

Génération de données synthétiques centrées sur le patient, aucune raison de risquer la ré identification dans l'analyse des données biomédicales

Amphi Schwarz, IMT 11:30 - 12:00

Pr Pierre-Antoine GOURRAUD, PhD MPH

Vendredi 9 Février 24, Toulouse

Professeur des Universités & Praticien-Hospitalier
 CHU Nantes, PHU 11 : Clinique des données, INSERM,
 CIC 1413, Nantes Université, INSERM, CR2TI

pierre-antoine.gourraud@univ-nantes.fr



- **COI :**

PA Gourraud est le fondateur (2008) (www.methodomics.com) et le co-fondateur de Big data Santé (2018) . Il est consultant et/ou intervenant pour de grandes entreprises pharmaceutiques ou de dispositifs médicaux. Ses activités sont toutes traitées par une contractualisation universitaire ou hospitalière (*AstraZeneca, Biogen, Boston Scientific, Cook, Docaposte, Edimark, Ellipses, Elsevier, Janssen, IAGE, Lek, Methodomics, Merck, Mérieux, Octopize, Sanofi-Genzyme*). PA Gourraud est administrateur bénévole des mutuelles d'assurances AXA (2021). Il n'a aucune activité de prescription de médicaments ou dispositifs médicaux. Il ne perçoit pas de rémunération complémentaire

- **COI :**

PA Gourraud is the founder of Methodomics (2008) and the co-founder of Big data Santé (2018). He consults for major pharmaceutical companies, and start-ups, all of which are handled through academic pipelines (AstraZeneca, Biogen, Boston Scientific, Cook, Docaposte, Edimark, Ellipses, Elsevier, Janssen, IAGE, Lek, Methodomics, Merck, Mérieux, Octopize, Sanofi-Genzyme). PA Gourraud is a volunteer board member at AXA not-for-profit mutual insurance company (2021). He has no prescription activity with either drugs or devices. He receives no wages from these activities.

Le problème de données de Santé

La circulation des données ...



La solution du jeu à la Toulousaine... *une combinaison de mouvements et de passes, mettant la vitesse d'exécution, Elle permet de créer des opportunités d'attaques.*

Les Données Personnelles de santé **NE CIRCULENT PAS**

... Car soumises :

- + au RGPD
- + au référentiel EDS CNIL,
- + à la pression des patients,
- + à la pression des soignants,
- + à la pression des établissements.

« 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. » Référentiel CNIL EDS



CNIL.
COMMISSION NATIONALE
INFORMATIQUE & LIBERTÉS

Un contexte:

La Clinique des données du CHU de NANTES, EDS, faciliter les usages secondaire de la donnée issue du soin

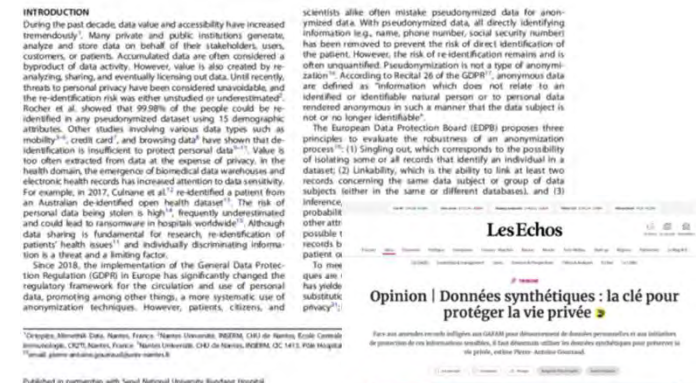
Mettre en qualité la donnée de santé pour la rendre réutilisable

La solution : Simuler des données anonymes à (très haute) valeur informative

Comment faire des données anonymes démontrables ?

Guillaudeau et al 2023 Nature Digital Medicine

<https://www.nature.com/articles/s41746-023-00771-5>



Triple Actualité des données Synthétiques

1. Actualité éthique : Avis « Plateformes de données de santé : enjeux d'éthique »

- **Avis 143** : Comité consultatif national d'éthique (CCNE)
- **Avis 5** : Comité national pilote d'éthique du numérique (CNPEN)
- Page 31 - Recommandation 8
 - <https://www.ccne-ethique.fr/publications/avis-143-du-ccne->



COMITÉ NATIONAL PILOTE
D'ÉTHIQUE DU NUMÉRIQUE

www.ccnpen.fr
COMITÉ CONSULTATIF NATIONAL D'ÉTHIQUE
POUR LES SCIENCES DE LA VIE ET DE LA SANTÉ

2. Actualité réglementaire : RGPD pour les usages données personnelles (de santé)

- (Vs.) **Données publiques (Open Source) - Confidentialité**
- **Contrôle des finalités, l'identité des usagers des données,**
- **Référentiel EDS** publié par la CNIL : « **Données Anonymes** »



CNIL
COMMISSION NATIONALE
INFORMATIQUE & LIBERTÉS

3. Actualité scientifique : Usages des données Synthétiques

- Lauréat de l'**Appel à projet Entrepôt de Données de Santé (EDS)** : « ODH 2.0 »
 - Projet technique Données Synthétiques Anonymes Avatars avec Octopize
- **“Patient-centric synthetic data generation, no reason to risk re-identification in biomedical data analysis”**

**OUEST
DATA
<HUB**

npj | Digital Medicine

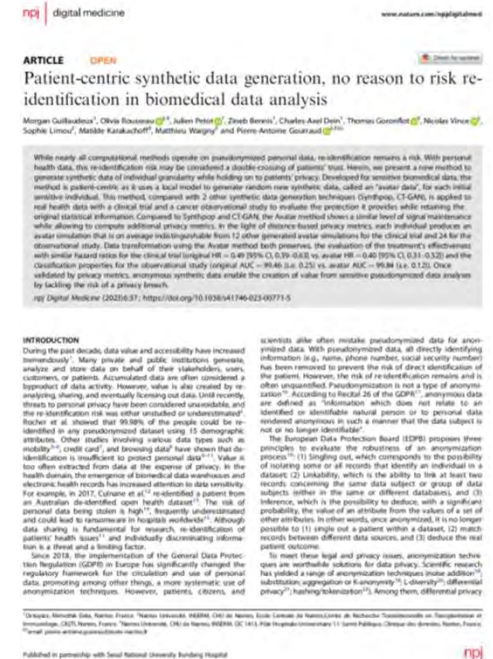
Guillaudeau et al 2023 Nature Digital Medicine

<https://www.nature.com/articles/s41746-023-00771-5>

La Méthode

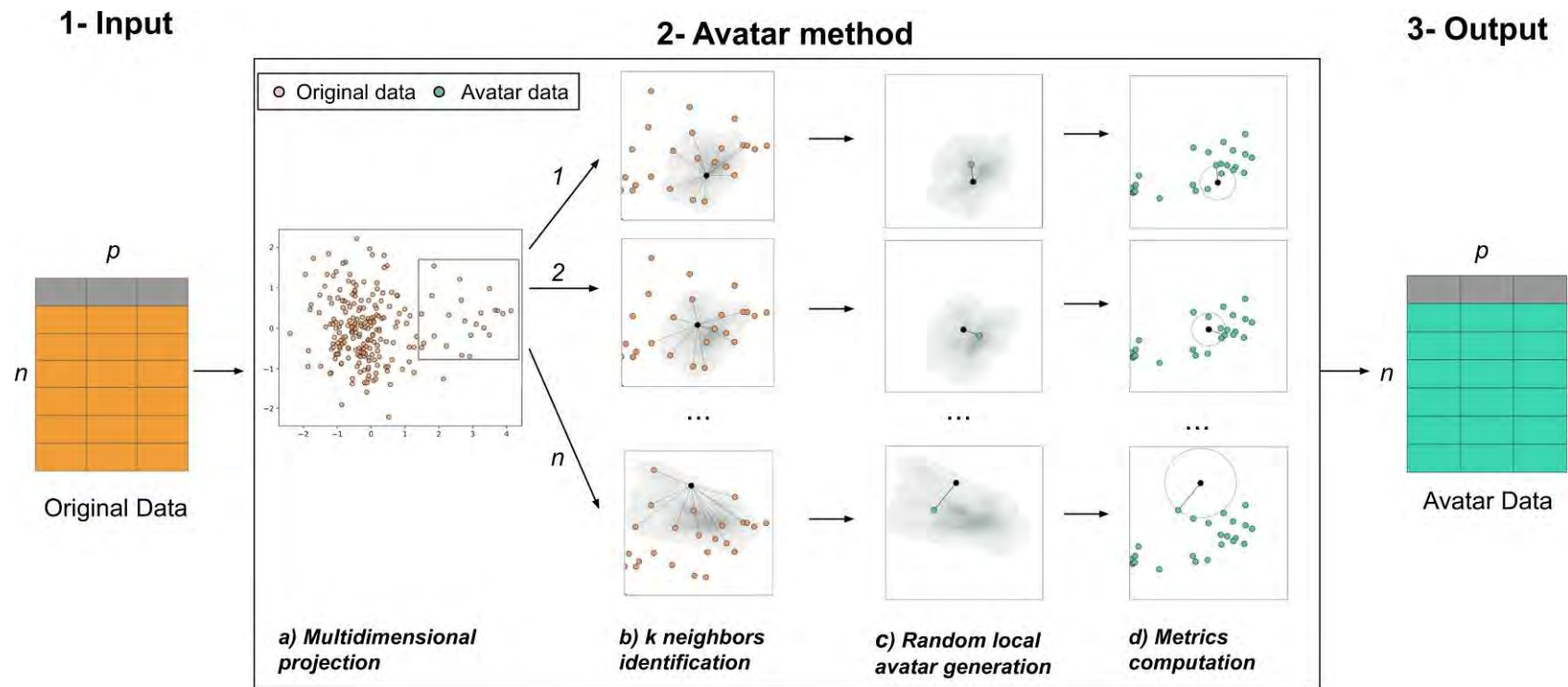
La Clinique des données du CHU de Nantes favorise les usages secondaires des données de santé.

Comment faire des données anonymes démontrables ?
 Guillaudeux et al 2023 Nature Digital Medicine
<https://www.nature.com/articles/s41746-023-00771-5>



Methodes: Données Synthétiques anonymes

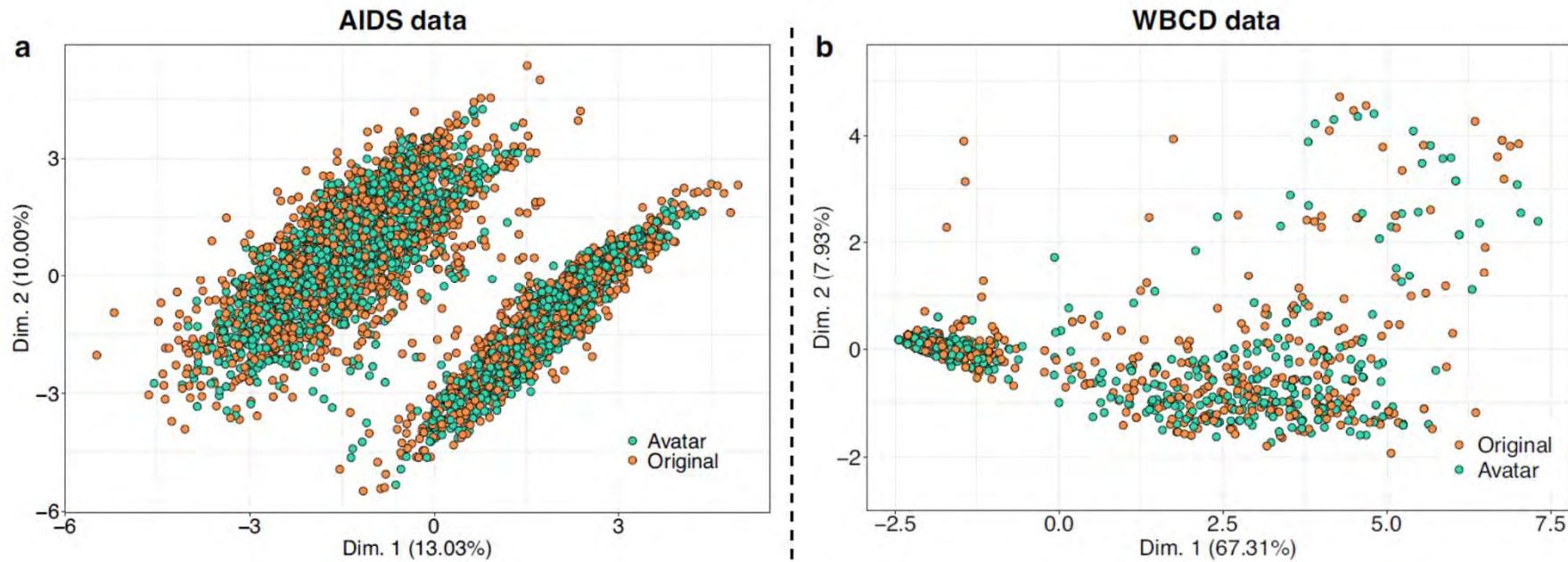
- Comment ça marche ? Un modèle privé non-paramétrique à usage unique



Guillaudeau et al 2023 Nature Digital Medicine - <https://www.nature.com/articles/s41746-023-00771-5> Figure 5

Results 1: Datasets conservations

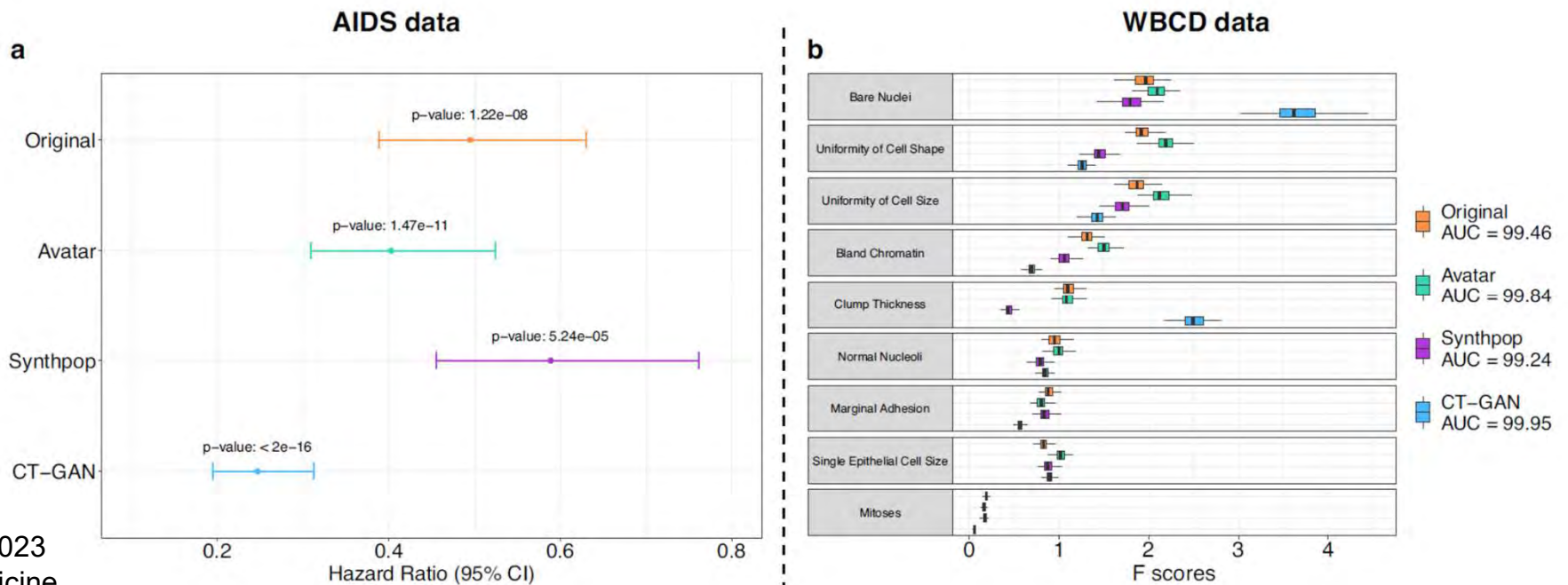
- Représentation multidimensionnelles semblables



Guillaudeau et al 2023
Nature Digital Medicine
<https://www.nature.com/articles/s41746-023-00771-5>
Figure 1 A&B

Results 2: Statistics conservations

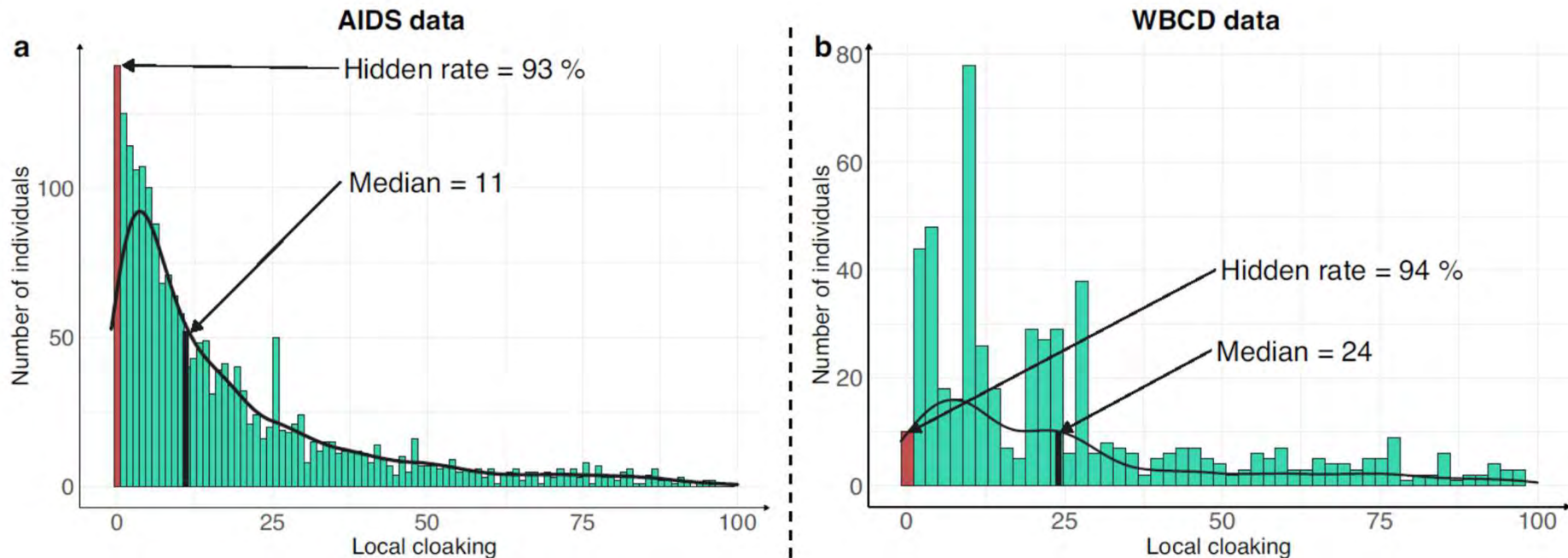
- Résultats d'ensemble similaires à l'analyse principale (si existante) + autres méthodes



Guillaudeau et al 2023
Nature Digital Medicine
<https://www.nature.com/articles/s41746-023-00771-5> Figure 2 A&B

Resultats 3 : Est ce que la valeur statistique des données compte vraiment en premier ?

- Les métriques qui attestent que les Avatars sont des données synthétiques anonymes au sens CNIL - c'est-à-dire au sens démontrable.





Guillaudeau et al 2023 Nature Digital Medicine <https://www.nature.com/articles/s41746-023-00771-5> Figure 3A & 3B voir 3C & 3D pour plus de garanties

Application dans le contexte des entrepôts des données de santé

La Clinique des données du CHU de Nantes favorise les usages secondaires des données de santé.

Un cas d'usage : Publication & Science Ouverte

- **Situation:** le Pr FX B et le DR S J demandent à la Clinique des données de l'aide pour répondre aux reviewers dernier obstacle pour la publication.
 - **“Does PaCO₂ correction have an impact on survival of patients with chronic respiratory failure and long-term non-invasive ventilation?”**
- **Commentaire 1 – rev 1 :** Etude avec manque de puissance statistique, hétérogénéité Clinique, possibilité de partager les données (en vue de méta analyse) ?
 - *Données pseudonymisées – Contrôle des usages* 
- **Solution :** Données synthétiques anonymes
 - Application du référentiel “EDS” 
- **Autre exemple :** JNNP en open source (licence MIT)
 - <https://jnnp.bmj.com/content/92/2/122>
 - <https://github.com/ICAN-aneurysms/RIA-predict/tree/master/notebooks>



Démonstration

The Avatar interface consists of several key components:

- Loaded Data:** A table showing input data with columns for 'Data' (e.g., age, sexe, v1_poids) and 'Data Type' (e.g., float, category).
- Avatarization parameters:** A form to configure the process, including 'Number of neighbors to take into account' (set to 4) and 'Number of components to use for neighbor selection' (set to 20).
- Privacy Metrics:** Displays metrics like Median Rule Range (MRR) and Local Anonymity (LAP) for both 'Original' and 'Avatar' data.
- Utility Metrics:** Displays metrics like Median Rule Range (MRR) and Local Anonymity (LAP) for both 'Original' and 'Avatar' data.
- Correlations:** Heatmaps comparing the correlation structure of the original data versus the generated avatars.
- Output:** A final table showing the generated avatars, with a green checkmark indicating successful completion.

Application dans le contexte des entrepôts des données de santé :

Projet RHU Médecine de précision pour les maladies chroniques

Exemples en maladie chronique : Sclérose en plaques et Transplantation



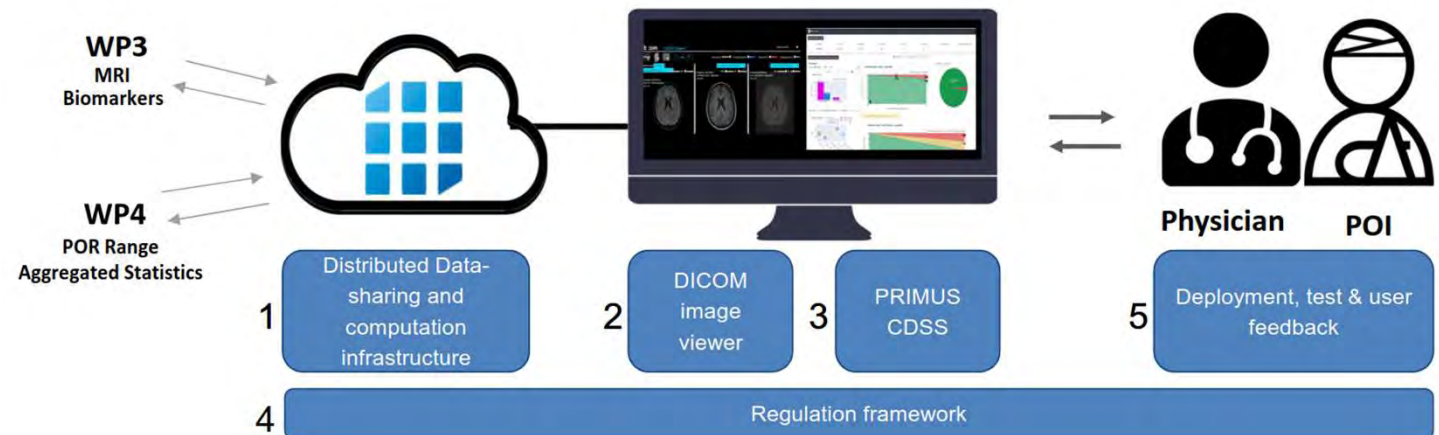
RHU KTD-Innov - Ref projet: IA-17-RHUS-0010 : Ce travail bénéficie d'une aide de l'Etat gérée par l'Agence Nationale de la Recherche au titre du PIA, portant la référence IA-17-RHUS-0010

RHU PRIMUS - Ref projet: ANR-21-RHUS-0014 : Ce travail bénéficie d'une aide de l'Etat gérée par l'Agence Nationale de la Recherche au titre du 3e PIA, intégré à France 2030 portant la référence ANR-21-RHUS-0014

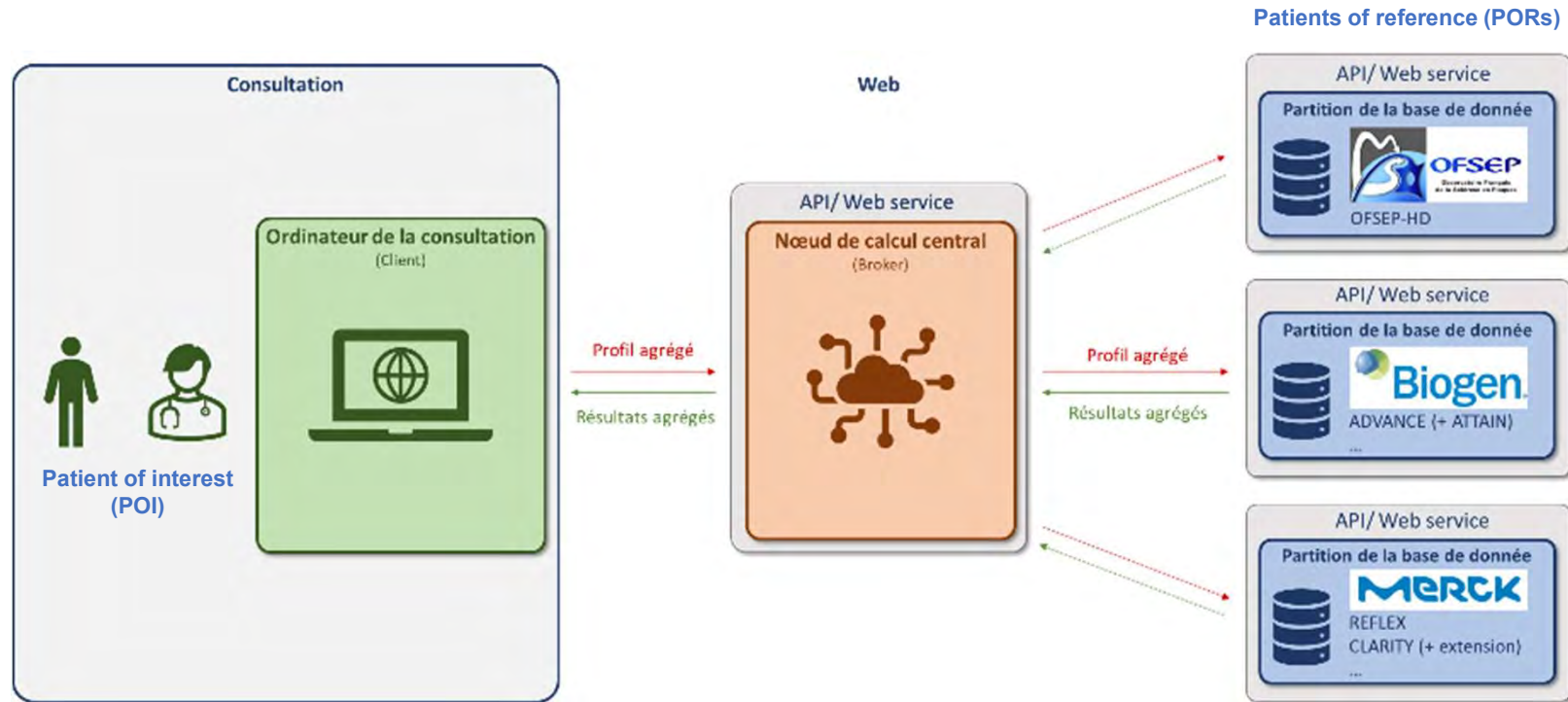
Projet de RHU Primus en Sclérose en plaques

Outils d'aide à la décision médicale (**CDSS**) qui permet au neurologue, lors d'une consultation avec son patient (Patient Of Interest – **POI**) de visualiser (**viewer**) la projection de l'évolution de la maladie du POI par rapport à son imagerie (**IRM**) et des algorithmes collaboratifs (**IA**) qui vont sélectionner des patients de références (**POR**) dans une base de données de références (**BD Ker PRIMUS**) construite à partir des essais cliniques de l'industrie et des données de vie réelles de plus de 12 000 individus en accord avec la réglementation (**RGPD**) et le Medical Device Regulation (**MDR**). Testé dans un essai clinique en prévision d'un **marquage CE**.

« Le bon traitement
pour la bonne
personne au bon
moment et prévenir la
survenue de
handicaps »



Projet de RHU Primus en Sclérose en plaques



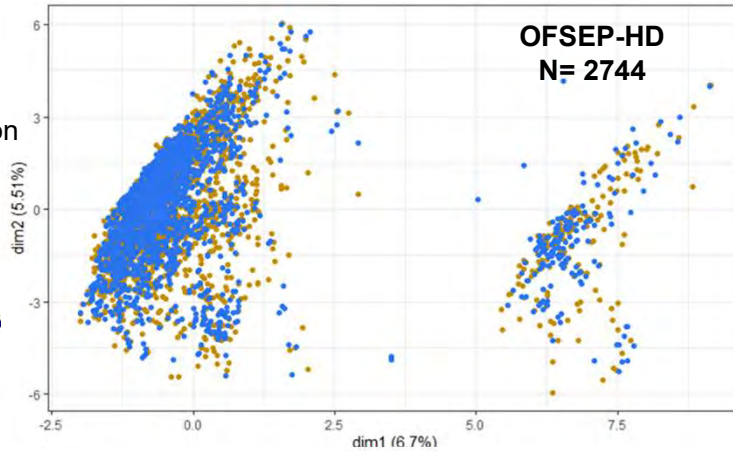
→ Base de données de référence distribuée: 1 partition = 1 propriétaire de données

Le logiciel d'aide à la décision du RHU PRIMUS, Dr. Stanislas DEMUTH, laboratoires des Prof. Pierre-Antoine GOURRAUD & Prof. Jérôme DE SEZE, Assises de l'OFSEP 2023

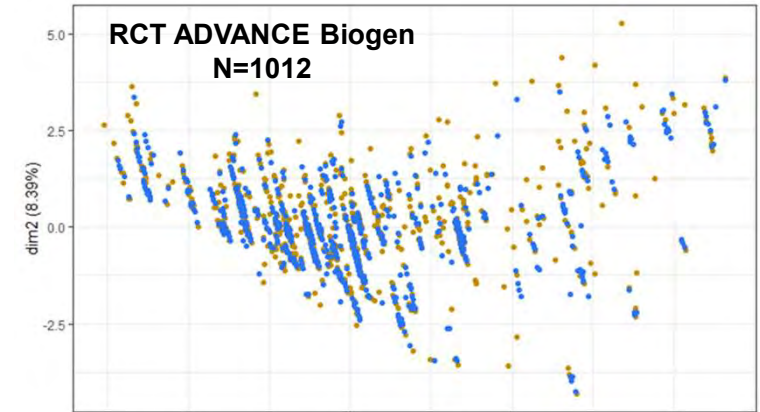
RHU PRIMUS: 4 reference synthetic Datasets MS Patients – Strategic alliance



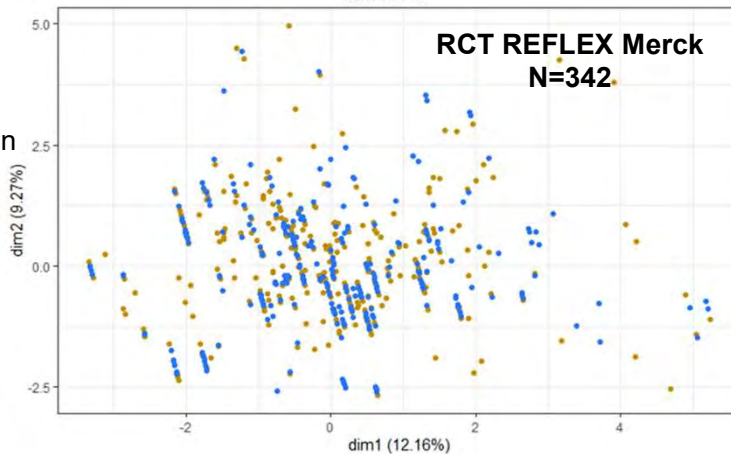
K = 5
NCP = 56 with categories reduction
HR = 83.0 %
Median LC = 3
HM = 0.04



K = 5
NCP = 18 with categories reduction
HR = 86.7 %
Median LC = 3
HM = 0.02



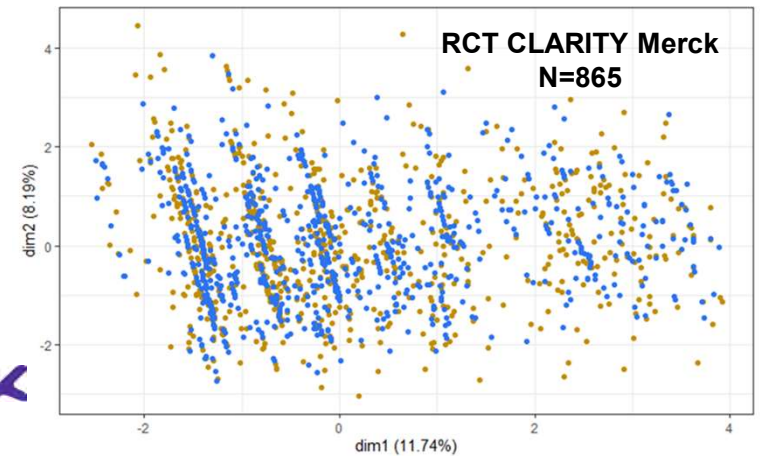
K = 5
NCP = 16 with categories reduction
HR = 82.1 %
Median LC = 2
HM = 0.04



veracity

- Avatar
- Real

K = 5
NCP = 19 with categories reduction
HR = 82.4 %
Median LC = 2
HM = 0.04

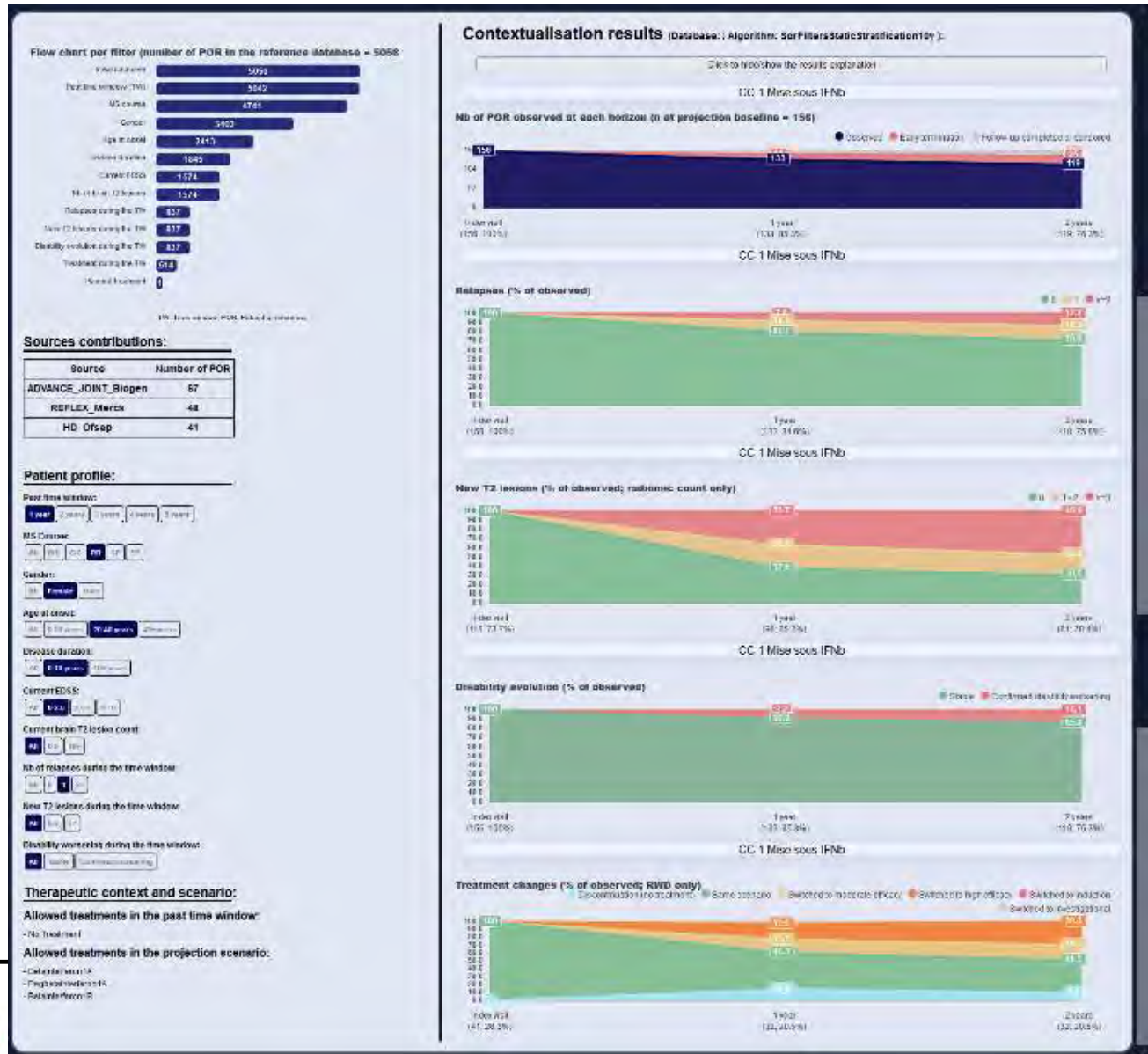


PRIMUS Alpha : Projection de toute la population de référence

Cas clinique 1 :
 ♀ 26a, première poussée, McDo+

Mise sous IFN bêta

Le logiciel d'aide à la décision du RHU PRIMUS, Dr. Stanislas DEMUTH, laboratoires des Prof. Pierre-Antoine GOURRAUD & Prof. Jérôme DE SEZE, Assises de l'OFSEP 2023

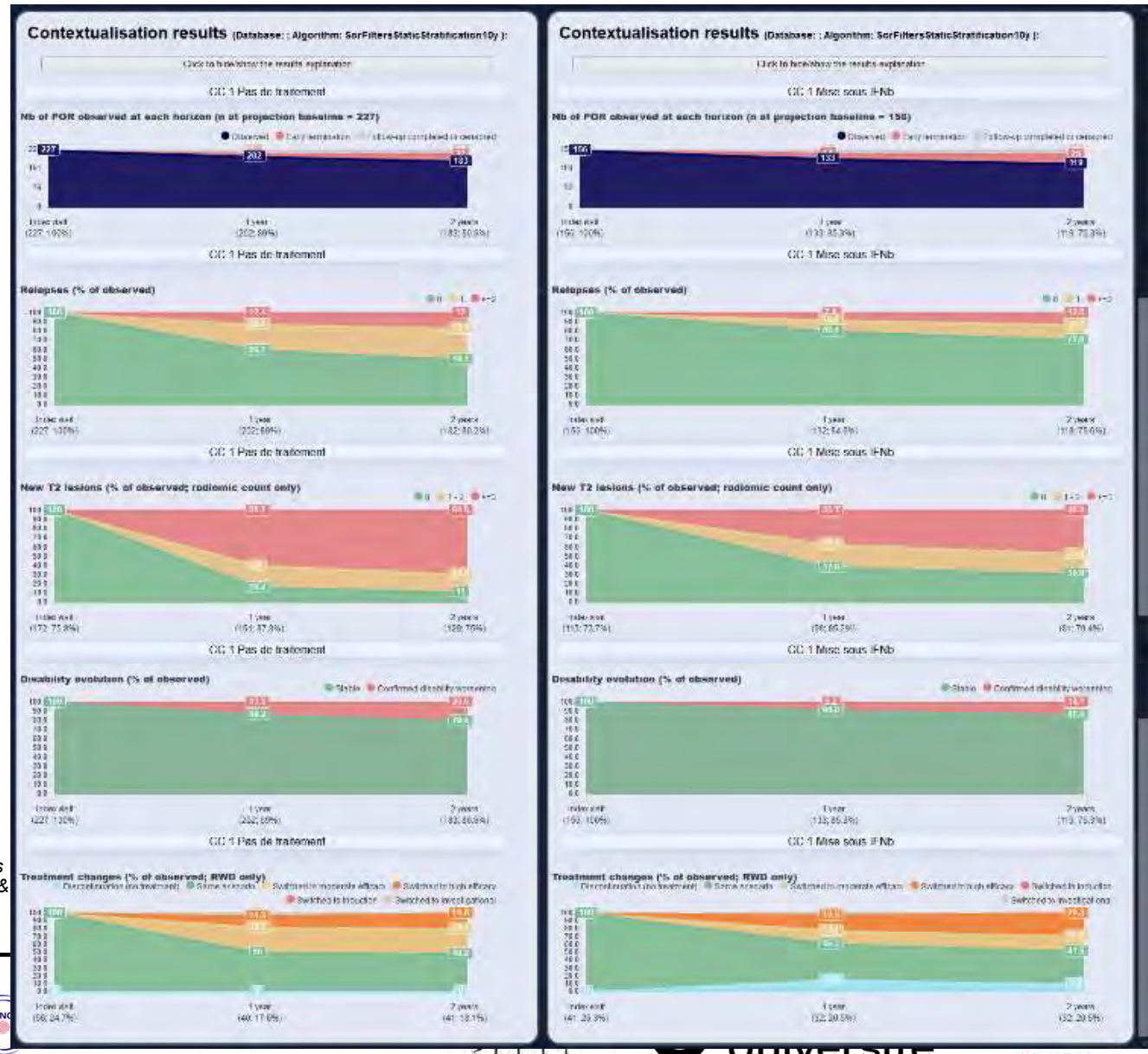


PRIMUS Alpha : Projection de toute la population de référence

Cas clinique 1 :
 ♀ 26a, première poussée, McDo+

Mise sous IFN bêta

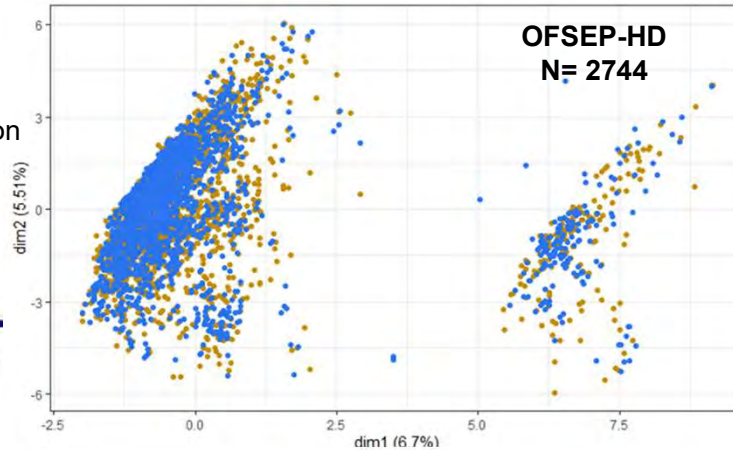
Le logiciel d'aide à la décision du RHU PRIMUS, Dr. Stanislas DEMUTH, laboratoires des Prof. Pierre-Antoine GOURRAUD & Prof. Jérôme DE SEZE, Assises de l'OFSEP 2023



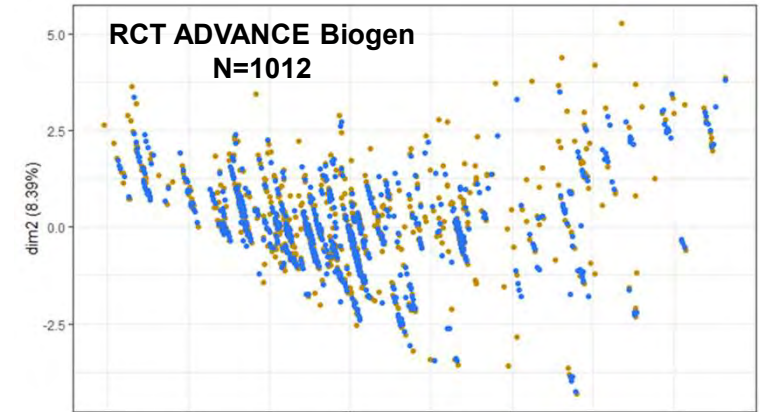
RHU PRIMUS: 4 reference synthetic Datasets MS Patients – Strategic alliance



K = 5
NCP = 56 with
categories reduction
HR = 83.0 %
Median LC = 3
HM = 0.04



K = 5
NCP = 18 with
categories reduction
HR = 86.7 %
Median LC = 3
HM = 0.02



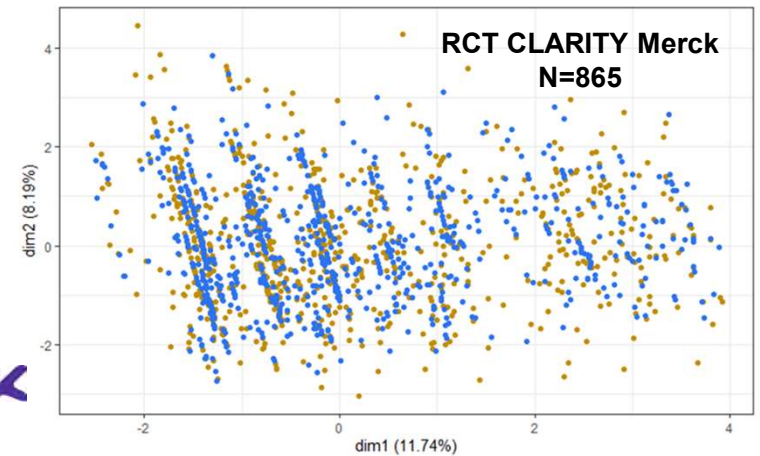
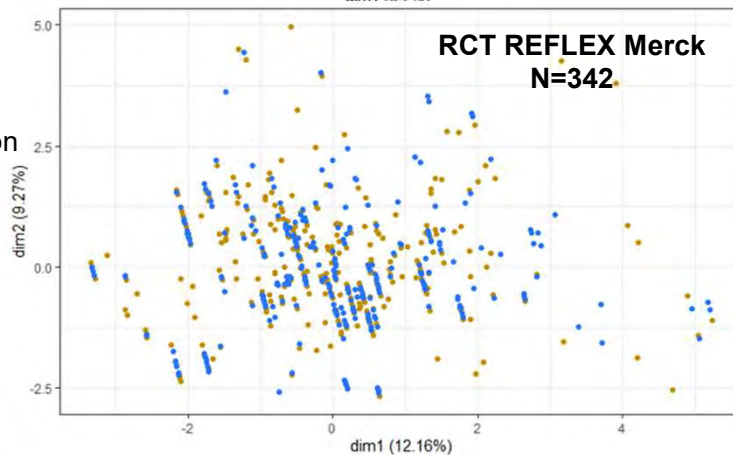
veracity

- Avatar
- Real

K = 5
NCP = 19 with
categories reduction
HR = 82.4 %
Median LC = 2
HM = 0.04

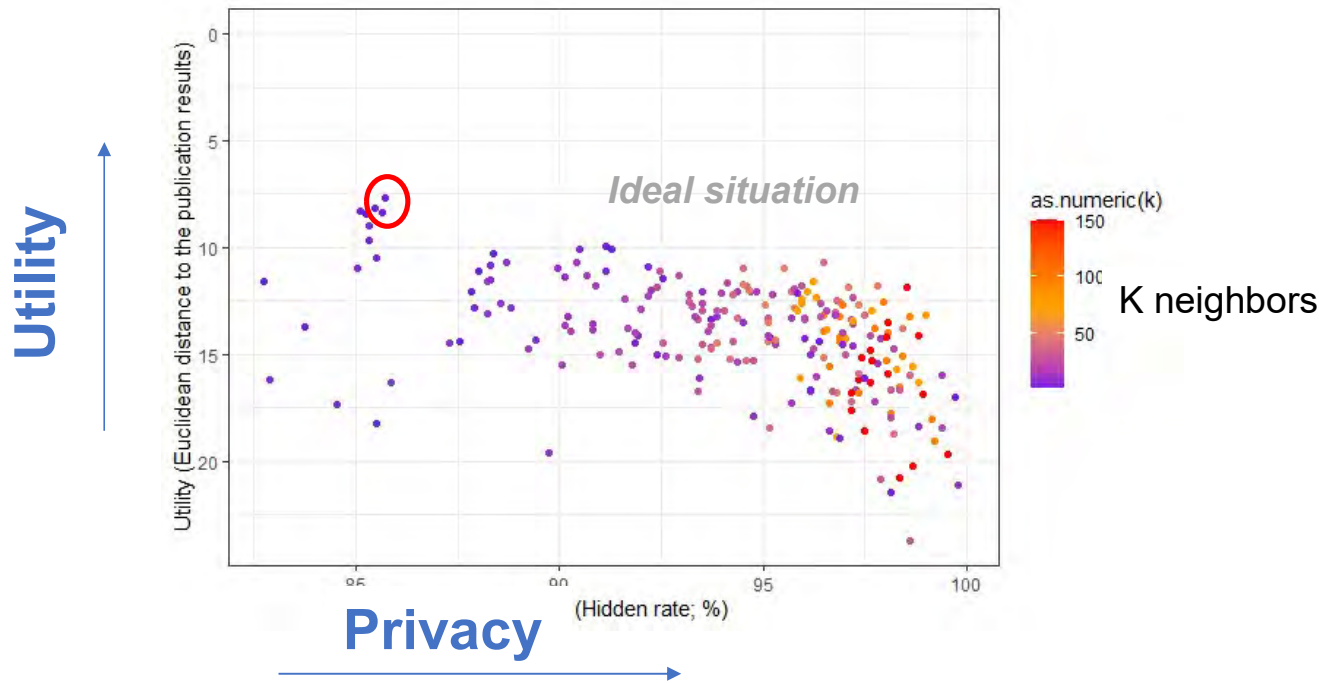


K = 5
NCP = 16 with
categories reduction
HR = 82.1 %
Median LC = 2
HM = 0.04

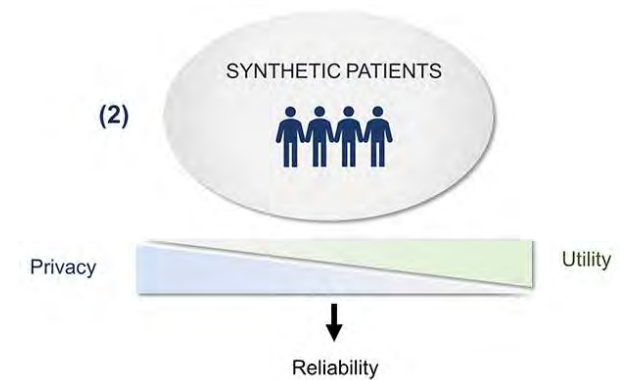


I) Avatarization of RCT data for multi-outcome analysis

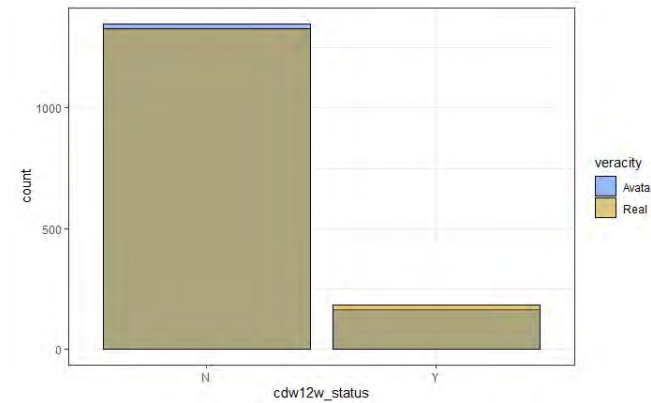
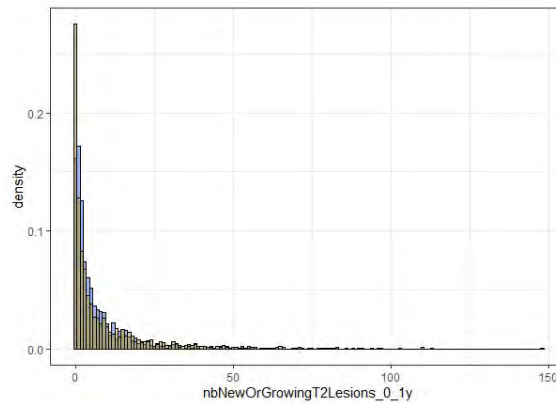
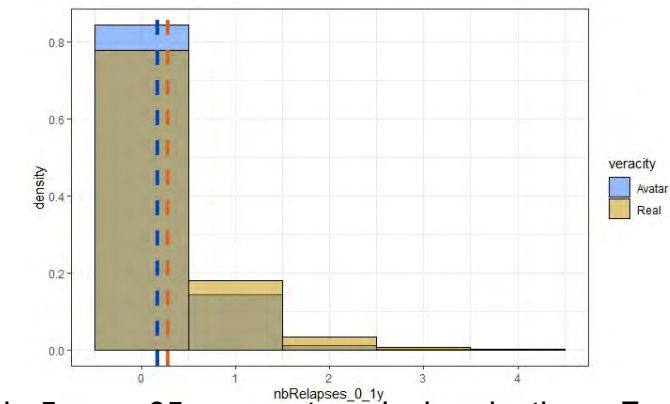
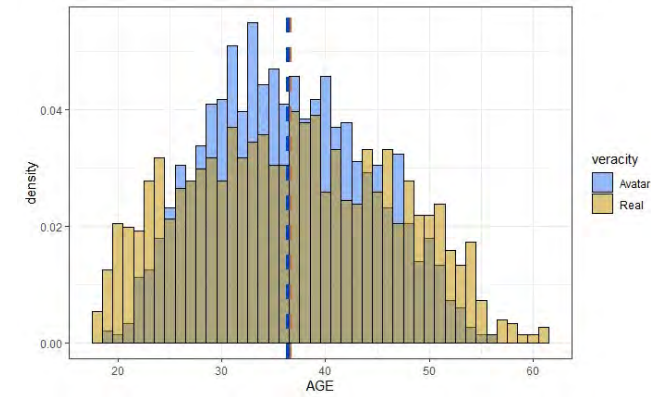
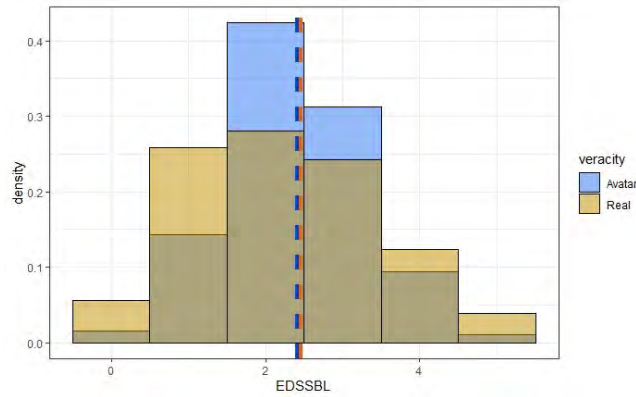
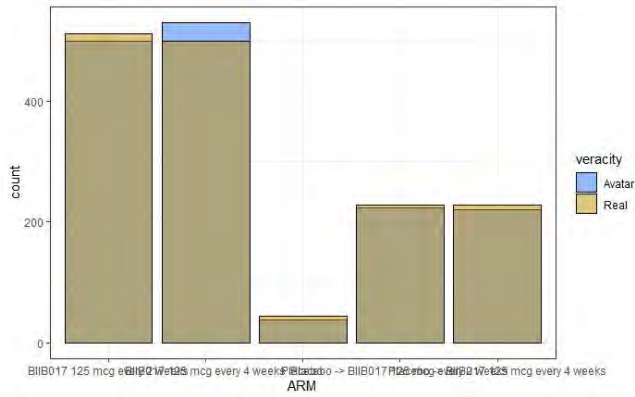
The best out
Utility and
Privacy



1512 observations,
25 variables (17 continuous, 8 categorical)
After dummy variables encoding: 62 columns



I) Avatarization of RCT data for multi-outcome analysis



k=5, ncp=65, use_categorical_reduction = F , seed = 1

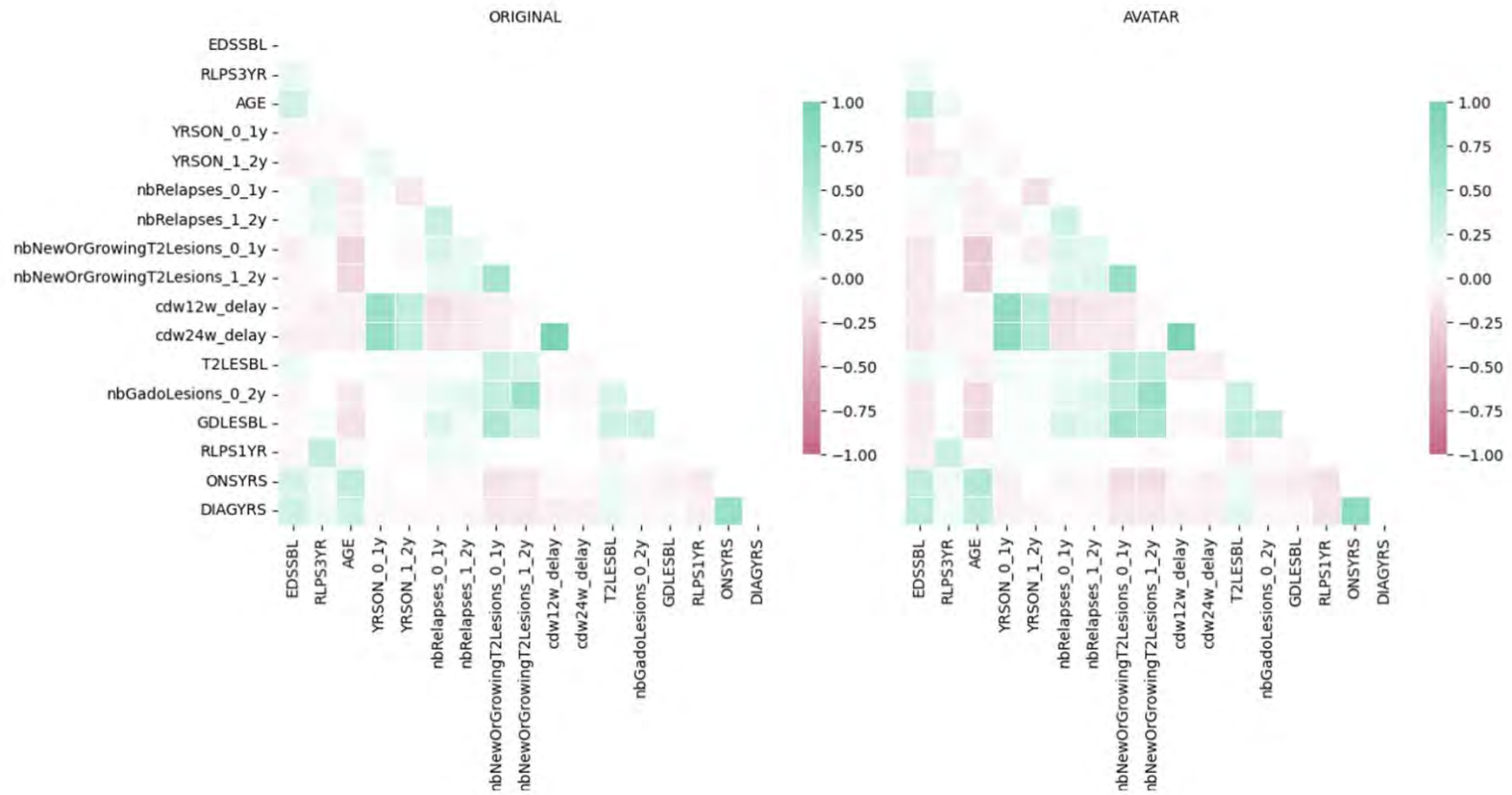
Real data Avatar

General Utility

I) Avatarization of RCT data for multi-outcome analysis

k=5, ncp=65, use_categorical_reduction = F , seed = 1

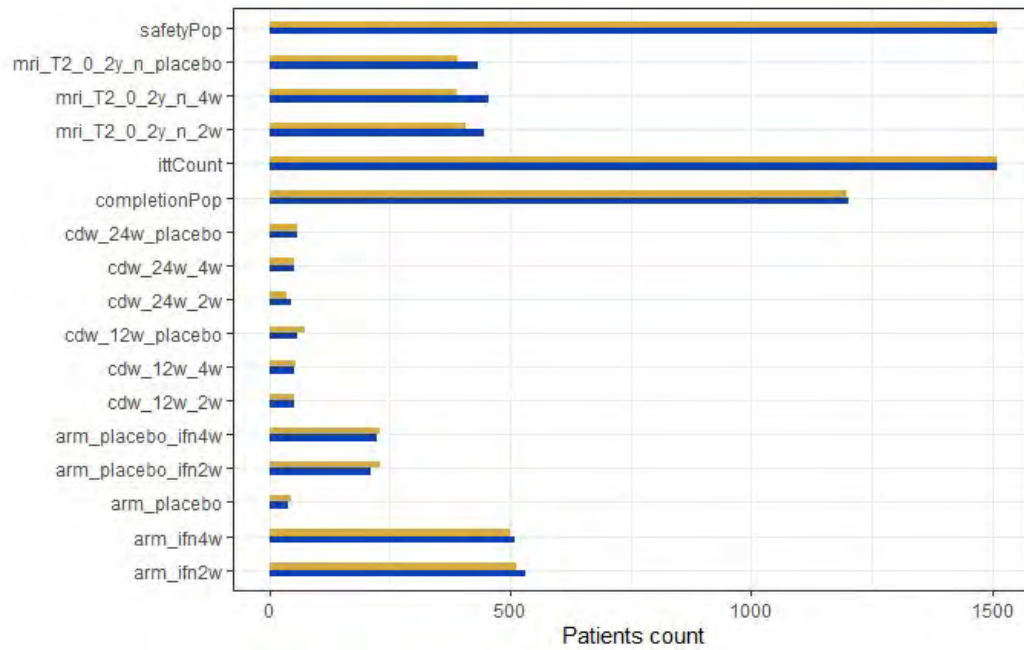
General Utility



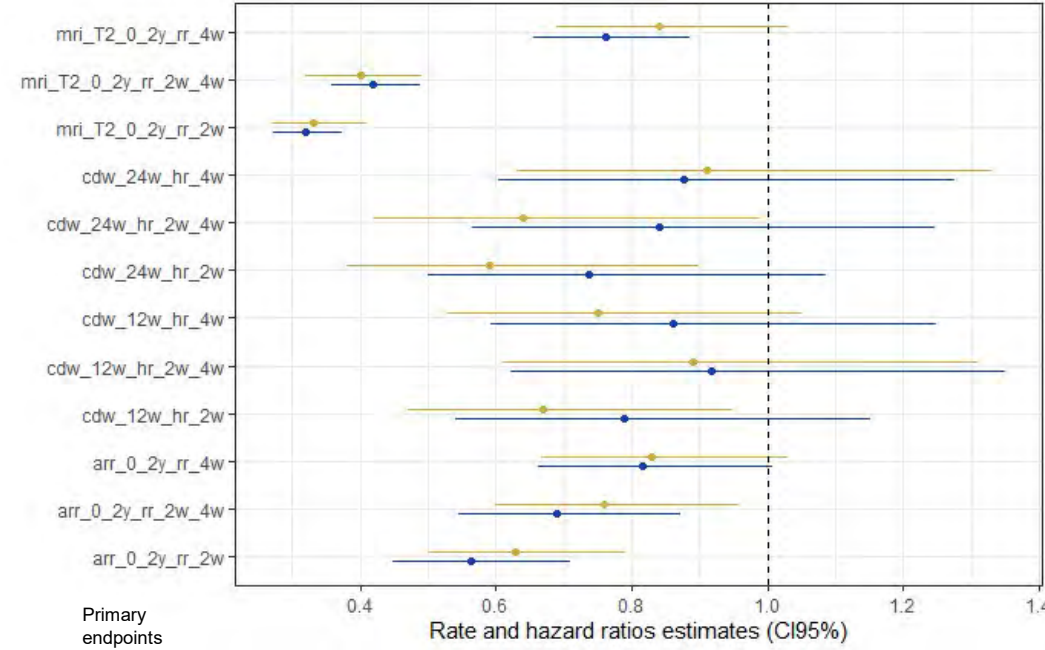
I) Avatarization of RCT data for multi-outcome analysis

k=5, ncp=65, use_categorical_reduction = F , seed = 1

Absolute counts



Association statistics and their 95%CI



General Utility



Real data



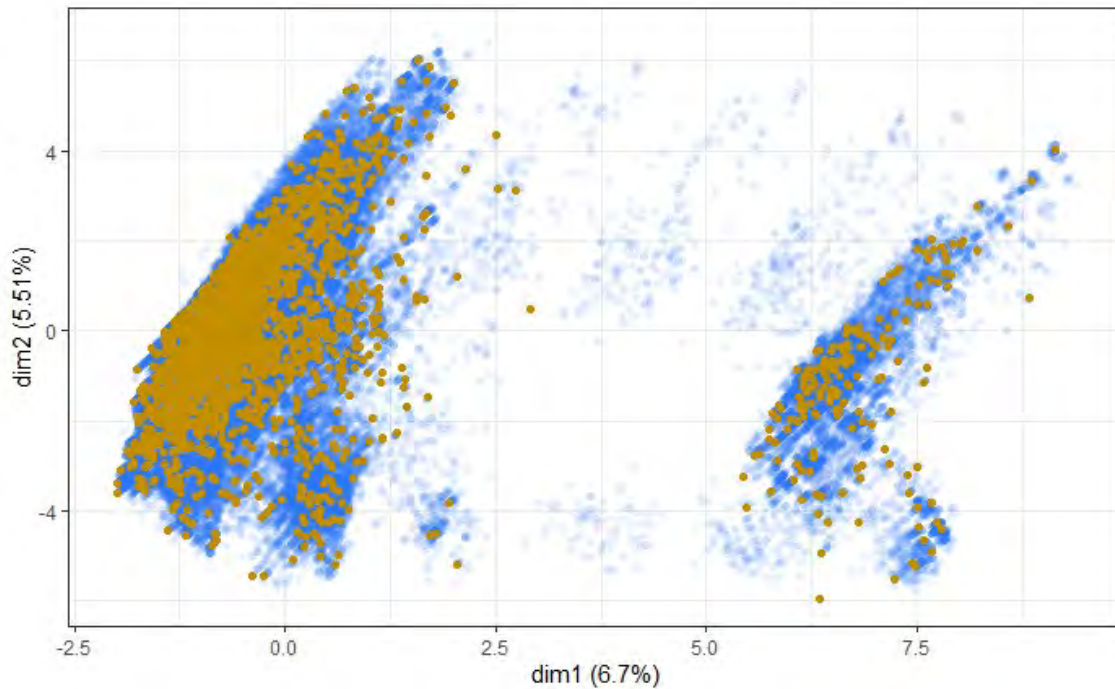
Avatars

Specific Utility

II) Avatarization of OS + RCT data for filter-based similar patient data visualization (independent data points)

FAMD projection

OFSEP-HD simplified to 2744 patients



From 1-1 to 1-100 Synthetic data generation

Avatarization parameters:

- Avatar
- Real
- K = 5
- NCP = 56 (100% of inertia)
- Column weights: All = 1
- Use Categorical Reduction: True

Metrics:

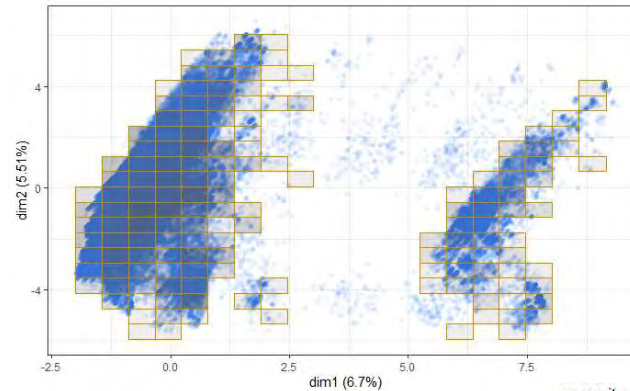
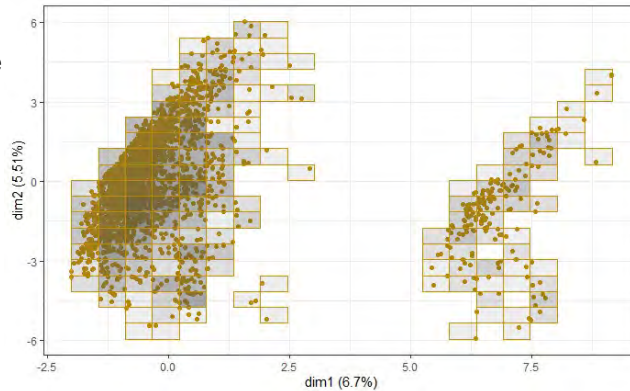
```
hiddenrate %>% summary()
medianLocalCloaking %>% summary()
closestDistanceRatio %>% summary()
hellingerMean %>% summary()
`
```

	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
hiddenrate	81.63	82.69	83.31	83.24	83.75	85.09
medianLocalCloaking	3	3	3	3	3	3
closestDistanceRatio	0.7305	0.7436	0.7488	0.7492	0.7550	0.7679
hellingerMean	0.03976	0.04069	0.04125	0.04132	0.04183	0.04354

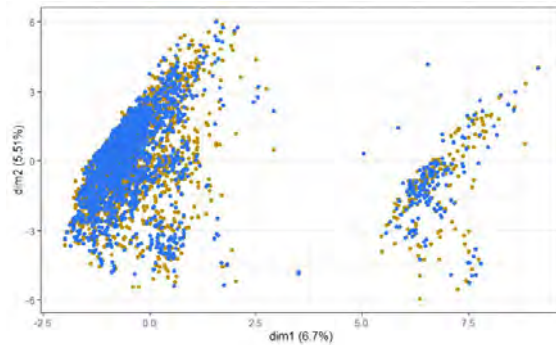
II) Avatarization of OS + RCT data for filter-based similar patient data visualization (independent data points)

Mask in reduced multidimensionnal space

PC1 binned 20 times,
PC2 binned 20 times,
PC3 binned 20 times

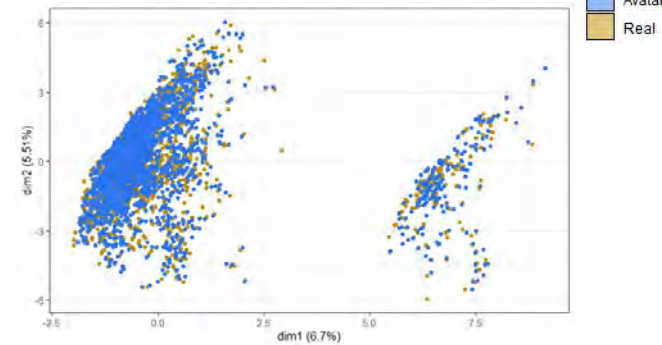


The exploratory 1:1 avatarization



→ 2733 avatars
resampled against 2744
real PORs

Resampled after 100:1 avatarization



Multiple avatarizations and resampling

II) Avatarization of OS + RCT data for filter-based similar patient data visualization (independent data points)

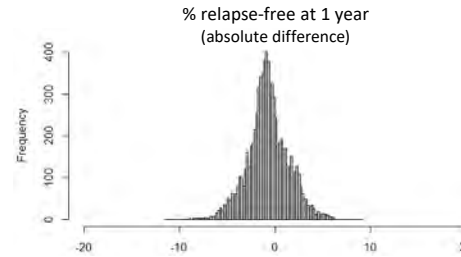
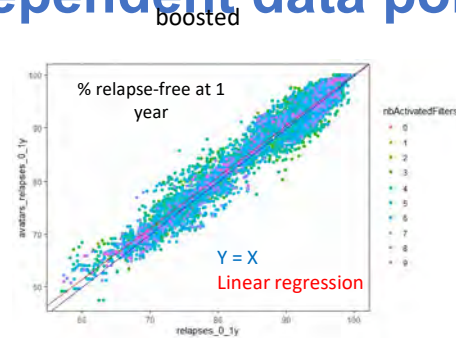
Multiple avatarizations

N queries = 4672

Potential uses of CDSS

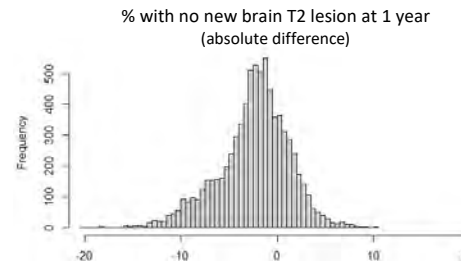
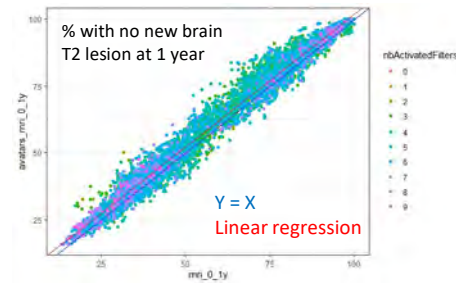
Strata of reference patients

Based on all subgroups with ≥ 100 real patients:

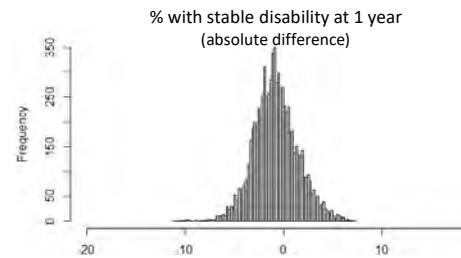
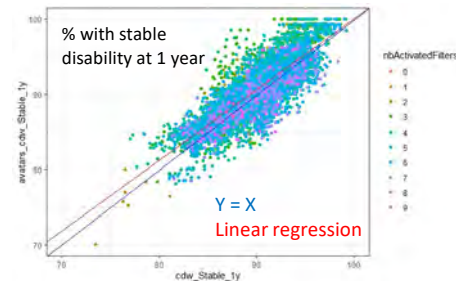


100:1 avatarization (+ mask-based resampling)

Error	Percentage of queries
Median error	1.4 %
Mean error	1.8 %
+/- 5% (absolute)	96.5 %
+/- 10% (absolute)	100 %

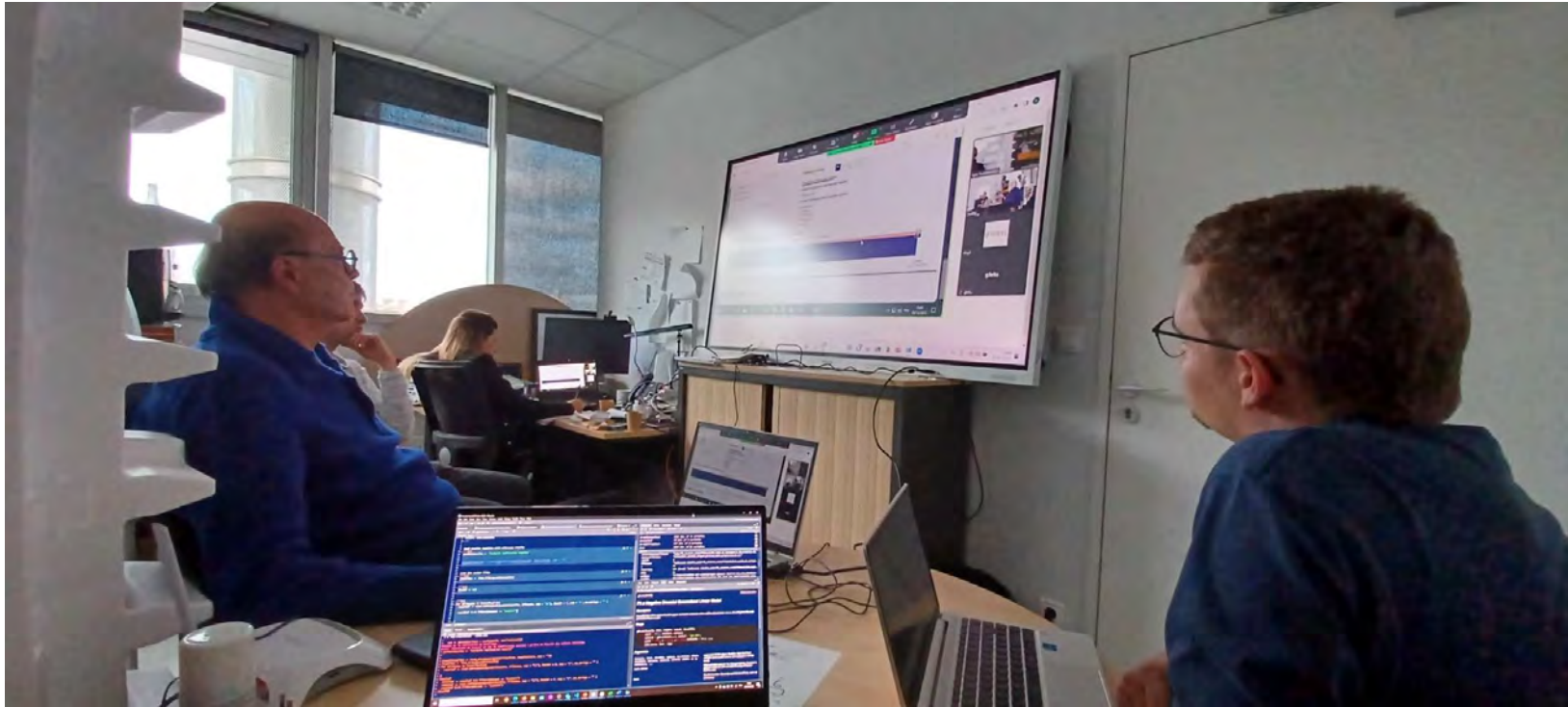


Error	Percentage of queries
Median error	2.5 %
Mean error	3.3 %
+/- 5% (absolute)	78.4 %
+/- 10% (absolute)	96.6 %

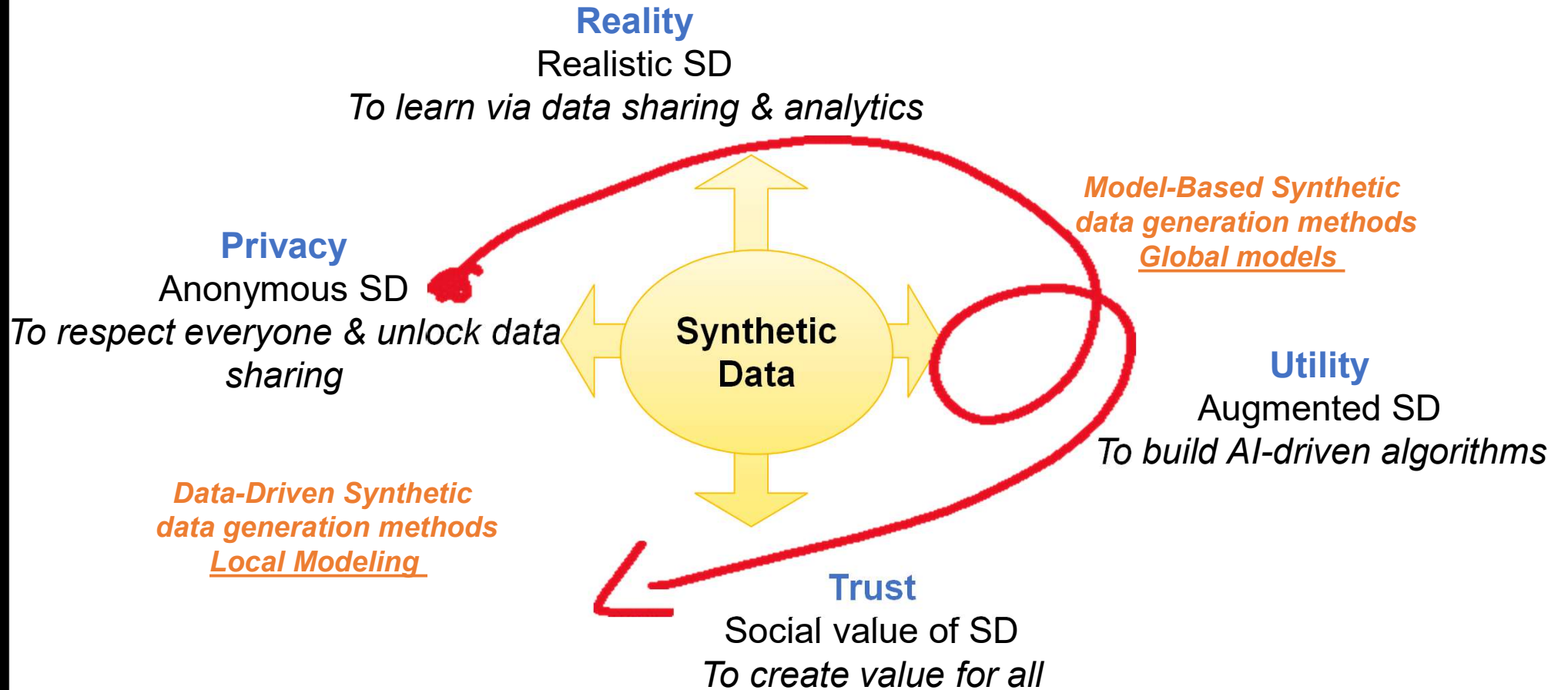


Error	Percentage of queries
Median error	1.6 %
Mean error	1.9 %
+/- 5% (absolute)	96.7 %
+/- 10% (absolute)	100 %

Tests utilisateurs et essai randomisé en cluster



Conclusion : une introspective graphique



Ouverture : Un glissement épistémologique ?

- **De la modélisation d'un phénomène ... par des modèles mathématiques**
 - Ex Maladie : un malade n'est qu'une réalisation de l'entité maladie
 - Ex Jumeaux digitaux : Machine, un organe
 - L'objet de la modélisation est méta-individuel
- **De la modélisation d'un patient atteint d'une maladie à la modélisation d'un individu**
 - Recours à une modèle faible local privé à usage unique
- **Deux orientations ? Confidentialité et Utilité**
- **Deux orientations des données synthétiques**
 - F Gross "Towards a Methodology for Systems Biology" 2017)
 - 1 → Model-based System Biology**
 - Miser sur le modèles / Les données synthétique (ex GAN)
 - 2 → Data-driven System Biology**
 - Miser sur les donnée « data-driven »
 - Agnostique de l'analyse qui en est fait, et paramétrable (donc explicable)
 - Modélisation privée et à usage unique.

Acknowledgement:



Ce travail bénéficie d'une aide de l'Etat gérée par l'Agence Nationale de la Recherche au titre du 3e PIA, intégré à France 2030 portant la référence ANR-21-RHUS-0014

Nous remercions nos partenaires, notamment l'OFSEP, Merck Santé et Biogen qui fournissent les données utilisées.



Prof. Gilles Edan
PRIMUS consortium

Prof. G. Edan,
Prof. P-A. Gourraud,
Prof. D. Laplaud,
Prof. L. Michel,
Prof. S. Vukusic
Dr. A. Kerbrat
Dr. E. Le Page
Dr. C. Dumas
Dr. R. Casey
N. Blanc
J. Martin-Gauthier
G. Jarre
S. Doyle
A. Auffret
F. Dubois
M. Payet



Prof. Pierre-Antoine Gourraud,
INSERM U1064, Nantes,
"Translational Immunogenomics of
Transplantation and Autoimmunity"

Prof. D Laplaud,
L. Berthelot,
A. Nicot,
F. Cornelis,
E. Dugast,
A. Garcia,
A. Serova-Erard

Prof. P-A. Gourraud,
Dr. S. Limou,
Dr. N. Vince,
Dr. M. Morin,
J. Martin-Gauthier,
C. Ed-Driouch,
O. Rousseau,
J. Paris,
I. Faddenkov,
A. Durand,
V. Mauduit,
V. Douillard,
I. Charles,
L. Boussamet,
S. Sayadi,
Nayane,
S. Demuth



L'équipe de la Clinique des Données - janvier 2024



Pierre-Antoine GOURRAUD
Responsable
Pôle de Santé Publique
Sep. 2015



Matthieu WARGNY
Responsable adjoint
Médecin de Santé publique
Nov. 2018



Thomas GORONFLOT
Epidémiologiste
Déc. 2017



Sandrine COUDOL
Epidémiologiste
Fév. 2018



Matilde KARAKACHOFF
Epidémiologiste
Déc. 2018



Soline BOBET
Epidémiologiste
Fév. 2023



Chloé DOUAREC
Ingénieur data
Nov. 2023



Adrien BAZOGE
Ingénieur TAL
Nov. 2023



Emilie VAREY
Chef de projets / Dpt DPI
Nov. 2021



Delphine TOUBLANT
Chef de projet SI
Déc. 2019



Laëtitia AUBERT
Master MEDS
Oct. 2022



Pacôme CONSTANT DIT BEAUFILS
Doctorant / Neurologue
Nov. 2021

Etudiants M2 RC ou BioInfo

Internes de Spécialité

Acknowledgments



Team 5 « NEMO »
NEuroinflammation, Mechanisms and therapeutic Options



Team 3 “iTHINK” integrative Transplantation, HLA, Immunology and geNomics of Kidney injury



Génération de données synthétiques centrées sur le patient, aucune raison de risquer la ré identification dans l'analyse des données biomédicales

Amphi Schwarz, IMT 11:30 - 12:00

Pr Pierre-Antoine GOURRAUD, PhD MPH

Vendredi 9 Février 24, Toulouse

Professeur des Universités & Praticien-Hospitalier
 CHU Nantes, PHU 11 : Clinique des données, INSERM,
 CIC 1413, Nantes Université, INSERM, CR2TI

pierre-antoine.gourraud@univ-nantes.fr

