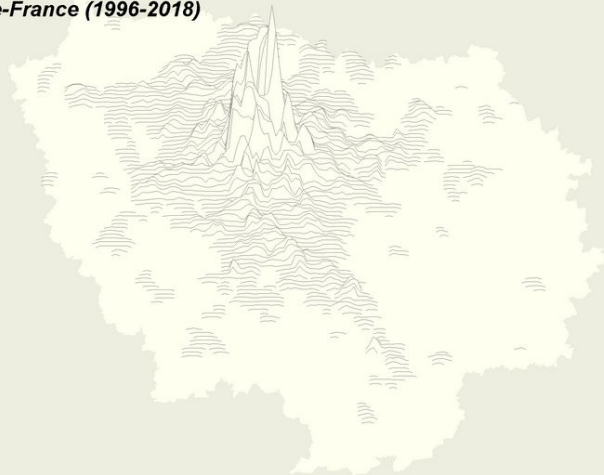
CASSMIR Présentation Projet et sources Echantillons Préparation de CASSMIR BIEN PTZ Meta Exploration About

## Présentation du site internet



Réalisé au sein de l'ANR WisDHoM par Thibault Le Corre (UMR 8504 Géographie-Cités), avec la contribution de Ronan Ysebaert (UMS 2414 RIATE), Timothée Giraud (UMS 2414 RIATE), Jean-Baptiste Durand (UMR 7300 Espaces) et Pierre Le Brun (UMR 7300 Espaces), le projet CASSMIR (Contribution à l'Analyse Spatiale et Sociale des Marchés Immobiliers Résidentiels) poursuit l'objectif de produire un jeu de données sur le marché du logement à la propriété en Ile-de-France, prêt à l'emploi, diffusable et réutilisable par la communauté scientifique.

### Les transactions immobilières en Ile-de-France (1996-2018)



# De “nouvelles” données privées en contradiction avec les principes de science ouverte ?

“**Nouvelles**” ? Nativement numériques et individuelles (mais parfois agrégées avant diffusion), différents types selon les domaines thématiques et les objets techniques qui les supportent...

Données	Mode de production	Source et accès	Composantes spatiales, temporelles et attributaires
<b>Mobilités et crise COVID-19 (1)</b> Présence sur réseau social Meta <i>(projet MAMA et DITES)</i>	Enregistrement automatique via Wifi/GPS	<b>Meta</b> <a href="https://dataforgood.facebook.com/dfg/about">https://dataforgood.facebook.com/dfg/about</a> , Mise à disposition <b>gratuite</b>	<b>Données zonales</b> (agrégation à l'échelle de « tuiles » de taille variable), pour une centaine de pays <b>3 périodes par jour</b> <b>Attributs</b> : effectif de population
<b>Mobilités et crise COVID-19 (2)</b> Flux de téléphonie mobile <i>(projet MAMA)</i>	Enregistrement automatique par bornage aux antennes	<b>Flux vision</b> <a href="https://flux-vision.orange-business.com/">https://flux-vision.orange-business.com/</a> Accès <b>payant</b>	<b>Données zonales</b> (agrégation à l'échelle des EPCI), France <b>Toutes les heures</b> pendant 10 périodes de 3 semaines <b>Attributs</b> : effectif de population, caractéristiques socio-économiques
<b>Locations immobilières temporaires</b> Plate-forme commerciale Locations Airbnb <i>(projet WHIsDHOM)</i>	Saisie par les hôtes et par les loueur.ses  Puis données webscrapées par l'entreprise AirDNA	<b>AirDNA</b> <a href="https://www.airdna.co/">https://www.airdna.co/</a> Accès <b>libre</b> à la visualisation (site web) Accès <b>payant</b> aux jeux de données	<b>Données ponctuelles</b> (résolution spatiale à 100 m près), pour <b>tout territoire</b> <b>Agrégation mensuelle</b> , depuis 2016 <b>Attributs</b> : descriptif logement, prix, disponibilité, revenus générés...

D'après Guérois et Madelin. 2023

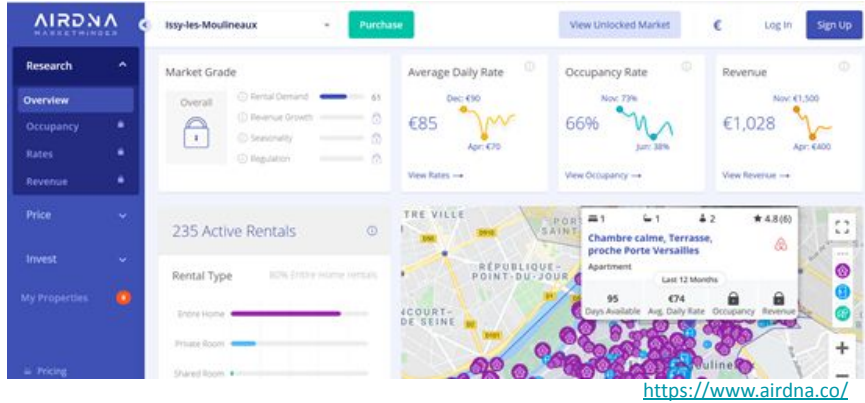
**Enjeux communs** pour la recherche (Arribas-bel 2014, Kitchin 2014, Li et al 2016...) :

- 1) données « accidentelles » pour scientifiques (manque de normalisation de l'info, pauvreté des métadonnées)
- 2) haute résolution spatio-temporelle, moindre richesse thématique (enjeu interopérabilité avec données de référence)
- 3) difficultés d'accès (sauf via achat)



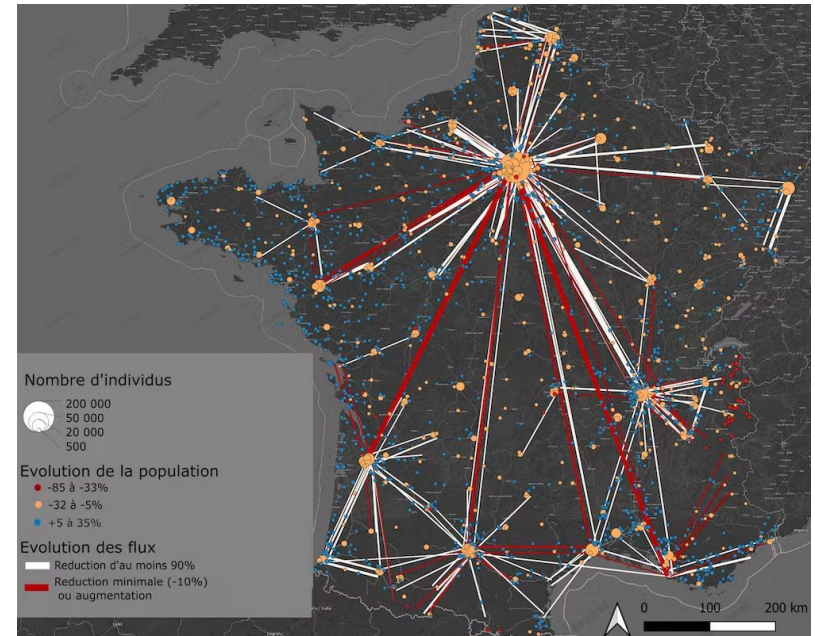
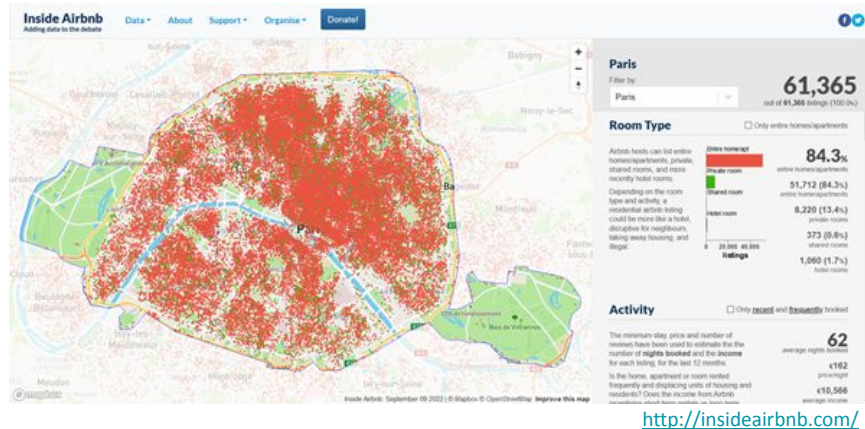
# Accès aux “nouvelles” données privées : à quel prix ?

## ● AirDNA vs InsideAirbnb



## ● Orange/Meta: 2 possibilités pour mesurer les changements des mobilités induits par la crise

spaco_id	Lat	Long	polygon na	Country	03.10.0800	03.10.1600	03.11.0000	03.11.0800	03.11.1600	03.12.0000	03.12.0800
131273	43.2927319	5.4202526	Marseille	FR	58229.00	57517.00	47689.00	57424.00	57195.00	47968.00	57618.00
132370	48.8566166	2.3429364	Paris	FR	209325.00	168074.00	133616.00	207550.00	169580.00	134606.00	208091.00
138794	43.5960126	1.4322893	Toulouse	FR	40491.00	36555.00	30546.00	39645.00	36528.00	30696.00	40123.00
155434	45.7549304	4.83637276	Lyon	FR	40265.00	36229.00	29076.00	39673.00	36652.00	29922.00	40167.00



# “Nouvelles” données privées : quel partage ? A minima, accès ouvert aux codes (pré-traitements + traitements)

## ● AirDNA

Préparation, analyse des données disponibles dans un site Web documenté ([https://riatestage.github.io/airdna\\_laurian/index.html](https://riatestage.github.io/airdna_laurian/index.html))



## Mise en oeuvre d'une chaîne de traitement de données AirDNA

Étudier l'activité de la plateforme Airbnb dans 3 communes franciliennes

Louis Laurian, M2 GAED - Géographie, aménagement environnement et développement, parcours Géoprisme, Université de Paris

### 1 Import et description des variables

- 2 Sélection des variables
- 3 Sélection des lignes
- 4 Traitement des valeurs non attribuées
- 5 Critique des variables
- 6 Création de nouvelles variables
- 7 Couches géographiques
- 8 Export des données

### 1 Issy-les-Moulineaux

- 1.1 Etat des lieux sur l'importance, la structure et la performance de l'offre Airbnb à Issy (2019)
- 1.2 Analyse temporelle 2015-2020
- 1.3 Analyse infra-communale
- 2 Pantin
- 3 Bagneux

## Préparation des données

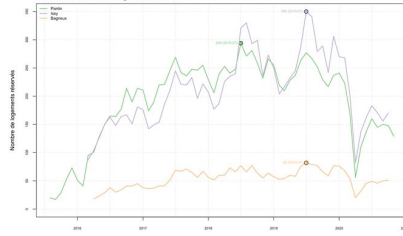
```
library(knitr)
library(rmdformats)
```

Cette partie vise à préparer et nettoyer les données des logements Airbnb fournis par AirDNA en vue de diverses analyses. La finalité de ce markdown est d'exporter les tableaux de données ainsi que les couches géographiques nettoyées de manière générique, de sorte à ce que ce code soit reproductible peu importe la commune d'entrée. Les deux paramètres à rentrer en début d'exécution du code sont le code INSEE ainsi que le nom de la commune.

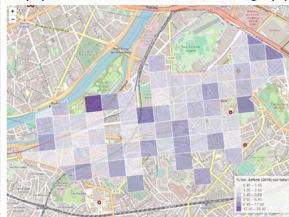
Afin de préparer les données pour une autre commune, le script `prep.R` condense toutes les opérations ci-dessous sous forme de fonctions, pour obtenir les mêmes résultats pour la commune sélectionnée, c'est-à-dire des fichiers consolidés ainsi qu'un geopackage regroupant les couches géographiques de la commune concernée. Ce script est appelé dans le markdown

## Synthèse à la commune

Evolution du nombre de logements réservés



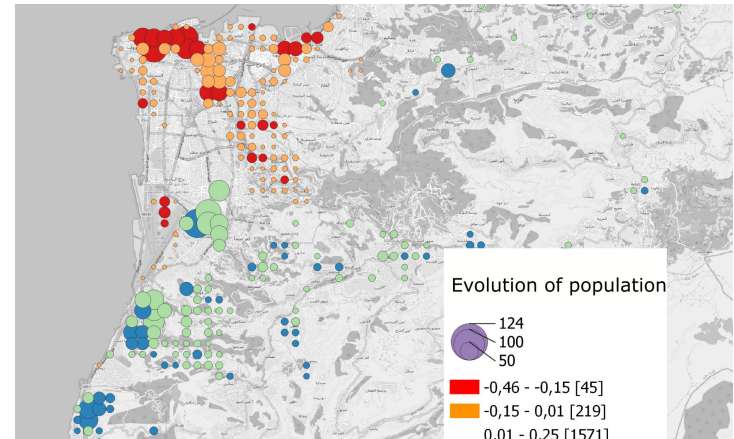
Issy : part locations actives Airbnb / total logts (%)



## ● Orange/Meta:

Les projets “Mama” et “Dites” ont permis d’engager un ingénieur (Joséphin Béraud) qui établit des chaînes de traitements sous R.

```
49 # Aggregation à la semaine en gardant la colonne période
50 data_week <- data %>% group_by(week, periode) %>% summarize(moy = mean(n_crisis))
51
52 ggplot(data_week, aes(x = week, y = moy)) +
53   #geom_rect(data=rectangle, inherit.aes=FALSE, aes(xmin = xmin, xmax = xmax, ymin = ymin, ymax =
54   # alpha=0.2)) #+ # rectangle pour griser les dates voulues (confinement par ex)
55   geom_line() +
56   facet_grid(~periode, scales = "free") + # facet selon la colonne periode, scales = free pour
57   labs(y = "Nombre de X", x = NULL, title = "TITRE",
58        subtitle = "Sous titre",
59        caption = "Note" + #, caption = \n ? la ligne") +
60   scale_x_date(breaks = "15 days", minor_breaks = "1 day", date_labels = "%d-%m-%Y") + # Affich
61   theme(legend.position="bottom", plot.caption = element_text(hjust = 0), legend.title=element_
62
63 - #####
64 # Carto
65
66 # Une fois la donnée traitée comme souhaité dans l'objet data :
67
68 #import des shp
```



## “Nouvelles” données privées : quel partage des données sources et des résultats ?

### ● *Airbnb* :

- Projet mené avec R. Ysebaert (RIATE) et M. Madelin (PRODIG) sur l’Ile-de-France :

Pas d’ouverture possible des données sources *AirDNA*, directement → projet de diffusion des résultats sous la forme d’un site web avec des représentations cartographiques agrégées et lissées à différentes échelles pour l’Ile-de-France → quel degré de dégradation de l’information ?

Pour la suite : échanges avec d’autres utilisateurs (institutionnels, académiques) sur les possibilités de mutualisation de l’accès aux données *AirDNA*.  
Voie alternative : requête auprès d’*InsideAirbnb* pour élargir la collecte des données de la métropole parisienne au-delà de Paris intra-muros

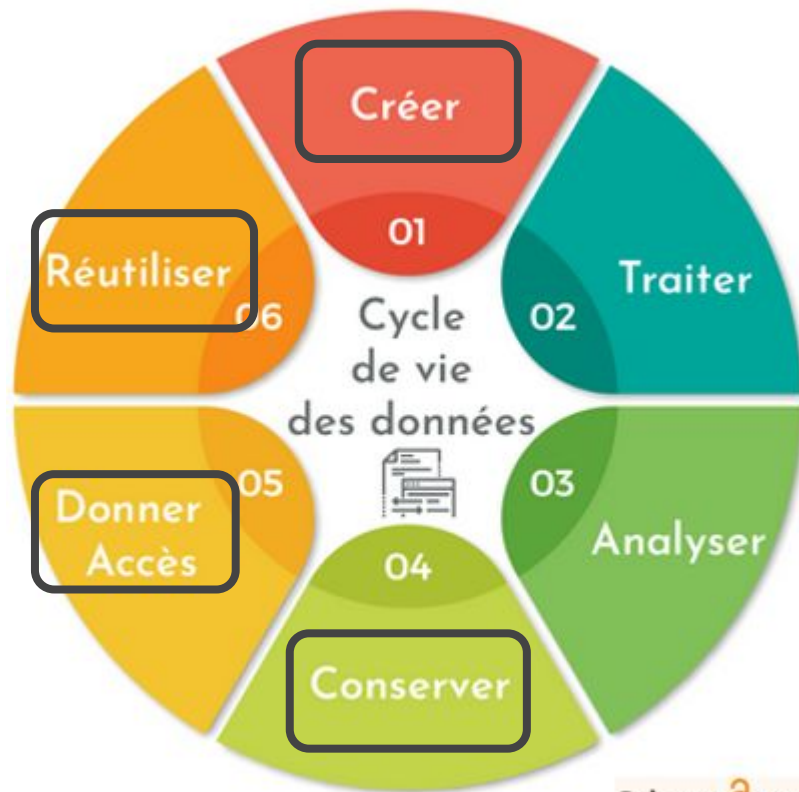
### ● Méta:

- Données stockées sur un NAS à géographie cités (2 TO) et disponibles pour les chercheurs de Géographie-cités. D’autres données sont disponibles mais non collectées/archivées (électricité, proximité relationnelle des utilisateurs à l’échelle des territoires, analyse des posts Facebook...).

<https://help.crowdtangle.com/en/articles/4302208-crowdtangle-for-academics-and-researchers>

- Valorisation des résultats auprès de la communauté scientifique “locale”, les résultats étant difficiles à analyser sans connaissance du territoire. Pour l’instant: Égypte, Inde, Japon, Singapour, Thaïlande, USA, France.

- Projet de mise en place d’un site présentant en “temps réel” les mobilités observées et leurs liens avec divers phénomènes. Accord en Inde avec Meta “Asie”.

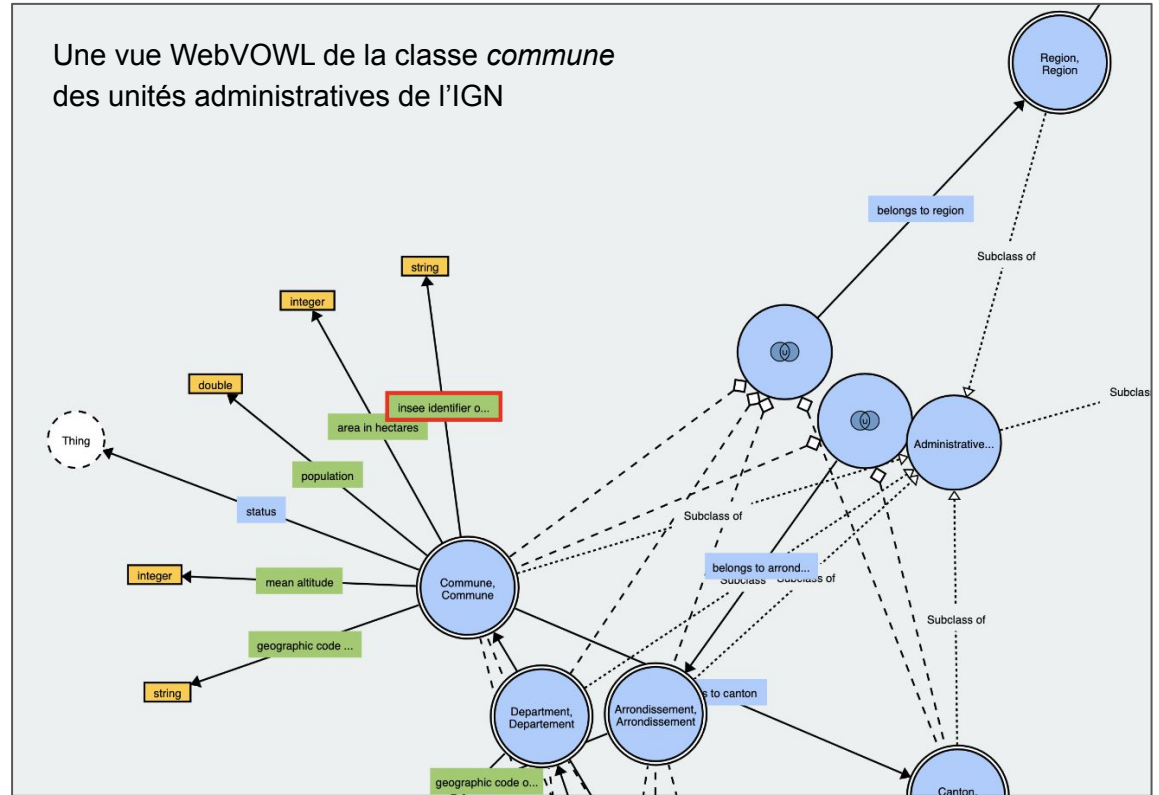




# Principes et réalités du / du FAIR, cas d'ontologie

## Ontologie axiomatisée

Ontologies formalisées avec des contraintes et un langage interprétable tel que OWL, qui permettent à une machine de raisonner sur les éléments représentés. OWLTime ou GeoSparql sont des exemples d'ontologies axiomatisées.



N. Abadie et G. Atemezing, *Ontologie des unités administratives de l'IGN. Revision: Version 1.1 - 2019-02-12*. Institut National de l'Information Géographique et Forestière, 2019. Disponible sur: <http://data.ign.fr/def/geofla/20190212.htm>



# Principes et réalités du / du FAIR, cas d'ontologie

Contexte : archéologie à Al-Ula



Diversité d'acteurs et de projets en présence

الهيئة الملكية لمحافظة العلا  
Royal Commission for AlUla



MINISTÈRE  
DE L'EUROPE  
ET DES AFFAIRES  
ÉTRANGÈRES

Etc. !

J. Charbonnier, Y. Kanhoush, E. Devaux, J. Gravier, V. Bernollin, et T. Hofstetter, « AlUla Old Town and Oasis – Oasis Farms: preliminary study and first results from the AlUla Cultural Oasis Project (Kingdom of Saudi Arabia) », in *13th World Congress on Earthen Architectural Heritage*, Santa Fe, USA, 2022.

# Principes et réalités du / du FAIR, cas d'ontologie

Contexte : archéologie à AI-Ula

Une base de données commune fondée sur le système Arches-HIP



Entités-relations-thésaurus



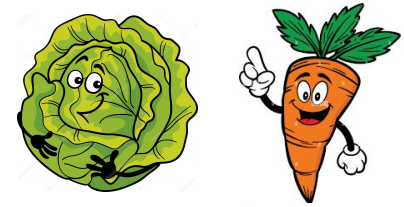
en 2006 ISO 21127, puis ISO 21127:2014

# Principes et réalités du / du FAIR, cas d'ontologie

Point de vue conceptuel



[@aaksa\\_project](#)

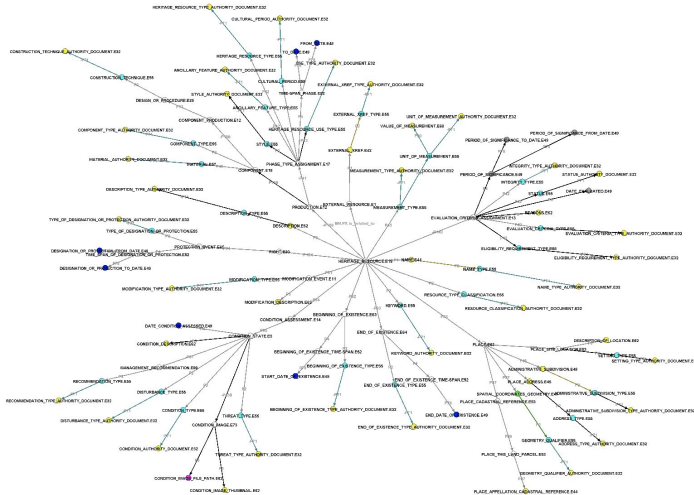


[@AF ALULA](#)

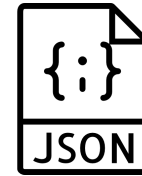
# Principes et réalités du / du FAIR, cas d'ontologie

## Point de vue technique

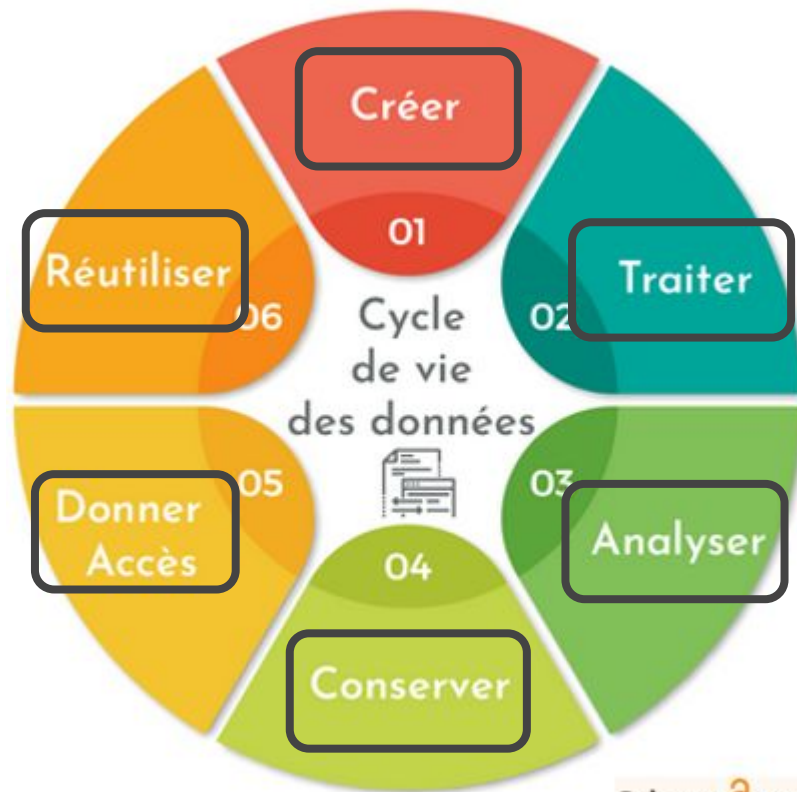
Graphe conceptuel de l'entité Heritage Resource - Arches-HIP



Une base de données orientée graphe

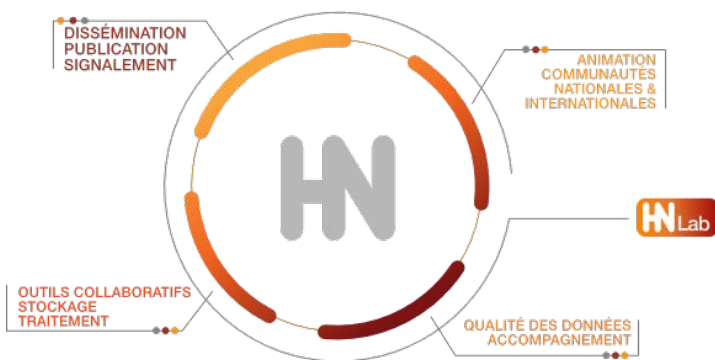


R. Gaston, C. Hiatt, R. Anderson, A. Peters, et A. Cox, « Arches-HIP Documentation ». Read The Docs, Getty Conservation Institute and World Monuments Fund, 2015. Disponible sur: <http://arches-hip.readthedocs.io/en/latest/>



## Mission : **Construire une infrastructure numérique de niveau international pour les SHS**

- Accompagnement des communautés scientifiques SHS en matière d'infrastructure numérique pour les données de la recherche.
- **Mise en œuvre d'une infrastructure numérique permettant aux SHS de développer, réaliser et préserver sur le long terme les programmes de recherche – leurs données et outils – dans un contexte de science ouverte et de partage des données (principes FAIR)**



### Plusieurs actions :

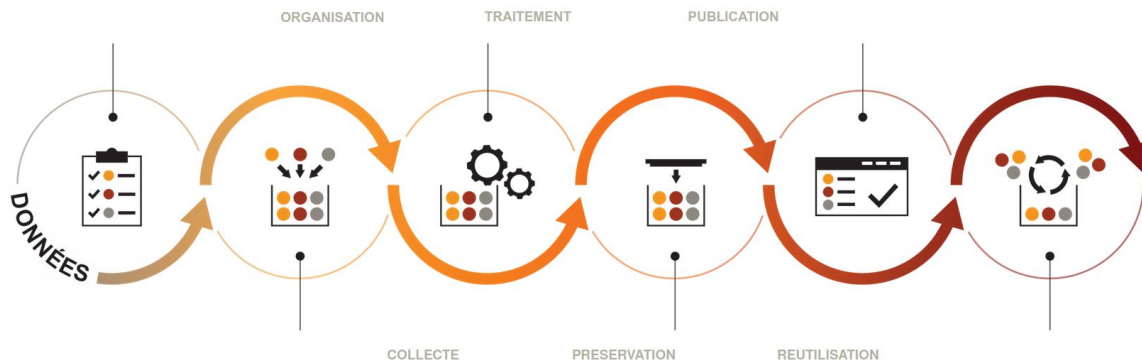
- Dissémination, publication, signalement
- Animation communautés nat. et internationales
- Huma-Num Lab
- Mise à disposition d'outils et services
- Formation



## Les outils et services d'HumaNum

**Services et outils mis à la disposition du monde académique français pour mutualiser, diffuser et stabiliser l'accès aux données et documents** produits dans les projets de recherche en SHS.

La mission première est d'**assurer la préservation du patrimoine scientifique des laboratoires** (corpus, bases de données, bases documentaires, systèmes d'information, enquêtes, données d'observation produites, en cours de production....). Cette mission sous-tend également une **stratégie économique visant à diminuer les coûts récurrents**, par la mise en commun d'une infrastructure en co-gérant des outils, instruments et systèmes de gestion des données.



Un services par étapes couvrant l'ensemble du cycle de vie des données


Quels outils ?  
Quels services ?  
Qui peut les utiliser ?  
Avec quels niveaux d'accès ?  
Quelle marge de manœuvre ?

**Parfois flou, car en évolution.**

## Outils HN et cycle de vie des données

KB

 Studio Server

 jupyterhub

 GitLab

 matomo



Mattermost

 ShareDocs



### ORGANISATION

Des services pour organiser le travail collaboratif autour de vos données.

- ShareDocs
- GitLab
- Kanboard
- Mattermost

### TRAITEMENT

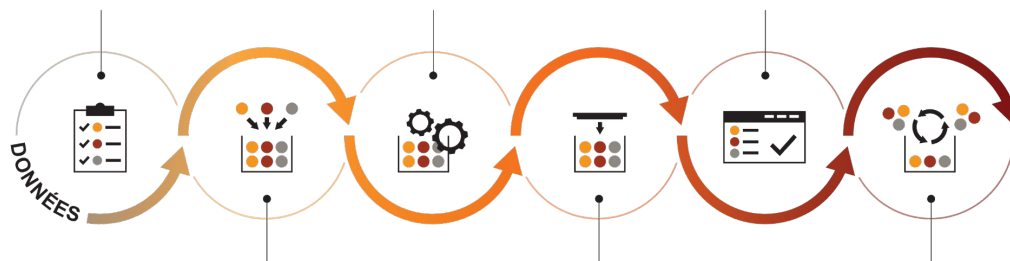
Des services et outils spécifiques pour le traitement et l'analyse de vos données.

- Calcul statistique et environnements R
- Logiciels d'enquête et d'analyse de données
- Reconnaissance de caractères
- Puissance de calcul (+ CC-IN2P3)
- ...

### PUBLICATION

Vos données peuvent être publiées depuis Nakala sur le web et signalées dans Isidore, moteur de recherche pour les SHS.

- Hébergement Web
- Machines Virtuelles
- Nakala
- Isidore



### COLLECTE

Des services de stockage sécurisé pour la collecte et la création de vos données.

- ShareDocs
- Huma-Num Box

### PRÉSERVATION

Huma-num vous accompagne pour le dépôt et la documentation de vos données dans Nakala, entrepôt pour les données en SHS.

- Nakala
- Huma-Num Box
- Préservation à long terme (+ CINES)

### RÉUTILISATION

Vos données entreposées dans Nakala et signalées dans Isidore sont réutilisables.

- Portail web
- API
- Triplestore
- OAI-PMH



 nakala

 isidore



# Ressources

Sur la publication en ligne et la réutilisation des données publiques (« open data ») :

<https://www.cnil.fr/fr/publication-en-ligne-et-reutilisation-des-donnees-publiques-open-data>

Sur l'application du RGPD à la recherche scientifique :

Le guide de l'INSHS du CNRS :

[https://www.inshs.cnrs.fr/sites/institut\\_inshs/files/pdf/Guide\\_rgpd\\_2021.pdf](https://www.inshs.cnrs.fr/sites/institut_inshs/files/pdf/Guide_rgpd_2021.pdf)

<https://www.cnil.fr/fr/recherche-scientifique-hors-sante-les-questions-reponses-de-la-cnil>

<https://cnpd.public.lu/dam-assets/fr/dossiers-thematiques/droit-image/CNPD-Lignes-directrices-droit-a-l-image-protection-donnees-personnelles.pdf>



**InSHS**

Les sciences humaines et sociales et la protection  
des données à caractère personnel dans le  
contexte de la science ouverte

**GUIDE POUR LA RECHERCHE**

VERSION 2

FÉVRIER 2021

