

8^e ÉDITION

JOURNÉES DU GFCO 2022

Biomarqueurs et analyses moléculaires en oncologie

Avec la participation
scientifique du



PLÉNIÈRE - REAL-WORLD DATA

Intégration des données pour l'émergence de nouveaux concepts et l'amélioration du parcours de soins

Pr. Pierre-Antoine Gourraud, Nantes

CHU Nantes, PHU 11 : Clinique des données, INSERM, CIC 1413, Nantes, France

Nantes Université, INSERM, CRT21 , Nantes, France

pierre-antoine.gourraud@univ-nantes.fr



LIENS D'INTÉRÊT

- **COI :**
- *He consults for major pharmaceutical companies, all of which are handled through academic pipelines (AstraZeneca, Biogen, Boston Scientific, Cook, Docaposte, Edimark, Ellipses, Elsevier, Methodomics, Merck, Mérieux, Sanofi-Genzyme, Octopize). PA Gourraud is the founder of Methodomics (2008) and the co-founder of Big data Santé (2018).*
- PA Gourraud is volunteer elected board member at AXA mutual insurance company (2021).
- He has no prescription activity with either drugs or devices.



Introduction et Plan

- « Intégration des données pour l'émergence de nouveaux concepts et l'amélioration du parcours de soins »
- **Intro et Détour par la génétique (épidémiologique)**
 - 3 concepts « approches agnostique », secondaires et intentionnalité
 - Une nouvelle activité hospitalière « la clinique des données »
- **« Parcours de soins - RWE » Quelles Données ? SNDS et EDS**
- **« Prendre soin des données » : Les entrepôts de de données de santé, les clinique des données dans les établissements**
 - Gouvernance légalité technique
 - EXEMPLES
- **Conclusion : Carrefour de données de Santé**



“Our generation will be called naïve about data”



- **“Nous sommes des naïfs de la donnée”**

PA Gourraud Y Coatanlem

Le Monde Oct 5th 2021

« Avec des protocoles d'accès plus souples, les données publiques pourront constituer un gisement de valeur du XXI^e siècle »

Tribune 05.10.2021

Yann Coatanlem

Président du club de réflexion Praxis

Pierre-Antoine Gourraud

Professeur à la faculté de médecine de l'université de Nantes

Les moyens existent de libérer l'exploitation des données tout en protégeant la confidentialité, notamment pour le croisement et le partage des fichiers de vaccinations et de tests, expliquent, dans une tribune au « Monde », Yann Coatanlem, président du club de réflexion Praxis et Pierre-Antoine Gourraud, professeur de médecine.

Publié aujourd'hui à 06h15 Temps de Lecture 4 min.

Tribune. Nous sommes des naïfs de la donnée ! Bien que plus ou moins conscients que les données sont au XXI^e siècle l'équivalent de la terre arable à l'ère agricole ou de la machine au XIX^e siècle, nous n'exploitons encore qu'insuffisamment les gisements d'opportunités dans ce domaine. Aujourd'hui, en pleine crise de Covid-19, le croisement et le partage des fichiers de vaccinations et de tests posent encore problème alors même que les enjeux de santé publique sont criants. C'est donc un véritable aggiornamento des politiques en la matière que nous appelons de nos vœux.

Dans les débats publics, les enjeux sont malheureusement souvent confondus : enjeux de confidentialité, d'usage (la finalité de l'analyse des données), d'usages secondaires (par opposition à l'intentionnalité première des données), de contrôle des usages (quelles données, pour faire quoi), de contrôle des usagers (par qui), de sensibilité (quelles sont les conséquences potentielles de l'interprétation des données). Cette confusion nuit à la transparence, à la collecte, à l'organisation, à la valorisation des données. Elle nuit finalement à la confiance requise pour que le développement économique se nourrisse de la création et de la diffusion des connaissances.

Le Monde



“Re-penser les données de Santé”

SANTÉ

Repenser la production des données de santé

DATA Les données de santé sont parmi les données personnelles les plus sensibles. Le déploiement de la plateforme Française des données de santé, le Health Data Hub, pour poser des questions médicales de souveraineté technologique, et la possible création d'un espace européen des données de santé en posent d'autres. Pour y répondre, nous devons repenser ce qu'est un « producteur de données de santé », et reconnaître que ces données sont « produites » par soignants et patients.

Dans le numérique, les données sont souvent réduites à un problème informatique dont la solution serait-ils de simples « producteurs » de données de santé, comme un arbre porterait du « fruit » ? Ces données seraient-elles assimilables à des biens matériels, en relation à leur valeur de transaction commerciale ? Cette vision des données, souvent prise pour du pragmatisme, ne vise pas à privilégier l'économique, induit en erreur tant les décideurs politiques que le grand public. La nature productive des données de santé est plus complexe que cela.

D'abord, d'un point de vue légal, les données de santé ne sont pas des biens patrimoniaux. Elles restent intransmissibles à une personne et chaque patient a le droit de s'opposer à ce que ses données qui le concernent soient utilisées dans certains lieux. Il faut donc discuter d'autres, par exemple pour la recherche. Mais elles ne leur appartenent pas au même titre qu'un objet qu'il posséderait et dont il pourrait disposer librement, il est donc plus juste de parler de données « qui concernent les patients » plutôt que de « leurs » données. Dans l'usage courant, l'adjectif possessif (« mes données ») tend à faire des données un objet marchand, sans l'influence du droit anglo-saxon de la data. Mais pour des données personnelles, et encore plus si elles touchent à notre santé, c'est plutôt le champ du droit à l'image qui est pertinent. Un patient exerce un contrôle sur les utilisations des données émanant de lui, comme il le ferait pour une photographie, sans avoir à justifier l'exercice de son droit d'opposition.

Faire de « producteur de données de santé » est de surcroît injuste pour la communauté des soignants. Tous à leur échelle, des médecins aux ingénieurs biomé-



Le déploiement de la plateforme Française des données de santé, le Health Data Hub, pose des questions de souveraineté technologique

dicans en passant par les aides-soignants, contribuent à « produire » les données de santé des patients qui résultent de leurs observations cliniques, de leurs analyses en laboratoire ou en biologie médicale, et de leurs jugements. Il faut donc bien parler à minima de « production » de ces données entre les patients et les soignants. Si les patients sont concernés, les soignants le sont aussi car ces données touchent à l'essence de l'exercice de l'art de la médecine.

Responsabilité et régulation

Patients et soignants sont donc concernés, contributeurs, coauteurs et coauteurs de la santé. Tous sont objets que sujet des données de santé, leurs rôles ne sont jamais vraiment passifs. Ainsi les « producteurs » de données de santé produisent, c'est plus un sens musical qu'agricole ! Cette production musicale, si nous la laissons « danser » entre mouvements et contributions multiples, comme l'évoquait Nietzsche à propos des vertiges des données sont la figure de prose du XXI^e siècle.

Bref, la création de valeur à partir des données de santé à un enjeu de technologie informatique est donc une double création, l'enjeu technique est nécessairement, prend le pas sur les enjeux de gouvernance des données, cher « et entre » les acteurs de la santé, l'hôpital comme en libéral. La technologie demeure un moyen au service d'un fin, sans chercher une recherche sur une supposée toute-puissance médicale.

À l'inverse, si l'enjeu technique est sous-estimé au risque de traiter ces données de santé avec le

général, et d'oublier le vécu ser-

viceant de chaque soignant, qui, « admis dans l'intimité des

soignants, les secrets qui lui sont confiés ». A contrario, le citoyen consent des usages souvent

difficiles possibles avec des données

de santé le concernant, leur donne

une seconde vie et espèrent sa

solidarité avec tous les malades.

La crise du Covid nous a rappelé

que la santé est une responsabilité

collective, incluant les données de

santé. Alors cessons d'être

saute et donnons à la France nos

« réserves d'effluents de données de

santé ». Du Health Data Hub à

chaque acteur de santé (au plus

proche du soin et de l'urgence de

données), seule une organisation

disciplinée permettra

d'accroître la confiance dans les

usages des données, et de conjuguer

les intelligences dans leur

exploitation pour mieux évaluer

et faire évoluer les services de

santé au service de la société. »

Dr. P. Fournier-Lesourd, Généraliste

diagnoste ne pas avoir de lien d'intérêt

en rapport avec le sujet traité. Il est

le fondateur en 2008 de Méthodoma

(www.methodoma.com) et le cofondateur de WeMoia en

2016 (www.wemoia.com). Il est

conseiller et/ou intervenant pour de

grandes entreprises pharmaceutiques

et de dispositifs médicaux. Ses

activités sont toutes traitées par une

contractualisation universitaire ou

hospitalière (AstraZeneca, Biogen,

Roche, Sanofi, Merck, Bayer, Boehringer,

Novartis, Amgen, AstraZeneca, etc.).

Il est administrateur bénévole des

matériaux d'assistance à la vie (MAV).

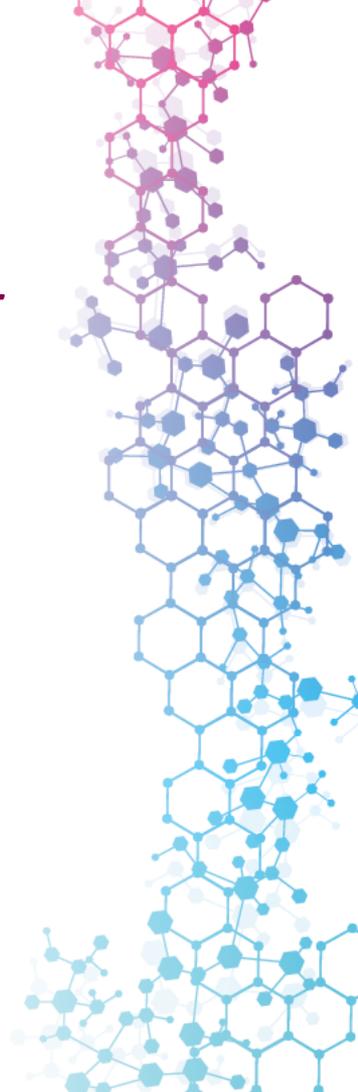
Il a une activité de prescripteur de médicaments

ou de dispositifs médicaux.

- « Patients et soignants sont donc co-concernés, co-contributeurs, co-acteurs et co-auteurs de la santé. »
- « Ainsi si les « producteurs » de données de santé produisent, c'est plus au sens musical qu'agricole ! »

PA Gourraud

Le Figaro 12 Septembre 2022



Intégration **des données** pour l'émergence de nouveaux concepts et l'amélioration du parcours de soins

Les données massives en santé

Introduction to the true nature of (medical) Data

- **Composed as a choral symphony**
 - By L Beethoven
- **Played**
 - ... By an symphonic orchestra
- **Written as a poem**
 - ... By Friedrich von Schiller
- **Officially interpreted**
 - .. By Herbert von Karajan
- **Translated in French English**
 - ... By many
- **Sung**
 - ... by anyone who dares
- **Et caetera ...**



• *Electronic Medical Records*

- **Acquired**
 - From patients
- **Written**
 - ... By caregivers
- **Produced**
 - ... by medical devices
- **Paid**
 - ... by Insurance companies
- **Stored**
 - ... by Care Institutions
- **Transformed**
 - ... by data scientists
- Et caetera ...



The "Ode to Joy"

***Personal Medical Data is not similar to material good
We need to take "good care" of it...***

Définitions : Données Massives en Santé

> Définition par source :

Données d'hospitalisation / consultation	Consommation de soins	Objets connectés	Produites à des fins de recherche	Tout autre type de données
<ul style="list-style-type: none"> • Diagnostic • Traitement • Biologie • Imagerie • Génétique • ... 	<ul style="list-style-type: none"> • Prescriptions, • Données de l'Assurance Maladie • ... 	<ul style="list-style-type: none"> • Médicaux : tensiomètre, implant cardiaque, ... • Non médicaux : montre connectée, pèse-personne... 	<ul style="list-style-type: none"> • Omics : génomique, transcriptomique... • Microbiote • Cohortes • Biocollections • ... 	<ul style="list-style-type: none"> • Données socio-démographiques • Géographiques • Météorologiques • ...

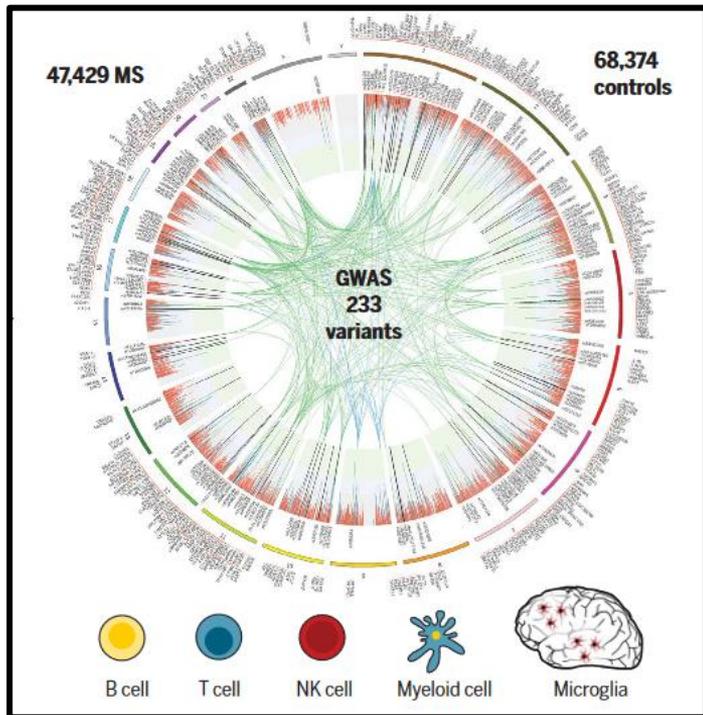
> Définition par structure :

Données structurées	Données non structurées
<ul style="list-style-type: none"> • Valeurs biologiques (hémoglobine, glycémie...) • Données génétiques • PMSI • ... 	<ul style="list-style-type: none"> • Compte-rendu : Texte « brut » -> TALN • Données de signal (ECG) • Images (scanner...) • ...

Intégration des données pour
l'émergence de **nouveaux concepts** et
l'amélioration du parcours de soins

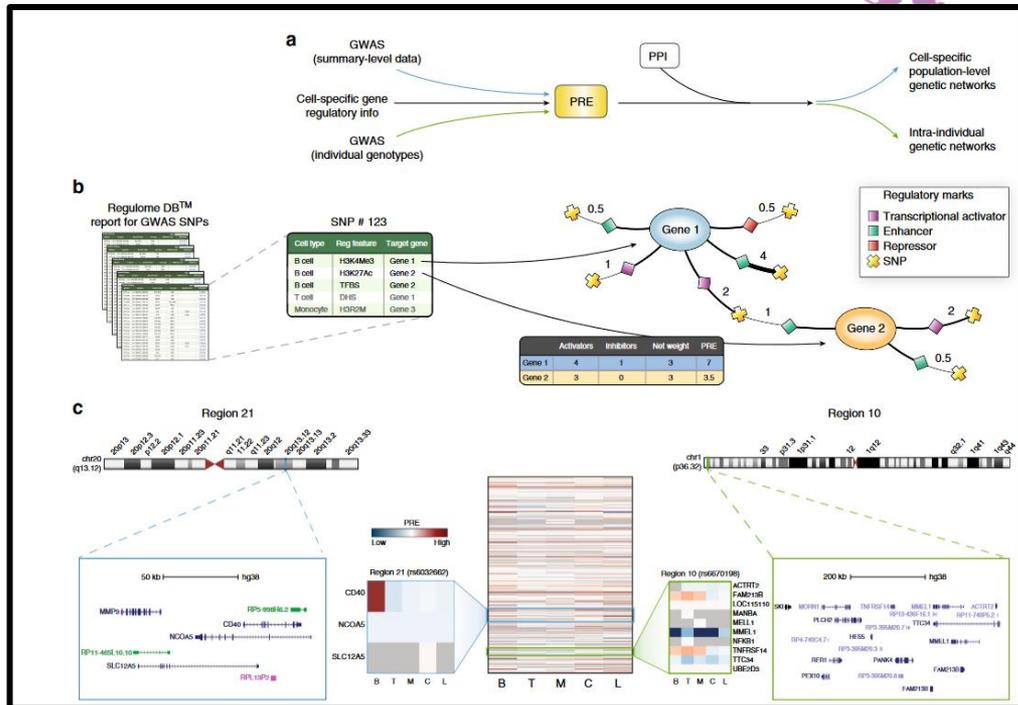
Détour par le Génétique

The International Multiple Sclerosis Genetics Consortium (IMSGC)



Science 2019: Genomic map

We analyzed genetic data of 47,429 MS and 68,374 control subjects and established a reference map of the genetic architecture of MS. This map includes 32 variants within the extended Major Histocompatibility Complex (MHC), 200 autosomal susceptibility variants outside the MHC, and one chromosome X variant.



Nature Communications 2019: Pathways

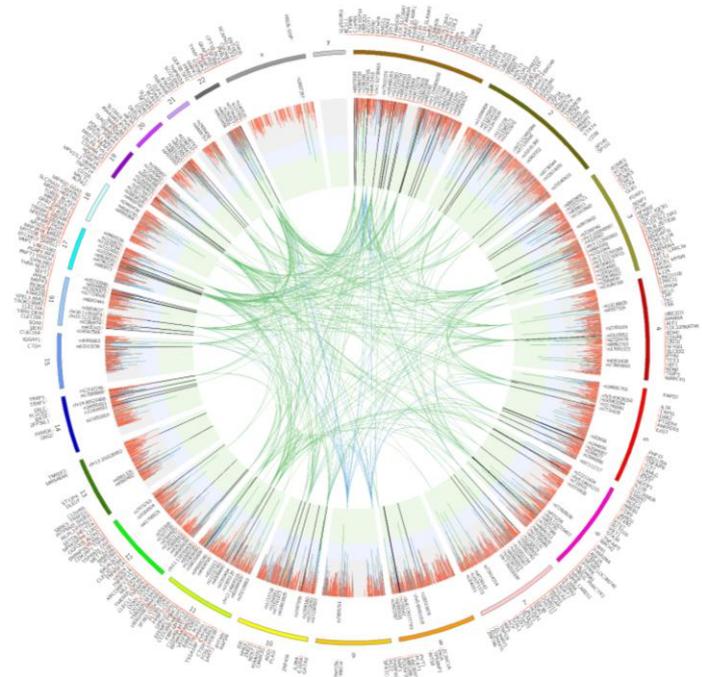
We conducted a cell-specific pathway analysis of the latest GWAS in MS. Our analysis identifies pan immune cell as well as cell-specific susceptibility genes in T cells, B cells and monocytes. Additionally, genotype-level data from 2,370 patients and 412 controls were used to compute intra-individual and cell-specific susceptibility pathways that offer a biological interpretation of the individual genetic risk to MS.

Update : 200+ Common Variants

- *“The Multiple Sclerosis Genomic Map: Role of peripheral immune cells and resident microglia in susceptibility”*
 - IMSGC et al. 2019 Science
- **Large-scale comparison**
 - 47,351 multiple sclerosis (MS) subjects
 - 68,284 control subjects
- **> 200 MS associated genomic regions**
 - 200 autosomal susceptibility variants
 - 1 chromosome X variant,
 - 32 independent associations in the extended MHC
 - Prioritize up 551 potentially associated MS genes

Establish a reference map of the 200+ genetic architecture of MS

Imsgc et al. Science 2019



3 “nouveaux “ Concepts vus par le détour de la génétique

1. *Rupture de la massification*
 2. « *Approches agnostique* » dans les données
 3. *Développements d'usage secondaires*
- Intentionnalité des données...*
- ...



Intégration des données pour l'émergence de nouveaux concepts et l'amélioration du **parcours de soins**

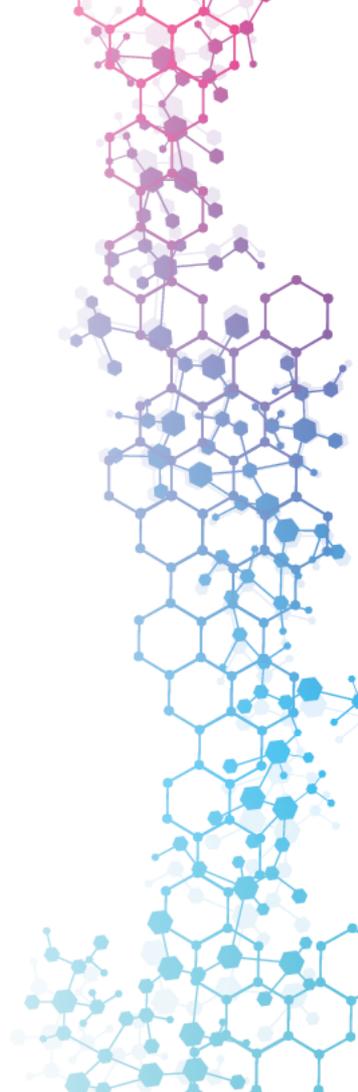
Accès aux données de Parcours de Soins...

Avec la participation scientifique du



Les données de parcours de Soins

- **Soins individuels**
 - o Le soin prodigué au patient est la raison d'être du partage de données entre soignants
 - o Portabilité (interopérabilité)
- **L'analyse du parcours de soins**
 - o Inter-établissements
 - Interopérabilité (comparabilité)
 - o Via prisme des données médico-administratives
 - Centralisées et centralisables
 - o Entrepôt de données de GHT, Coordination entre EDS de type Ouest Data Hub



Nantes Université in West of France



Technical Deployment -
public- private
partnership



2016 : From Local Biomedical Data Warehouses

- **Development et optimization:**
 - 1- *Data Governance*
 - 2- *Legal framework*
 - 3- *Technical Issue.*



Pr Marc CUGGIA



Périmètre plus large que les entrepôts de données issues du soin ré-exploités à des fins de recherche

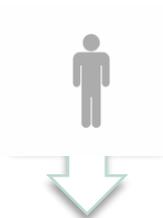
		RESPONSABILITÉ LÉGALE de la source de données : portée par le CHU de Nantes	
		OUI	NON
INTENTIONNALITÉ Données collectées à visée de soin	NON	<p>Entrepôt de données de santé Données médico-administratives et de soins</p> 	 <p>Système national des données de santé</p> <p>+ Projets « exceptionnels » : Météo France, SOS médecins, Argos...</p>
	OUI	<p>Cohortes du CHU de Nantes :</p> <p>BRUGADA, ICAN, EASY, VISIOCORT, EXAN, COVER, VALIDate, CoHPT, IT-DIAB, CORONADO...</p>	<p>Cohortes Nationales ou Internationales Bases de données externes</p>    



La Clinique des Données « Centre des données clinique »

- Service de CHU de Nantes
- Cellule Epidémiologie Clinique du CIC INSERM 1413
- Déploiement approuvé par la CNIL en juillet 2018
 - Dont Entrepôts de données Biomédicales issues du soin
 - 407 projets pris en charge.

Governance Matrix		Legal Resp. (CHUN)	
		Yes	No
Research Intentionality	Yes	Registry Cohorts	Nat. or Int. Databases
	No	BDW	SNDS



1,5 millions
Documented
Patients



540 millions
Structured
Data



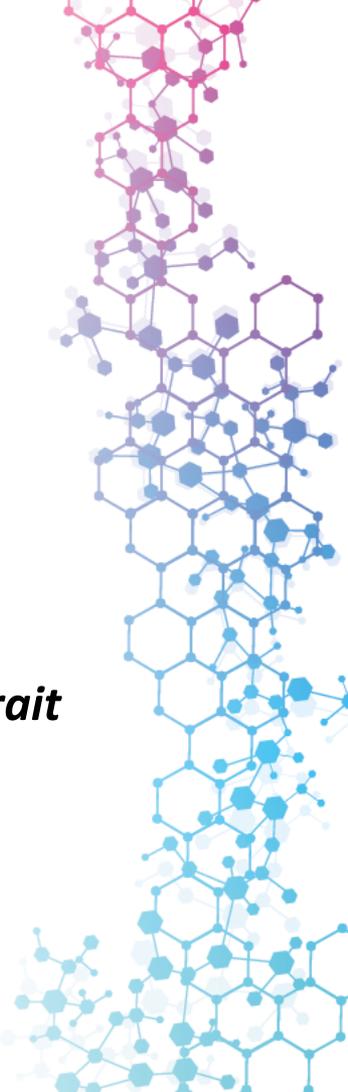
34 millions
Textual
Documents



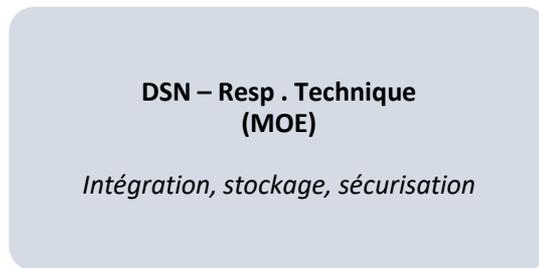
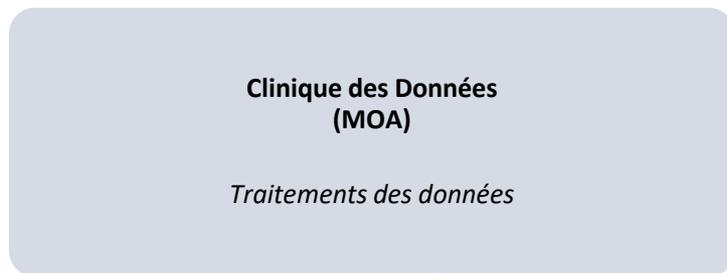


Les données : c'est le vecteur de transformation du XXIème siècle

« A terme, il n'y aura plus un projet de recherche en santé qui ne pourrait pas bénéficier d'une extraction de données d'entrepôt hospitalier. »



Une gouvernance Data qui doit être un catalyseur



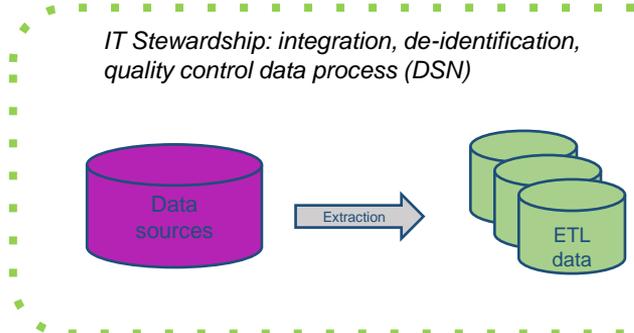
Cellule épidémiologie clinique du CIC 1413

Service de Santé Publique – PHU11

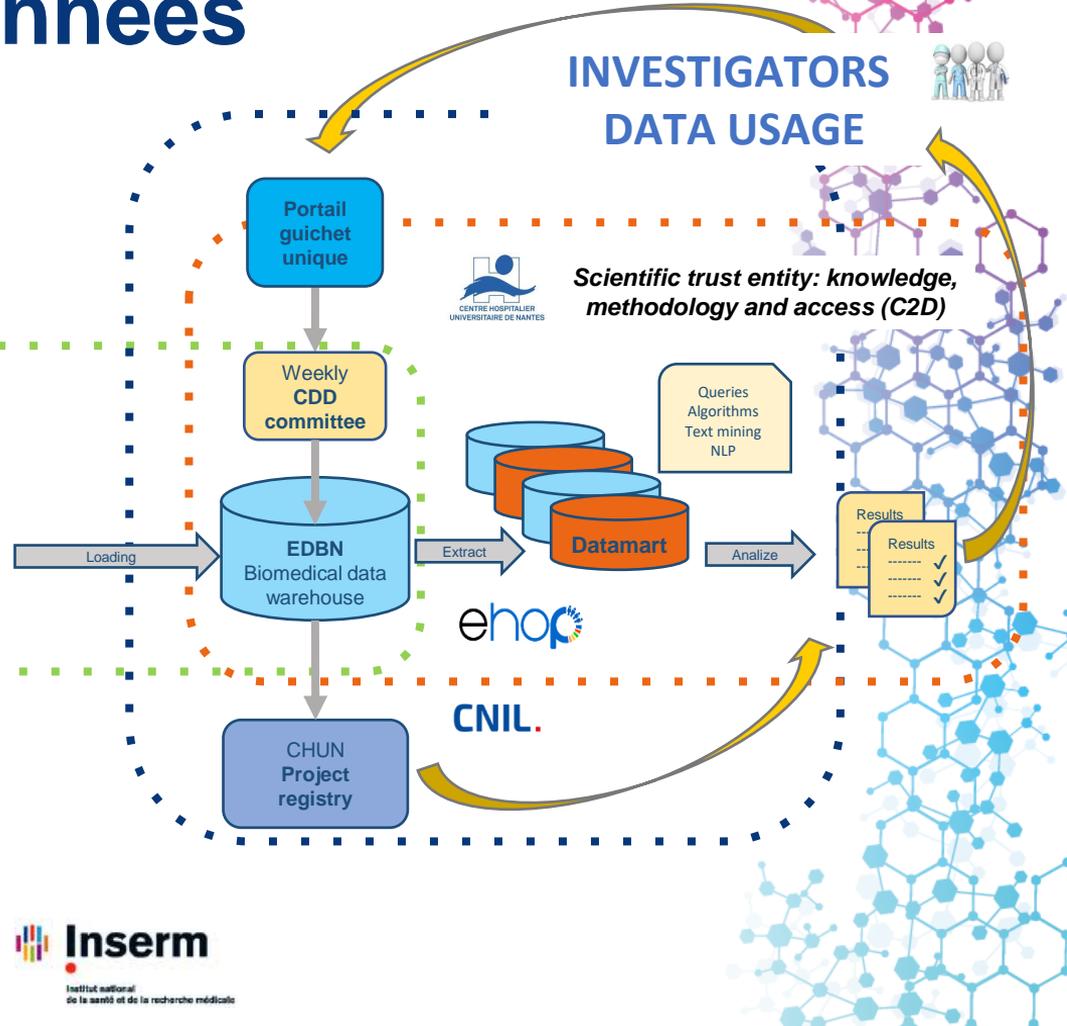
**Expertise = assistance à ...
MOE , MOA, Légal et financier
Voir à « l'urbanisation »**

La Clinique des Données

- Service de CHU de Nantes
- Cellule Epidémiologie Clinique - CIC INSERM 1413
- Déploiement approuvé par la CNIL en juillet 2018
- > 400 projets pris en charge, 120 publications, 8.3ETP



Regulatory framework: use and access policies, security and privacy (DRCI) : Governance & Legal Framework



Intégration des données pour l'émergence de nouveaux concepts et l'amélioration du parcours de soins

Un exemple

Illustrations

Projet SEVASAR



		RESPONSABILITÉ	
		OUI	NON
INTENTIONNALITÉ	OUI		
	NON	X	



#1

Identification des patients atteints du variant anglais du COVID-19

Dr P. Le Turnier
Service Maladies infectieuses et tropicales

- Projet national, problématique d'identification des patients avec le variant anglais hors CHU



Support
CdD

- Screening :

Identification population cible grâce à la recherche textuelle sur l'EDS, mise à disposition d'un listing d'IPP



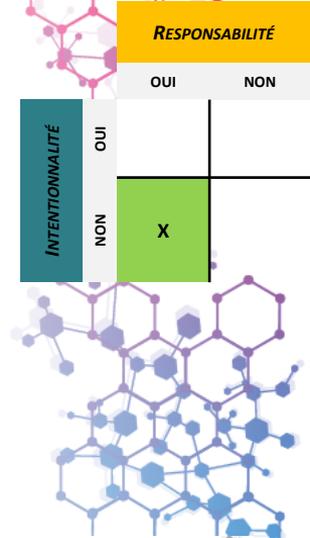
Illustrations

Projet SEVASAR



#1 Identification des patients atteints du variant anglais du COVID-19

Dr P. Le Turnier
Service Maladies infectieuses et tropicales



3. Recherche eHOP par interface graphique : résultats

Contexte: Ensemble de données "Globale"

Statistiques de l'étude

- Nombre de patients: 1 481 838
- Nombre de documents: 32 250 999
- Afficher les données sensibles:

Etude: [dropdown]

Lien: http://ehopoppprdr/ehop/main?id_study=1

Demandeur: Aucun

Créateur: Aucun

Dates d'accès: accès sans restriction

Type: prescreening

Accès global: oui

Autorise l'accès aux données sensibles: oui

Vue mat.: non

Recherche rapide

344 Patient(s)

1 112 Document(s)

17 (0.05 %) Opposé(s) réutilisation

17 (0.05 %) Opposé(s) recontact

Requête exécutée en 2.53s.

Concepts

#row	Id Pat	Id Pat Etude	Age actuel	Sexe	Décédé ?	Actions
1	3281	-	59 ans	F	-	
Période séjour						
UF/UM Document(s) Sign. doc Age Patient/Document Actions						
Le 08/02/2021 2 document(s)						
	UF 2072	Synthèse patient	08/02/2021	58 ans		Afficher
	UF 2083	Compte rendu examens biologique	08/02/2021	58 ans		Afficher
2	3973	-	67 ans	F	-	
Période séjour						
UF/UM Document(s) Sign. doc Age Patient/Document Actions						
Du 31/03/2021 au 15/04/2021 3 documents(s)						
	UM 2088	PMSI 5935771 1	01/04/2021	66 ans		Afficher
	UM 3710	PMSI 5935771 2	01/04/2021	66 ans		Afficher

Illustrations

Projet SEVASAR



#1

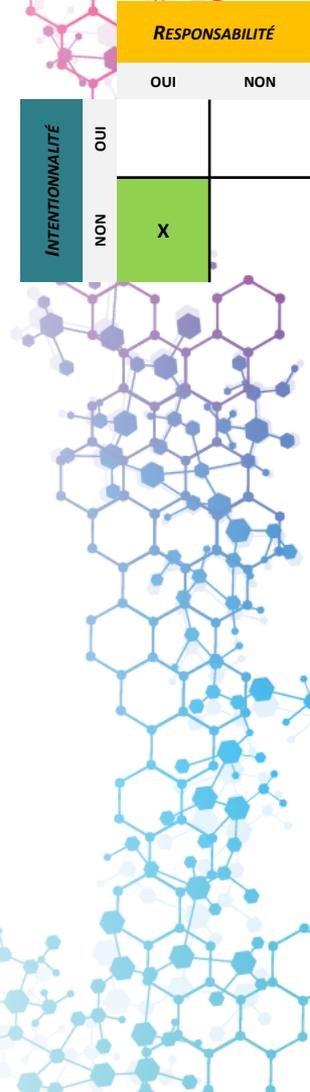
Identification des patients atteints du variant anglais du COVID-19

Dr P. Le Turnier
Service Maladies infectieuses et tropicales



Recherche eHOP par interface graphique : CR également accessible

Nature:	Ecouvillon nasopharyngé
	RECHERCHE DE VIRUS PAR BIOLOGIE MOLECULAIRE
Recherche du SARS CoV-2 (COVID 19):	
Test moléculaire rapide IDNOW COVID 19 (ABBOTT)	
Résultat:	POSITIF
	Présence d'ARN viral compatible avec une excrétion virale significative; patient à considérer comme contagieux.
SARS CoV2: recherche de la mutation E484K	
Kit VirSniP SARS CoV-2 spike E484K (TIB MOLBIOL)	
Mutation E484K:	NON
SARS CoV2: recherche de la mutation N501Y	
Kit VirSniP SARS CoV-2 spike N501Y (TIB MOLBIOL)	
Mutation N501Y:	OUI
Interprétation:	Détection d'un variant dit anglais
CONCLUSION:	Mise en évidence d'un coronavirus SARS CoV-2 (COVID 19)



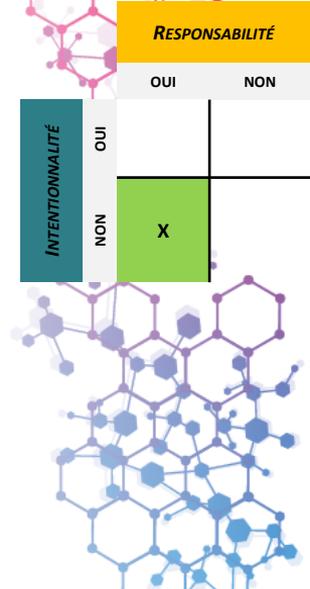
Illustrations

Projet SEVASAR

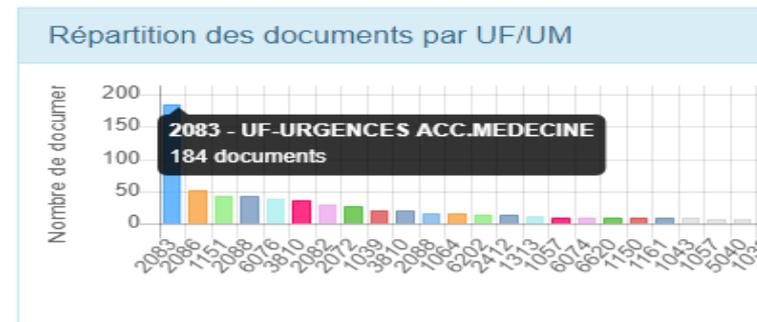
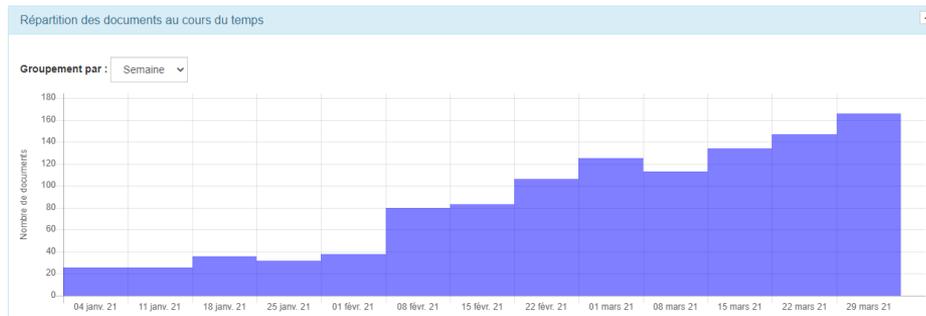


#1 Identification des patients atteints du variant anglais du COVID-19

Dr P. Le Turnier
Service Maladies infectieuses et tropicales



Recherche eHOP par interface graphique : Distribution temps et UF/UM



Les “carrefours de données” de Santé ?

Articulation entre modes d'accès nationaux – accompagnement local

- Circuit d'information individuelle systématisée
- Proximité de la production des données : « En santé, un data scientist est d'abord un expert du contexte dans lequel naît la donnée avant d'être un expert des méthodes de traitement de ces données”

Médiation par un tiers-expert

Structure locale labellisée

- « les données parleraient d'elles-mêmes »
- “Jamais seul face aux données”
- HUB local - Clinique des données
– Centres de Données Cliniques



Un réservoir d'innovation...

- TAL, données synthétiques, apprentissage fédéré, chiffrement homomorphe, Pilotage par la multidata

Enjeu de transformation – nouvelle épidémiologie de données.

On parle très vite en millions... Données en « vie réelle » - qui requièrent plus de méthodes



PLÉNIÈRE - REAL-WORLD DATA

Intégration des données pour l'émergence de nouveaux concepts et l'amélioration du parcours de soins

Pr. Pierre-Antoine Gourraud, Nantes

CHU Nantes, PHU 11 : Clinique des données, INSERM, CIC 1413, Nantes, France

Nantes Université, INSERM, CRT21 , Nantes, France

pierre-antoine.gourraud@univ-nantes.fr



 Nantes
Université



Intégration des données pour l'émergence de nouveaux concepts et l'amélioration du parcours de soins

Un bonus

Avec la participation
scientifique du



La confiance des patients

« Quand on analyse des données , il n'y a plus de justification à faire courir un risque de ré-identification aux patients.

« L'analyse de données personnelles dans un but d'analyse non personnelle devient un enjeu éthique »



Actualité Légale

« Naïf de la donnée »

Rôle d'alerte et d'éveil de la CNIL



https://www.cnil.fr/sites/default/files/atoms/files/referentiel_entrepot.pdf

Page 14 : Exportation de données hors de l'entrepôt et hors des espaces de travail

SEC-EXP-1 A l'exception des données relatives aux procédures de ré-identification SEC-REI-1 à SEC-REI-3, **seuls des jeux de données anonymes peuvent faire l'objet d'une exportation hors de l'entrepôt ou d'un espace de travail**. Le processus d'anonymisation doit produire un jeu de données conforme **aux trois critères définis par l'avis du G29 n° 05/2014** ou à tout avis ultérieur du CEPD relatif à l'anonymisation. **Cette conformité doit être documentée et démontrable**. À défaut, si ces trois critères ne peuvent être réunis, une étude des risques de ré-identification devra être menée et documentée.

SEC-EXP-2 Les exports de données doivent être soumis à **la validation préalable d'un responsable** afin d'en avaliser le principe, notamment au regard de l'exigence SEC-EXP-1.

SEC-EXP-3 Les exports doivent faire l'objet d'une surveillance automatique ou manuelle par un opérateur spécialisé afin d'en vérifier le caractère anonyme. Dans le cas où cette surveillance est automatique, tout export identifié comme non conforme doit faire l'objet d'une remontée d'alerte et d'une mise en quarantaine dans l'entrepôt, puis doit être vérifié manuellement par un responsable spécifiquement formé et spécifiquement habilité.

SEC-EXP-4 Les systèmes mis en place dans l'entrepôt relatifs à la production d'indicateurs et au pilotage stratégique de l'activité d'un établissement de santé ne doivent permettre que des restitutions anonymes, y compris en tenant compte des fonctionnalités de filtrage et de sélection de ces restitutions. Ce processus de restitution doit être conforme aux trois critères définis par l'avis du G29 n° 05/2014 ou à tout avis ultérieur du CEPD relatif à l'anonymisation. Cette conformité doit être documentée. À défaut, si ces trois critères ne peuvent être réunis, **une étude des risques de ré-identification devra être menée et documentée**.

Comparison of FAMD

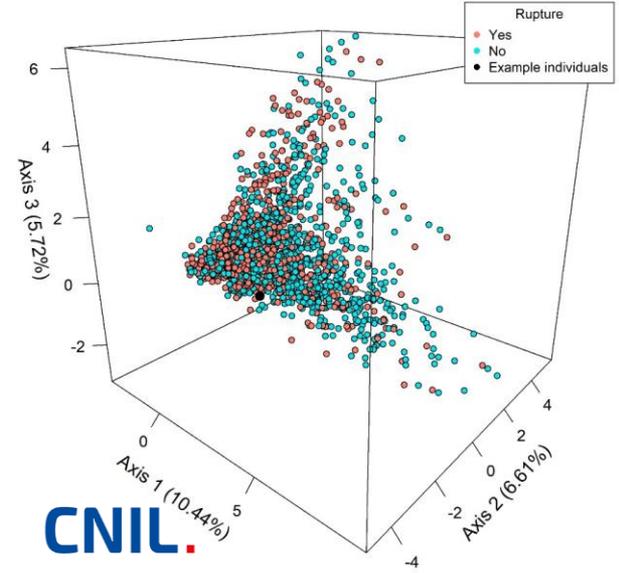
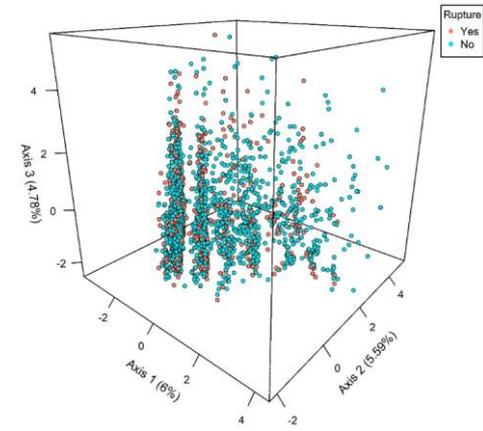
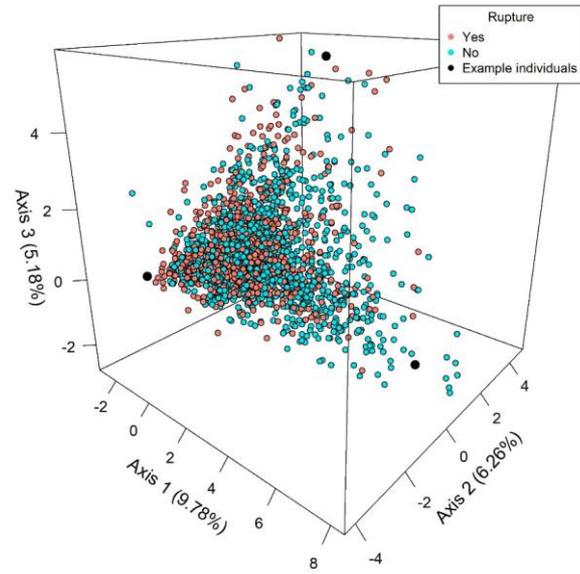


AVATARS

(synthetic version of data the original dataset)

SENSITIVE DATA
(Original Pseudonymous data)

SIMULATED DATA
(mathematically simulated representation of the original dataset)



SENSITIVE DATA (Original Pseudonymous data)	
Risk of re-identification	+++
Preparation/curation	+
Informative Value	++++

SIMULATED DATA (mathematically simulated representation of the original dataset)	
Risk of re-identification	0
Preparation	++++
Informative Value	+

AVATARS (synthetic version of data the original dataset)	
Risk of re-identification	0
Preparation/curation	+
Informative Value	++++



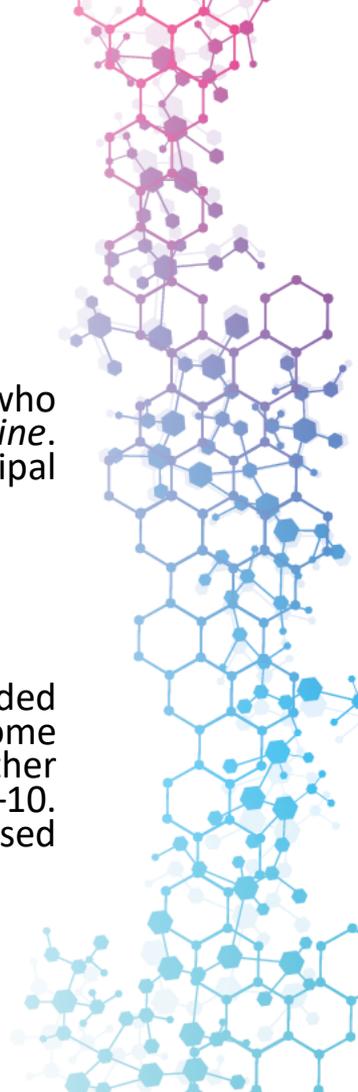
Uses Cases: 2 Typical Biomedical (tabular) datasets

(AIDS) Clinical trial

- The AIDS dataset includes 2139 patients and 26 variables for HIV-infected patients who participated in a clinical trial published in 1996 in the *New England Journal of Medicine*. The clinical trial had four arms and was analyzed by Hammer et al. (1996)²⁹. The principal endpoints used were survival and a 50% drop in CD4+ cell counts.

Wisconsin Breast Cancer Diagnosis (WBCD): prediction issue

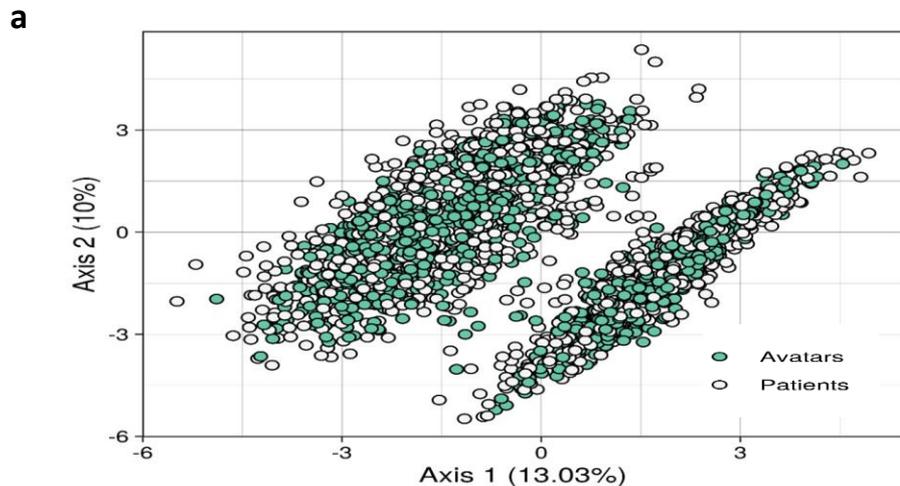
- The WBCD dataset comprises 683 observations and 10 variables. It can be downloaded from the University of California Irvine machine-learning repository³⁰. The outcome corresponds to the tumor severity: benign (n=444) versus malignant (n=239). The other nine features are built from imaging specific annotations and are graduated from 1–10. Feature selection (F-score computation) and a support vector machine (SVM) were used to predict the severity of a patient's breast cancer diagnosis as per Akay et al. (2009)³¹.



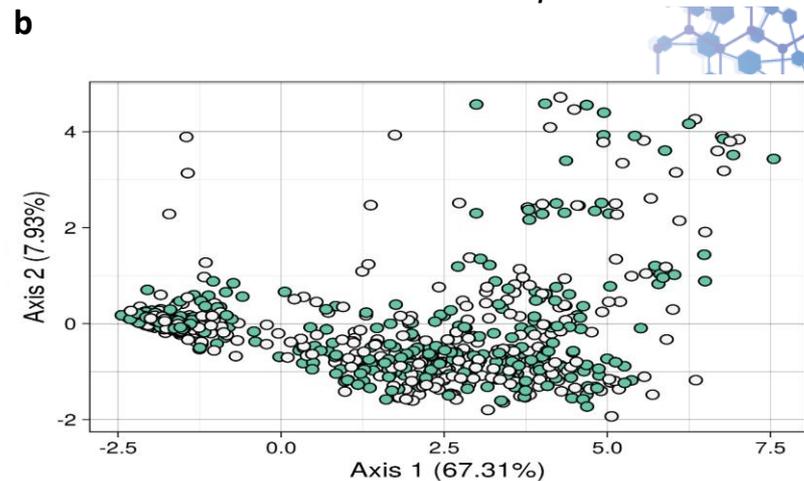
Results 1: Datasets conservations

- Similar multidimensional representation

AIDS - Clinical Trial



WBCD – cancer prediction



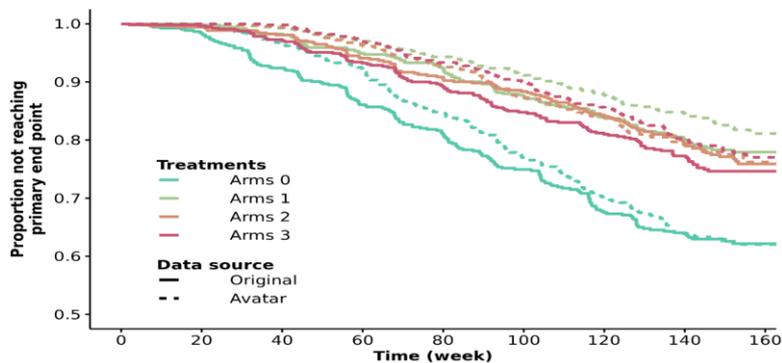
Results 2: Statistics conservations

- Similar results of the main analysis (if exsiting) associated to the data set

AIDS - Clinical Trial

Avatars $p= 1,5 E-9$ vs $p= 1,2E-8$

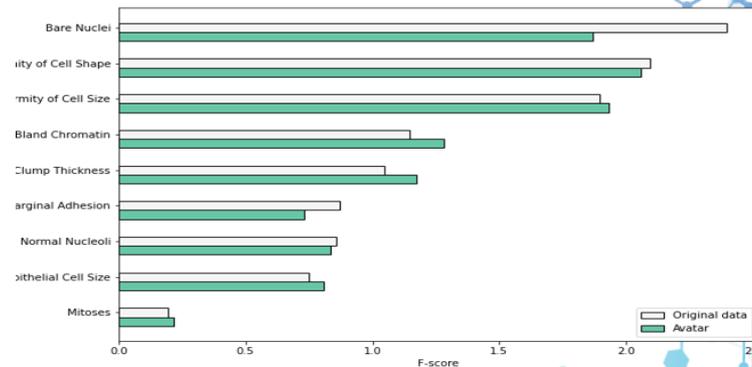
c



WBCD – cancer prediction

AUC= 99,84 vs AUC= 99,46

d



Results 3: Statistics conservations

- Similar results of the main analysis (if existing) associated to the data set

AIDS - Clinical Trial

Avatars $p= 1,5 E-9$ vs $p= 1,2E-8$

	Hazard Ratio	Pr(> z)	lower .95	upper .95
Avatar arms1	0.4	<0.001	0.31	0.51
Original arms1	0.49	<0.001	0.39	0.63
Avatar arms2	0.52	<0.001	0.41	0.67
Original arms2	0.52	<0.001	0.41	0.67
Avatar arms3	0.5	<0.001	0.39	0.63
Original arms3	0.59	<0.001	0.47	0.73

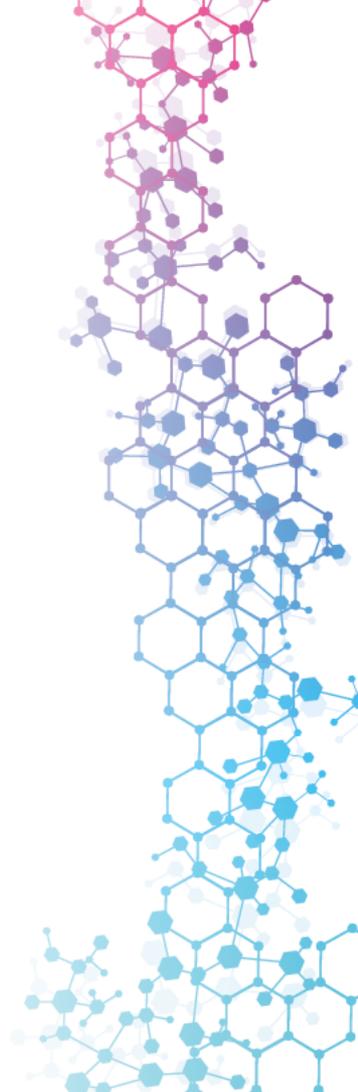
WBCD – cancer prediction

AUC= 99,84 vs AUC= 99,46

	avatar	original
acc	99.024390	92.682927
auc	99.186864	99.940312
npv	97.368421	90.140845
ppv	100.000000	94.029851
sens	98.473282	94.736842
spec	100.000000	88.888889

? Est-ce que la valeur statistique des données compte vraiment ?

Peut-être pas en premier...



La confiance des patients

«Quand on analyse des données , il n’y a plus de justification à faire
courir un risque de ré-identification aux patients. »

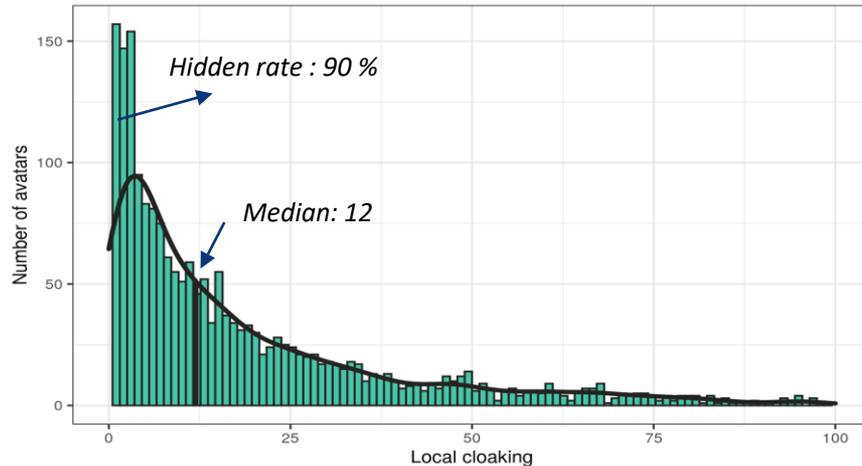


Results I: Privacy matters most

- How well protected sensitive observations are ?

12 avatars in average

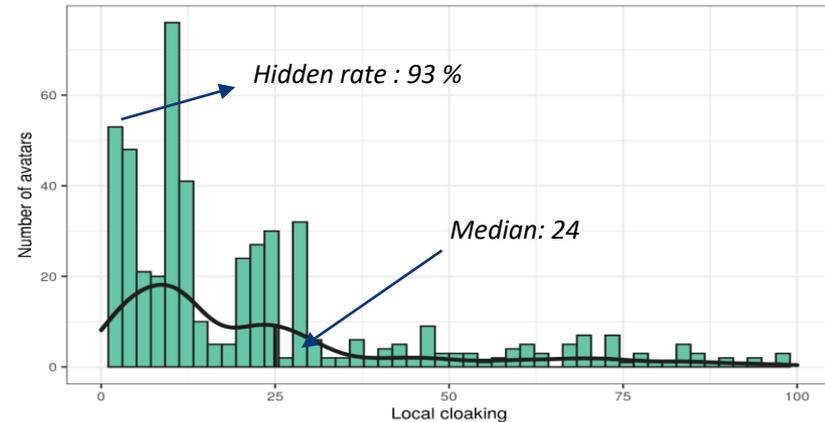
– 10 % with their avatars at their side



AIDS - Clinical Trial

25 avatars in average –

5,7% with their avatars as the closest



WBCD – cancer prediction

Merci de votre attention

