

“La Clinique des données”

Pr Pierre-Antoine GOURRAUD

« Un check-up de l’IA pour la santé numérique, c’est grave docteur(s) ? »

Séminaire Aristote - jeudi 27 février 2020,

École Polytechnique, Amphithéâtre Arago
91120 Palaiseau

ATIP-Avenir Team 5 «Translational Immunogenomics of Transplantation and Autoimmunity »,
ITUN - CRTI - UMR Inserm 1064 -CHU de Nantes

Pôle Hospitalo-Universitaire 11 : Santé Publique, Santé au Travail et Pharmacie,

Hôpital St-Jacques - CHU de Nantes



« Un check-up de l’IA pour la santé numérique, c’est grave docteur(s) ? »



→ **Passez donc à la « Clinique des données »**

1. Quelques notions clés
2. Un nouveau service – et ses motivations: “la clinique des données”
3. La technique au service de la gouvernance des données : “les avatars”

Quelques notions en filigrane

Partie 1

A two-fold re-evolution...

Ever lower cost of data access

- Information become technically available
- Central role of search (how to access what)

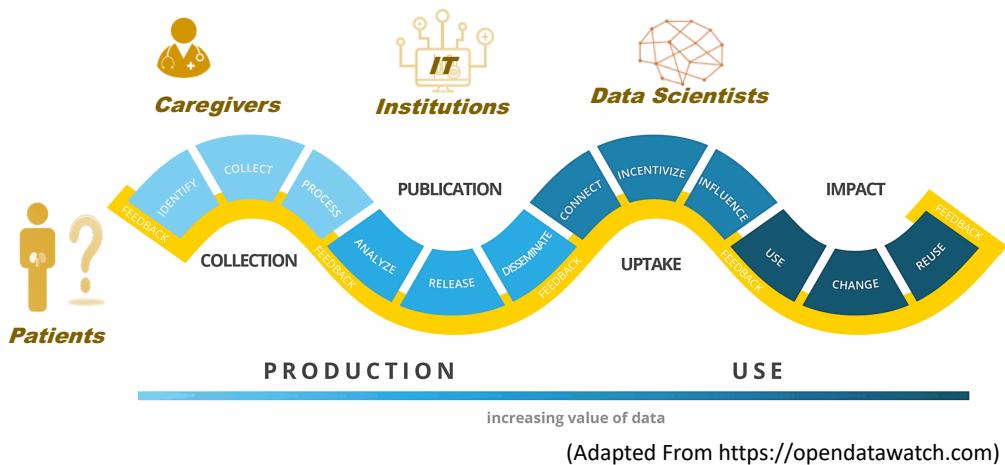


Ever lower cost of computation

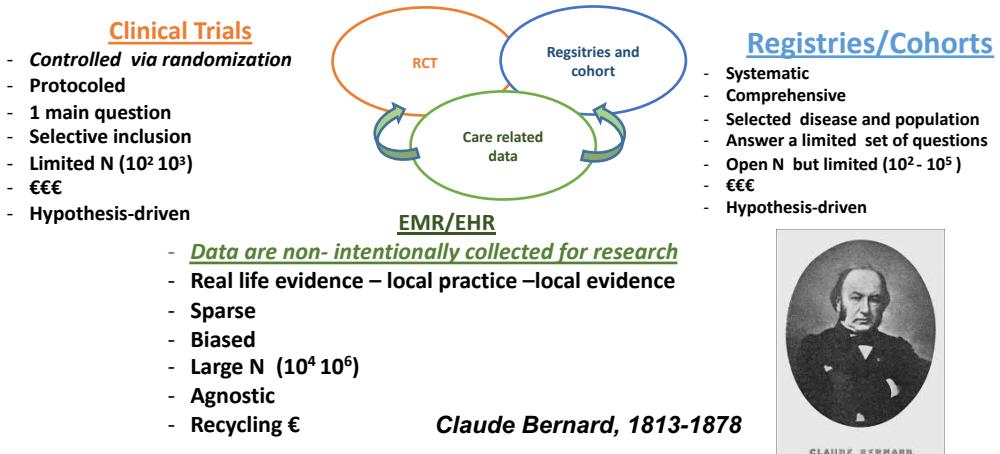
- Computation
- On-demand
- Just computations...



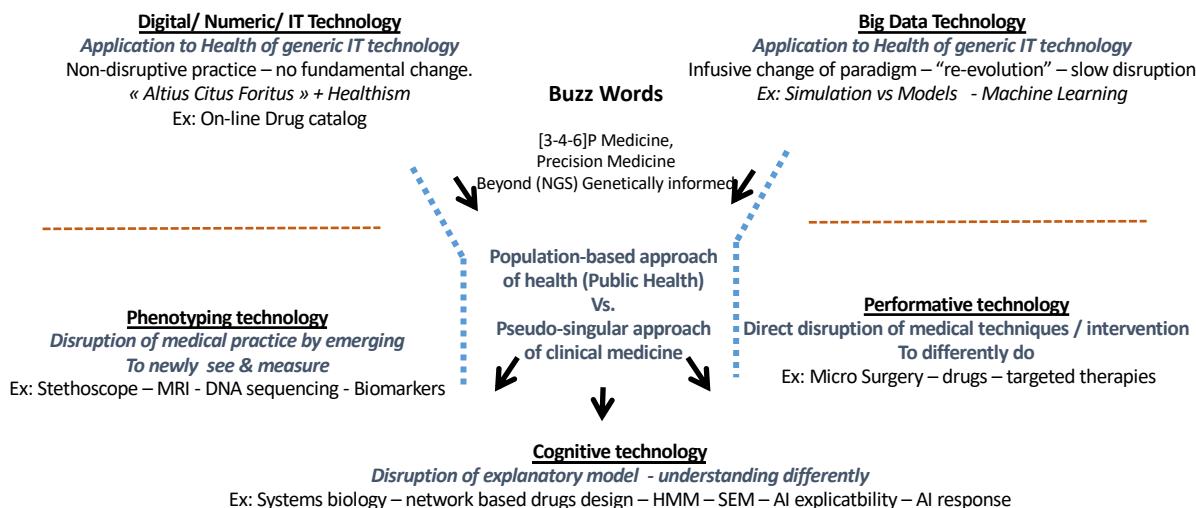
Multiples stakeholders and contributions to the data chain value.



A shift in health Data for research



All innovation in health now deals with IT A framework

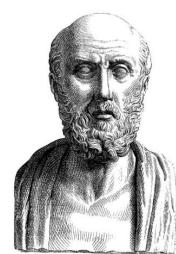


Data exceptionalism in health « not matter how big... ”

- Technically complex, Legally ruled, Socially sensitive
 - Health data at the crossroad of patients, health care systems, and care givers (vs.genomics)
 - V : Variety ; N : Numbers ; C : Categories ; P : Probabilistic
 - T : Time is needed, duration matter
 - Health data are sensitive – **Privacy**

“Primum Non Nocere”

- **Data come from patients...**
- **Patients come first**
- **It applies to ... Data protection in health research**
 - **First do no harm**
 - **First respect patient privacy**



Hippocrates of Kos;
 c. 460 – c. 370 BC



La « cliniques des données » et l'entrepôt de données biomédicales du CHU de Nantes

Partie 2



Un réflexion sur l'enjeu de données

Constat partagé : Verrou des données

« Bonnes » données, un préalable à de bonnes analyses

Les données: c'est l'expertise du domaine

Compréhension et qualité des données sont un enjeu croissant à mesure que la disponibilité technique est résolue par les outils numériques.

Les données méritent une « clinique » au CHU de Nantes

Réfléchir, Traiter, Intégrer = Consultation, Open-Space, Formation

- Dans son rôle d'accompagnateur méthodologique des projets de recherche et d'évaluation , la « Clinique des Données » devient médiateur de l'accès à l'entrepôt de données du CHU de Nantes (source de données)
 - Fait la requête et fournit une extraction utile = décompte ou datamart pseudo-anonymisées(principe de parcimonie)
 - Garant de la protection du patient-usager-payeur-citoyen (Transparence et contrôle des usages et finalités)


CENTRE HOSPITALIER
UNIVERSITAIRE DE NANTES

Trois dimensions à articuler



1. Gouvernance garantissant le respect éthique, réglementaire et scientifique de chaque projet de recherche (=structure de confiance vis-à-vis des patients et des soignants)

2. Information des patients, pour chaque exploitation des données avec opposition possible, par affichage site web CHU

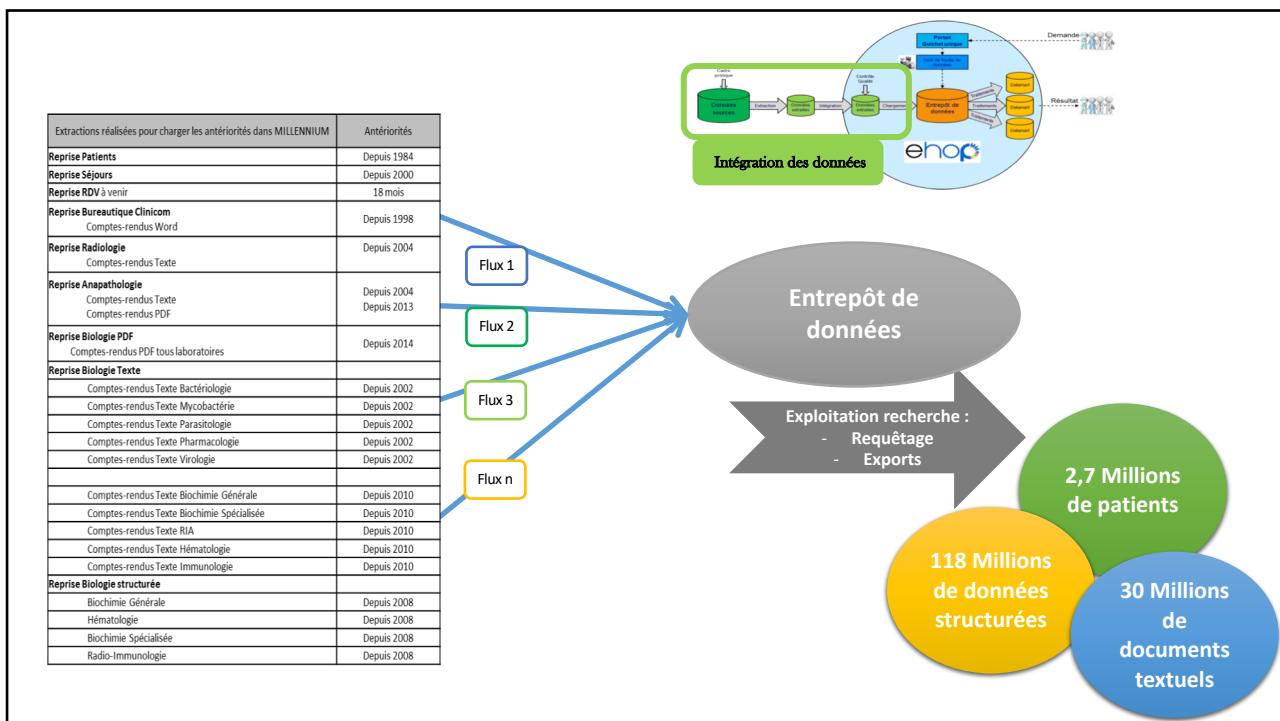
- Communication grand public (avril 2018→)

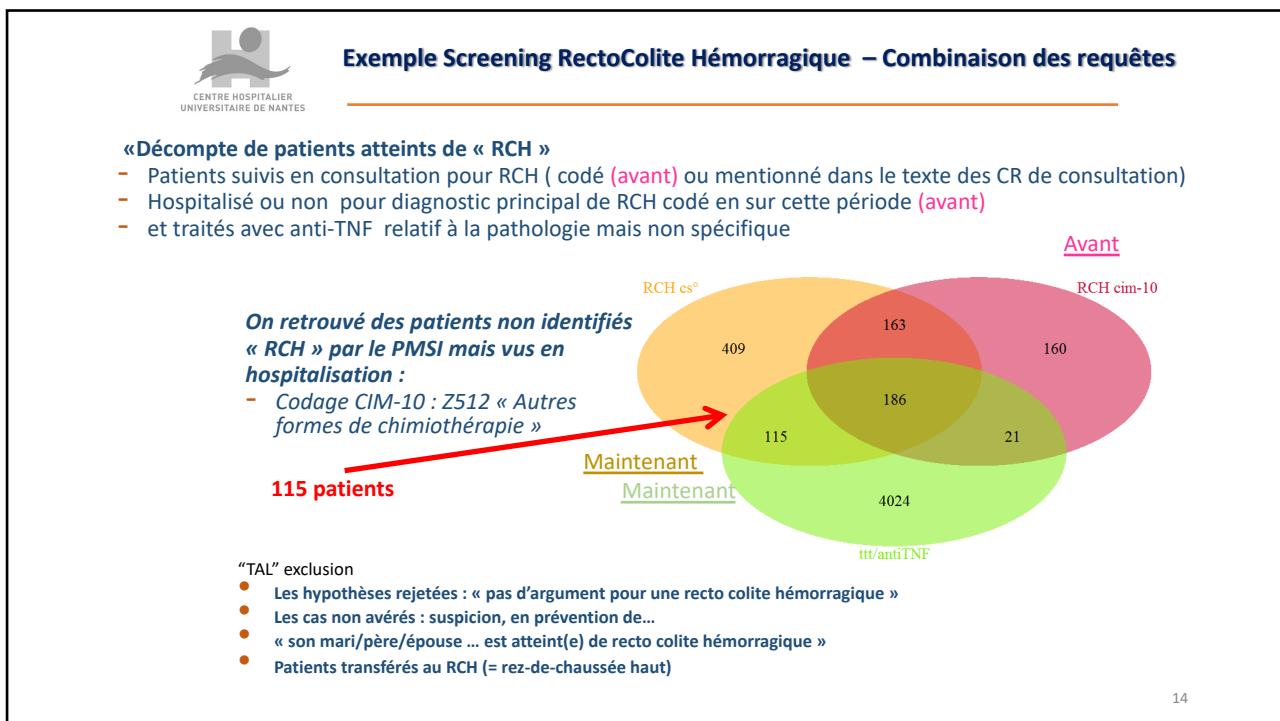
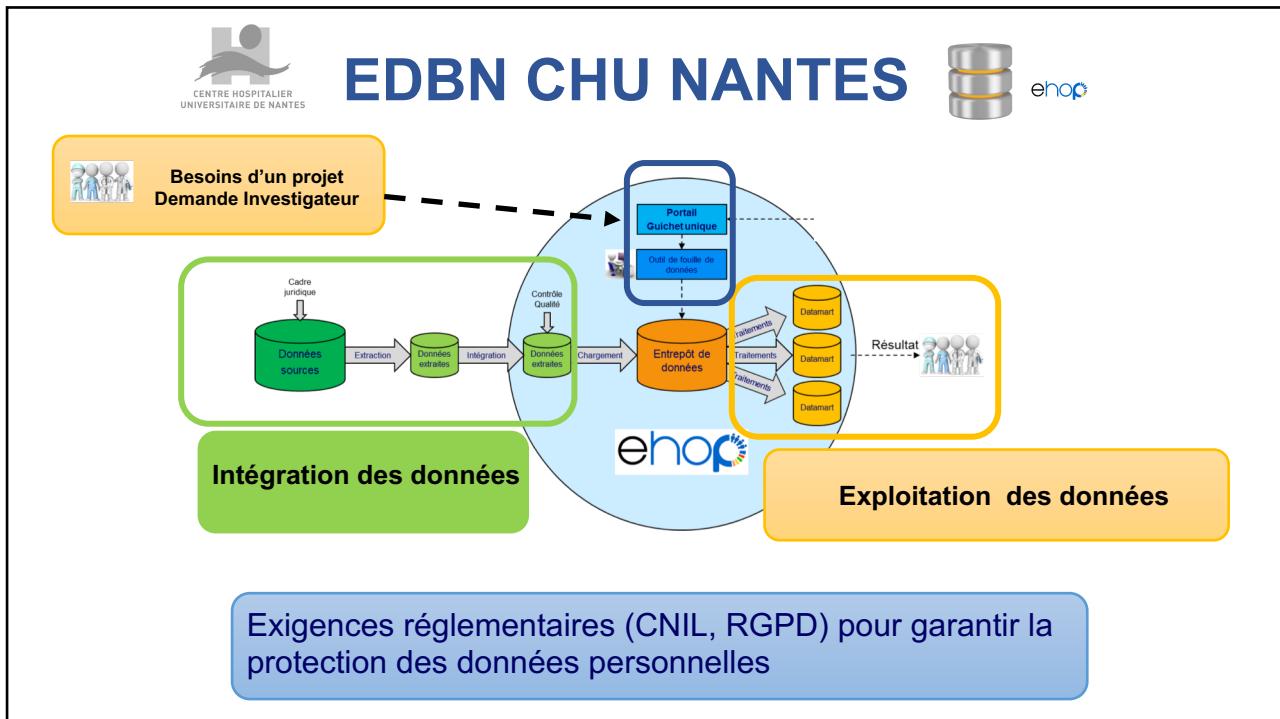
3. Outil des gestion données générées par le soin au CHU de Nantes

- Ehop (Pr M Cuggia - Rennes)= Transformation/intégration de données + Entrepôt de données type « Big Data » + Requêteur « intelligent »
- GIRCI GO en 2016-2018 (réseau des RiCDC)- HUGO


Juillet 2018



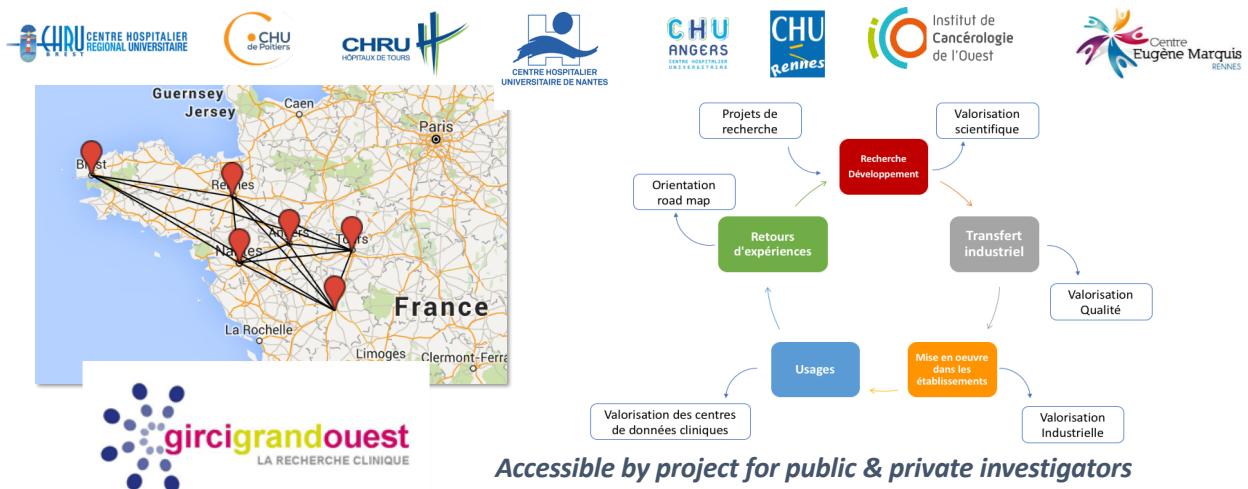




Perspective : Les CHU “carrefours des données” de Santé

- 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
- Enjeu de transformation – nouvelle épidémiologie de données.
 - On parle très vite en millions...
 - Données en « vie réelle » - qui requièrent plus de méthodes

Regional Context - Millions of Medical Data from University Hospitals



Biomedical data warehouse on the move



4,6 millions patients



112 millions clinical documents



A single IT for 8 institutions at the crossroad of healthcare system for 11 millions people



5+ millions Visits



1.3 billions structured data

The 3 ingredients hierarchy of Biomedical Data Warehouse

3- A common IT infrastructure

When data meets IT – interoperability is not guaranteed

2- A regulatory framework to inform patients

Nantes Biomedical warehouse was approved on April 2018

1- Governance

Mediated access to the wealth of data
Internal for research, screening etc

Stratégie Interrégionale Projet GAVROCHE



Une question clinique

La variabilité glycémique permet-elle d'affiner le pronostic vital des patients admis pour insuffisance cardiaque aiguë ?

Une source de données massive inter-entreprises du Grand Ouest

+ de 20 000 patients admis, + de 100 000 mesures de glycémie, > 8 ans d'activité

Des données hétérogènes mais accessibles

- Structurées : biologie (DXLAB) et diagnostics (PMSI)
- Non structurées : clinique extraite du DPI par traitement automatique du langage



=> Un parmi 4 projets fédérateurs du GIRCI Grand Ouest

Contexte national

Projet



- Accès facilité aux données de santé de sources multiples (hôpitaux, ville, pharmacie, biologie...)
- Environnements sécurisés
- Valorisation du « patrimoine » français de données de santé
- Hub national et hubs locaux

→ Entrepôt(s) HUGO

La technique au service de la gouvernance des données : “les avatars”

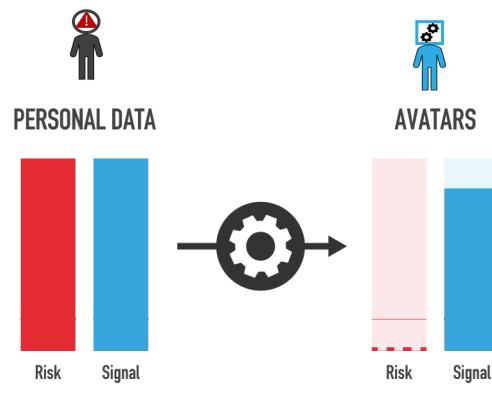
Partie 3

Risk of re-identification in (pseudo-)anonymous data

“Primum non nocere” = lower risk
 Systematic introduction of the use of avatars for each project of the “Data Clinic”



Public- private PARTNERSHIP

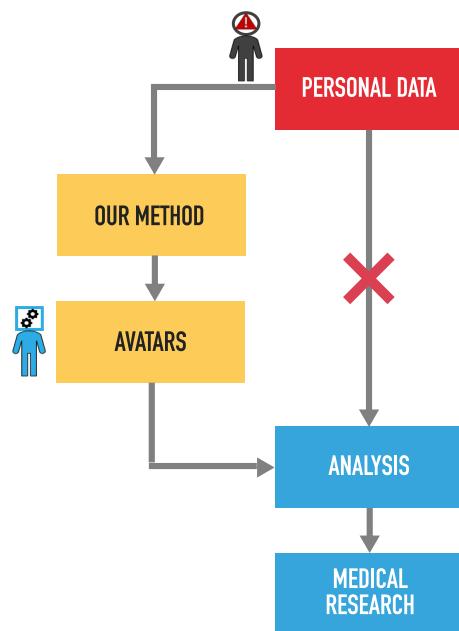


WeData method transforms personal data into signal-retaining avatars with privacy priority

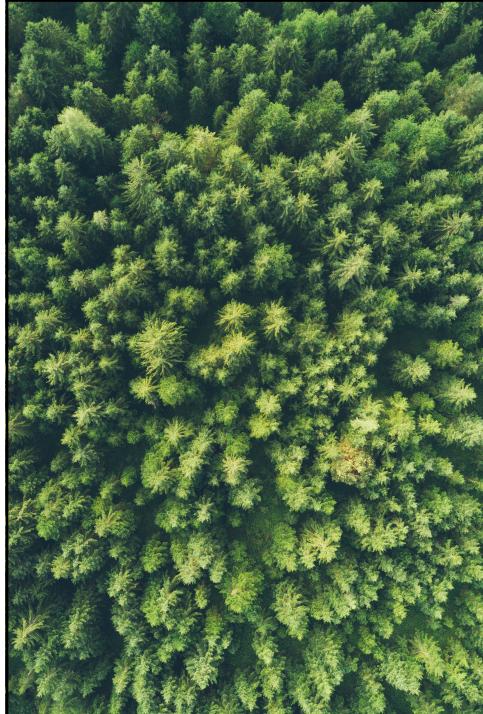


21

Enable analyses to draw comparable conclusions without compromising privacy



22



ALGORITHM

"Blend into the crowd"

- Based on "look-alike" approaches
- Each individual's **surroundings** are used to create an avatar
- **Locality** sensitive method
- An Avatar is a **synthetic** individual

23

EXAMPLE

- Application of the method

Example of avatars dataset

Individual	age	height	color
...
165	12 9	147 129	Blue Blue
166	13 11	158 140	Green Red
167	9 12	132 146	Blue Red
168	14 8	139 124	Red Blue
169	6 12	140 150	Blue Blue
...

• Original sensitive data

• Avatar

EXAMPLE

- Application of the method

Example of one individual avatarized multiple times



Individual	age	height	color
...
165	12 9	147 129	Blue Blue
166	13 11	158 140	Green Red
167	9 12	132 146	Blue Red
168	11 8	139 124	Red Blue
169	6 12	140 150	Blue Blue
...

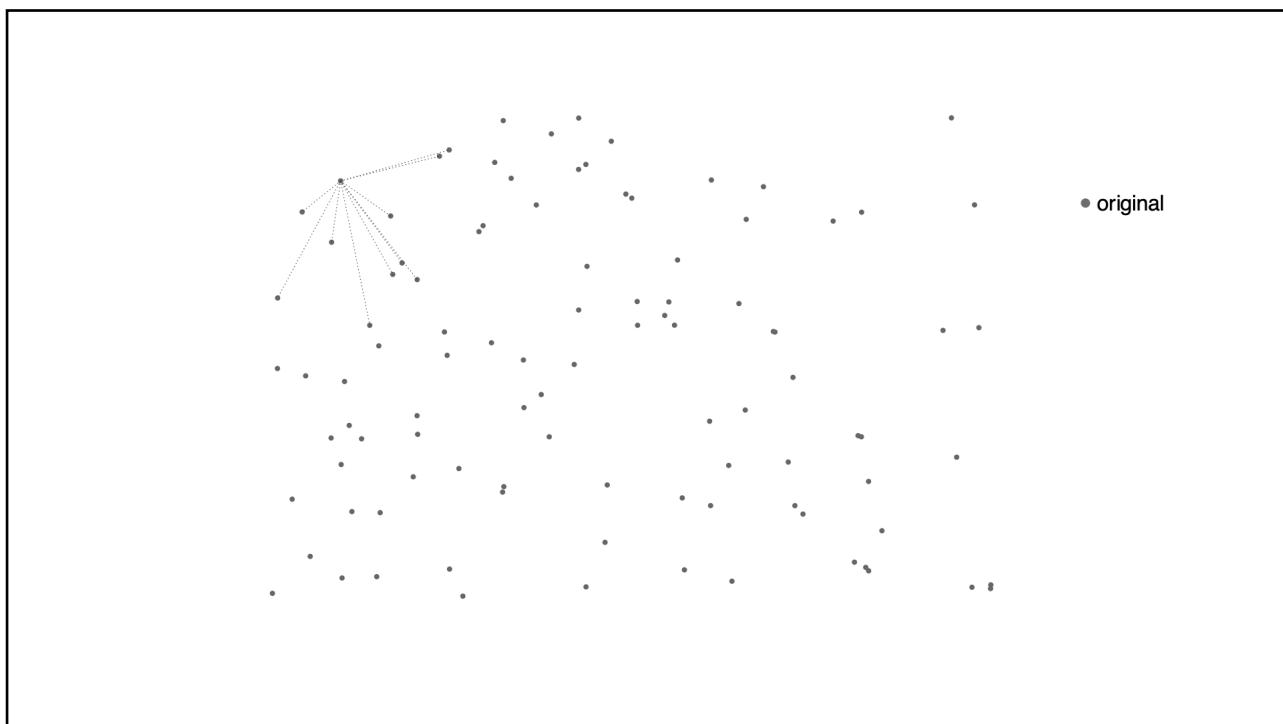
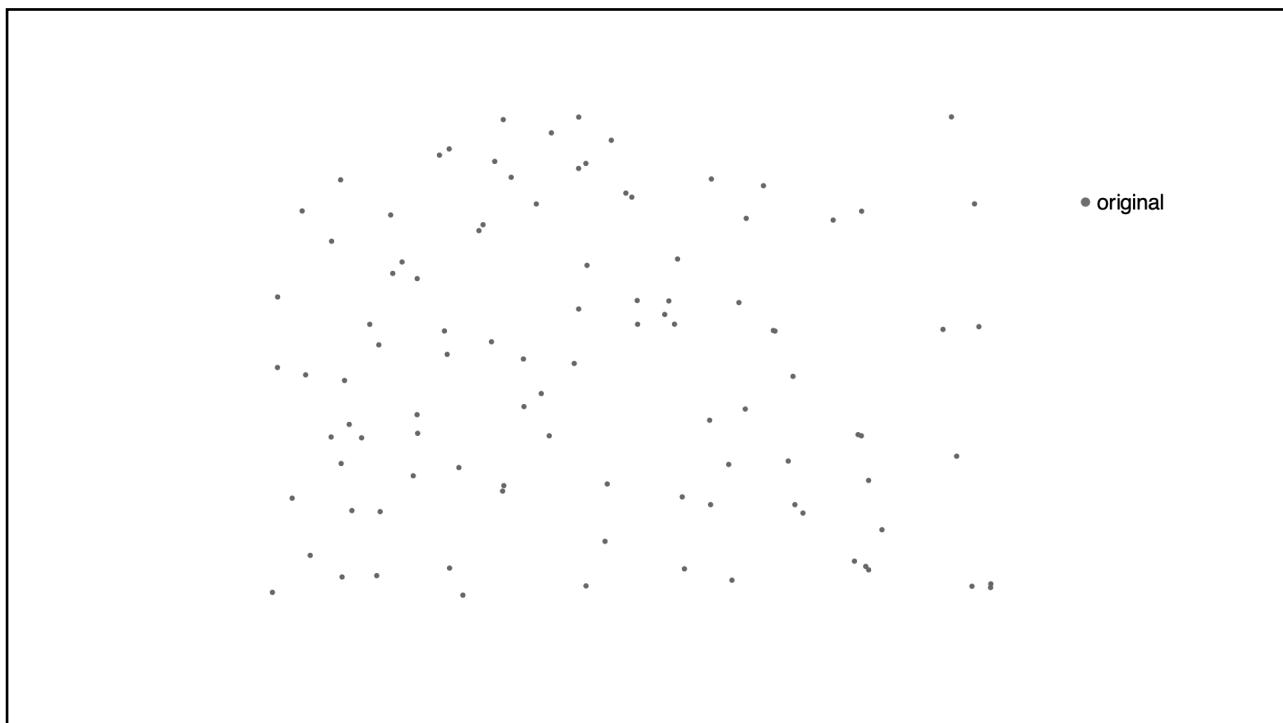
26

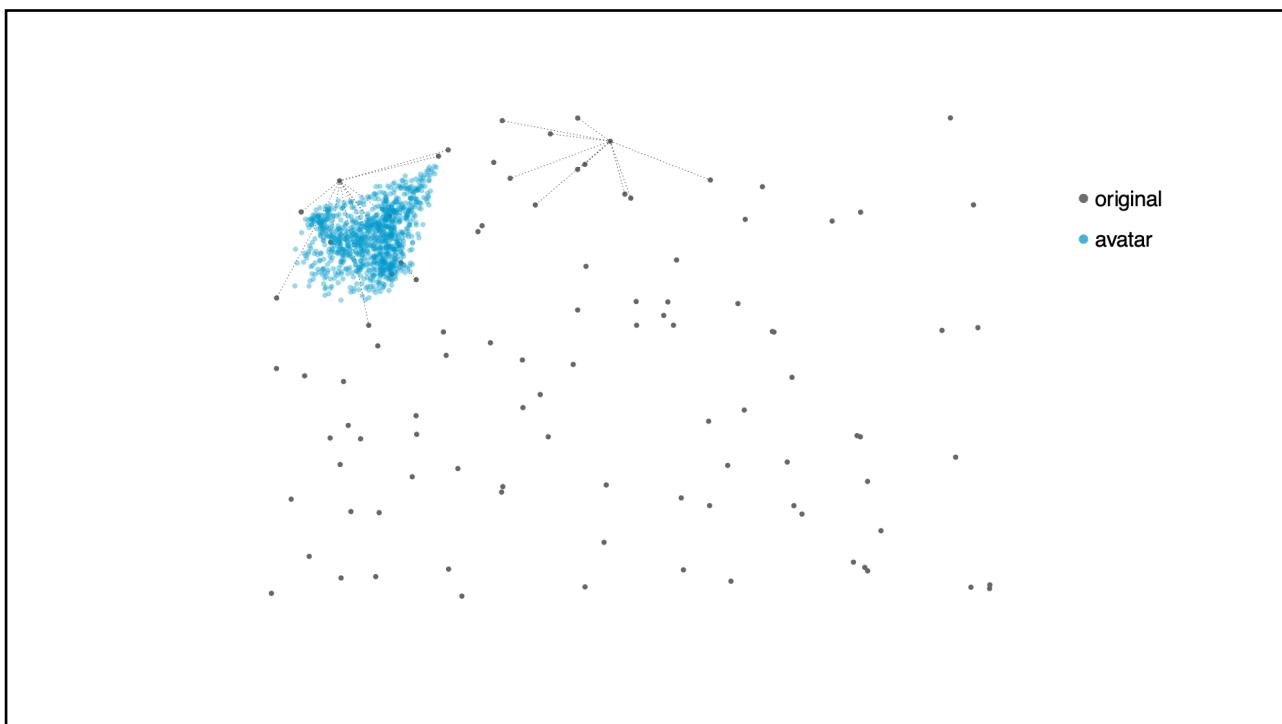
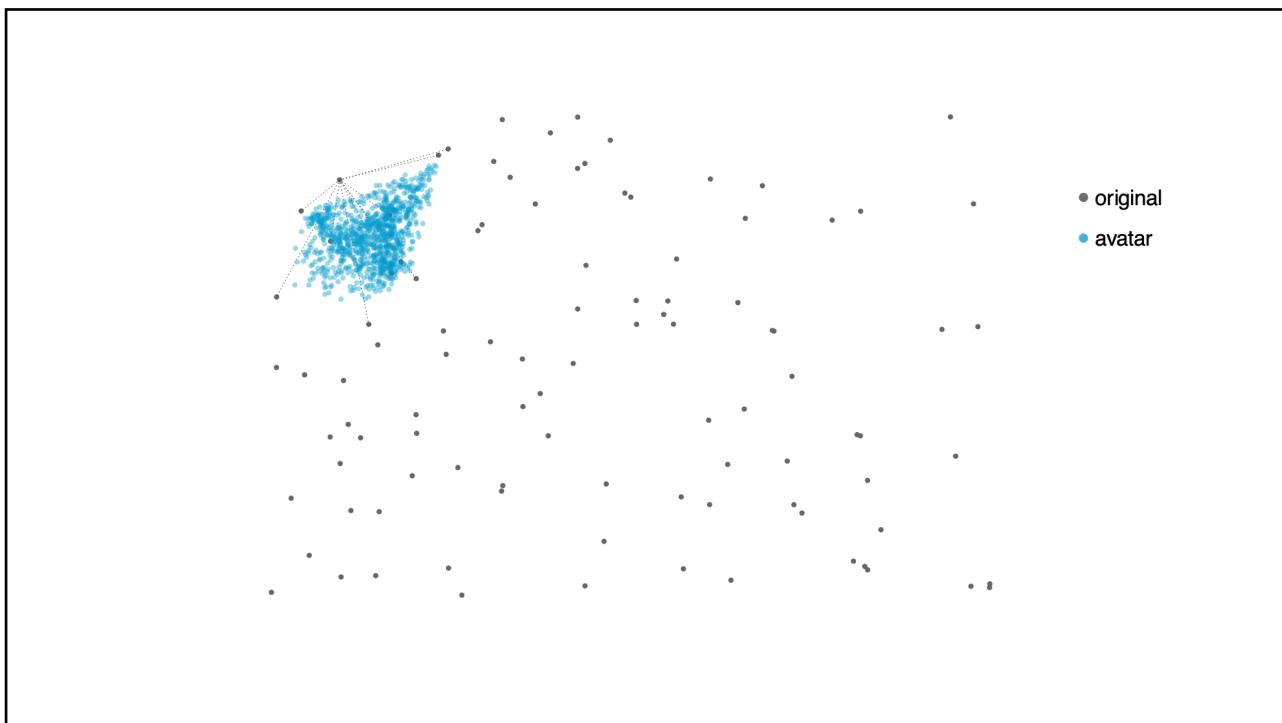
EXAMPLE

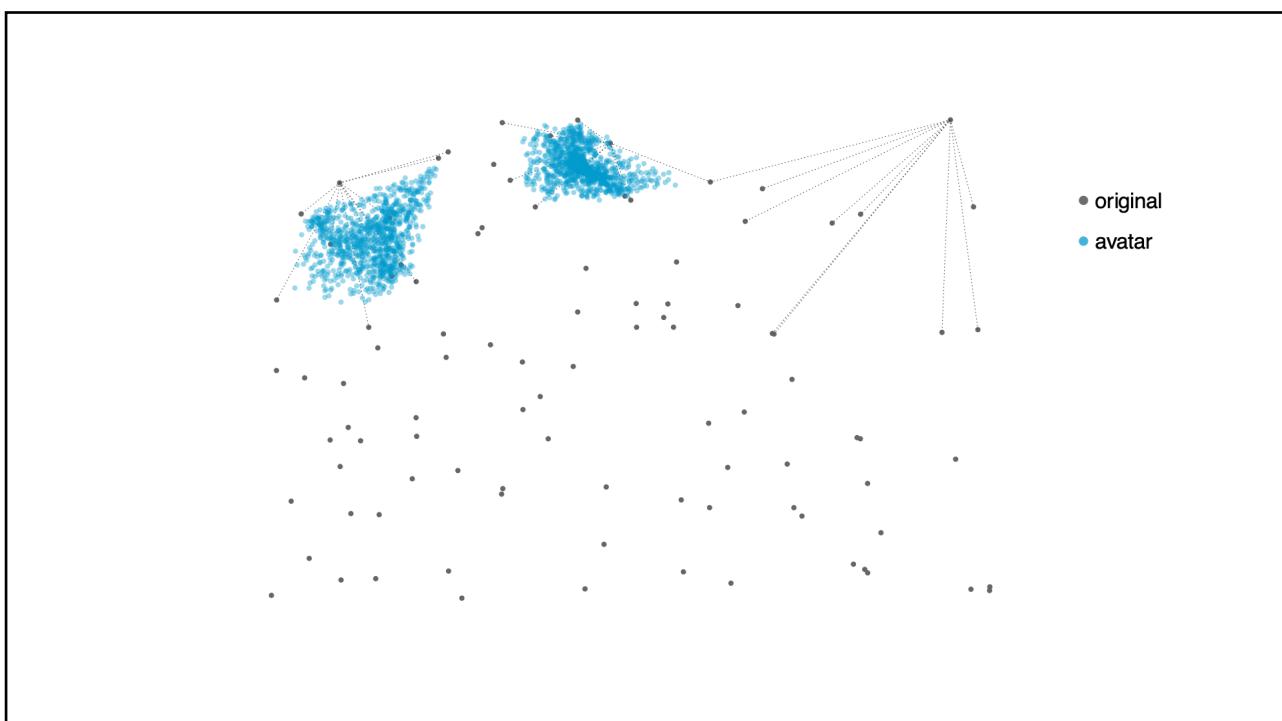
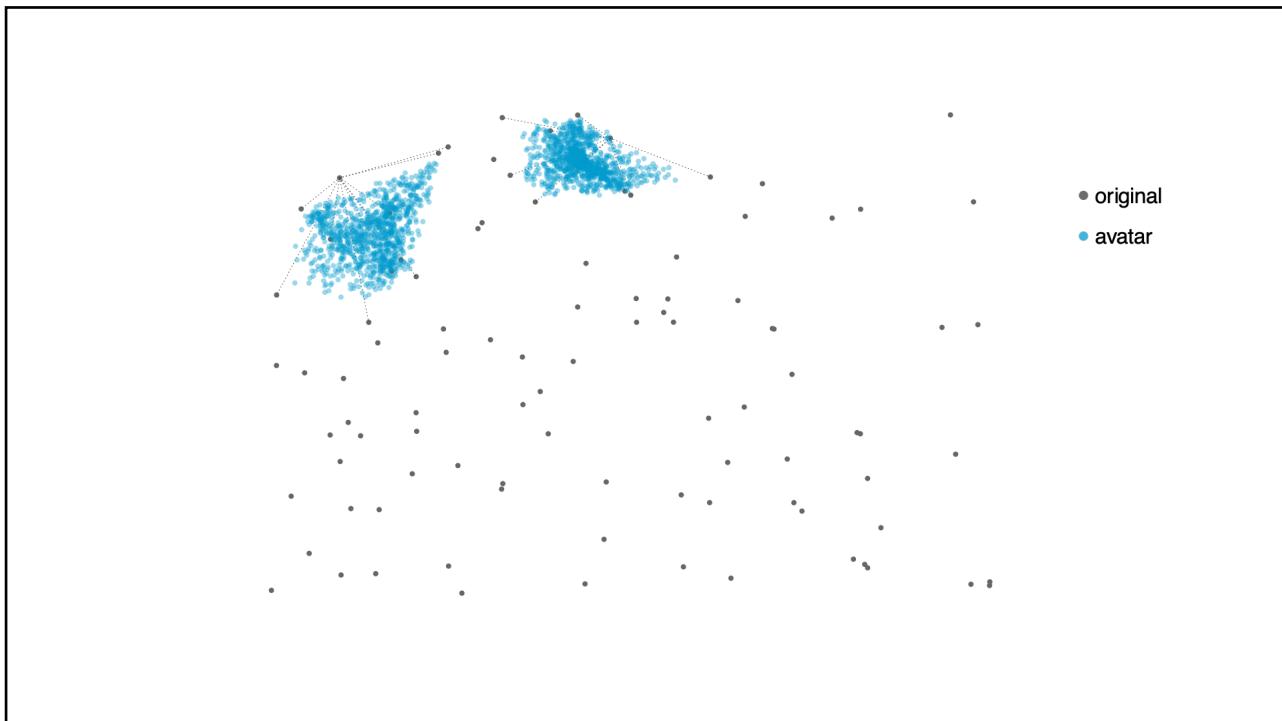
- Application of the method

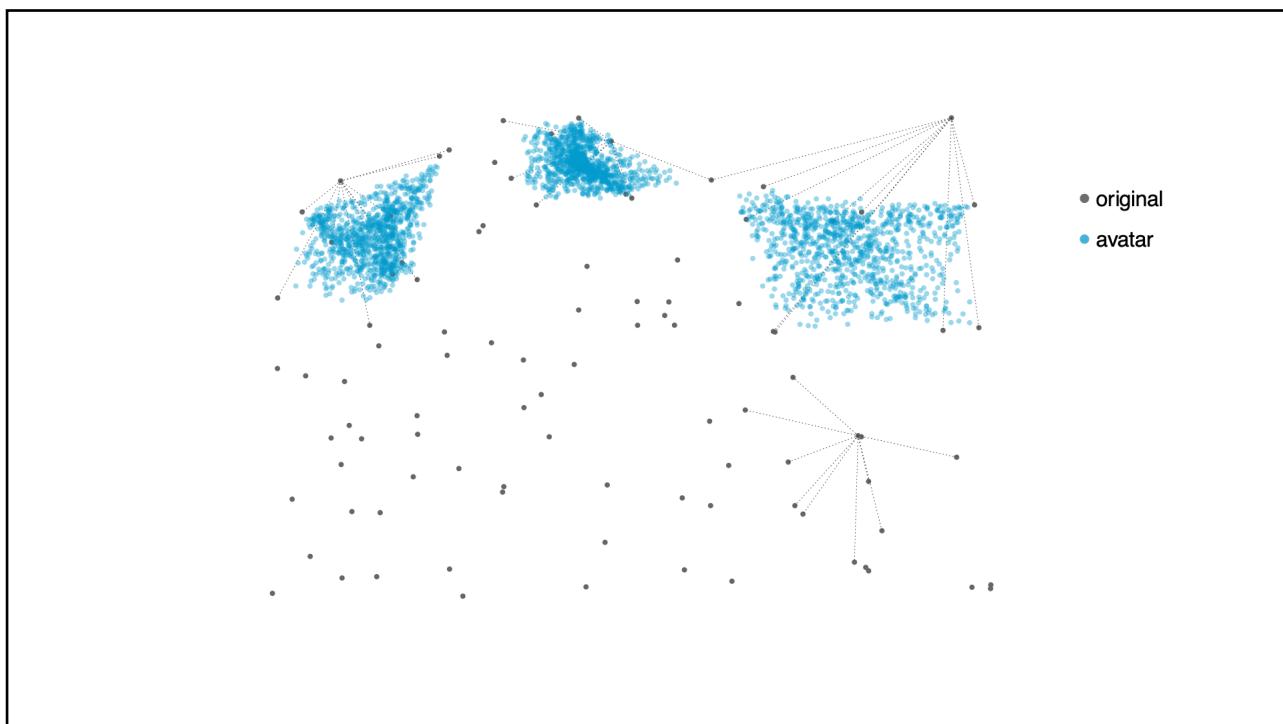
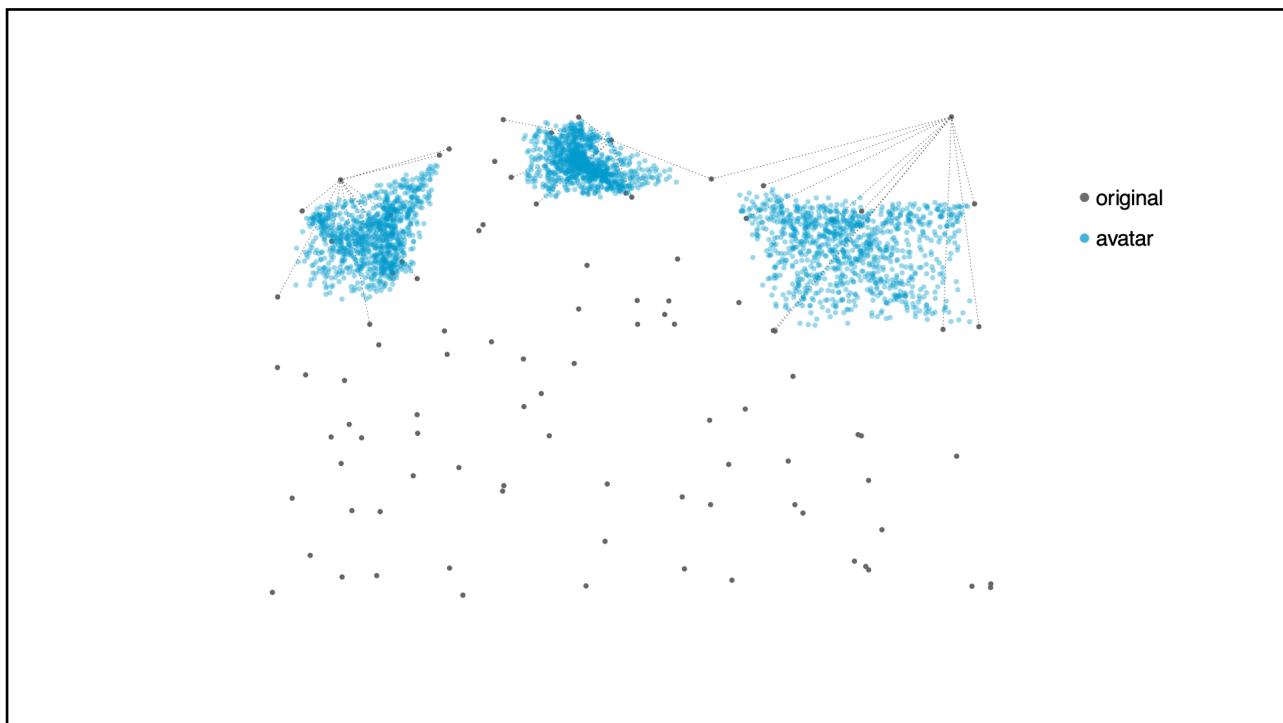
Example of one individual avatarized multiple times

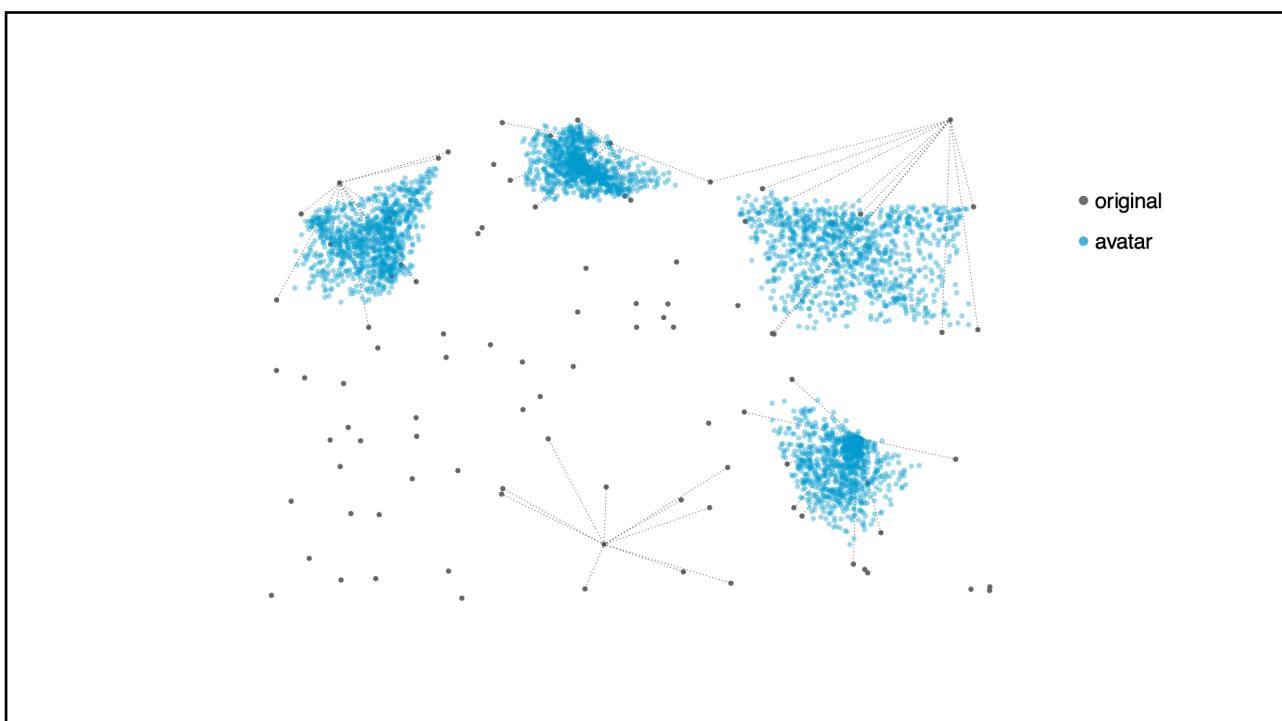
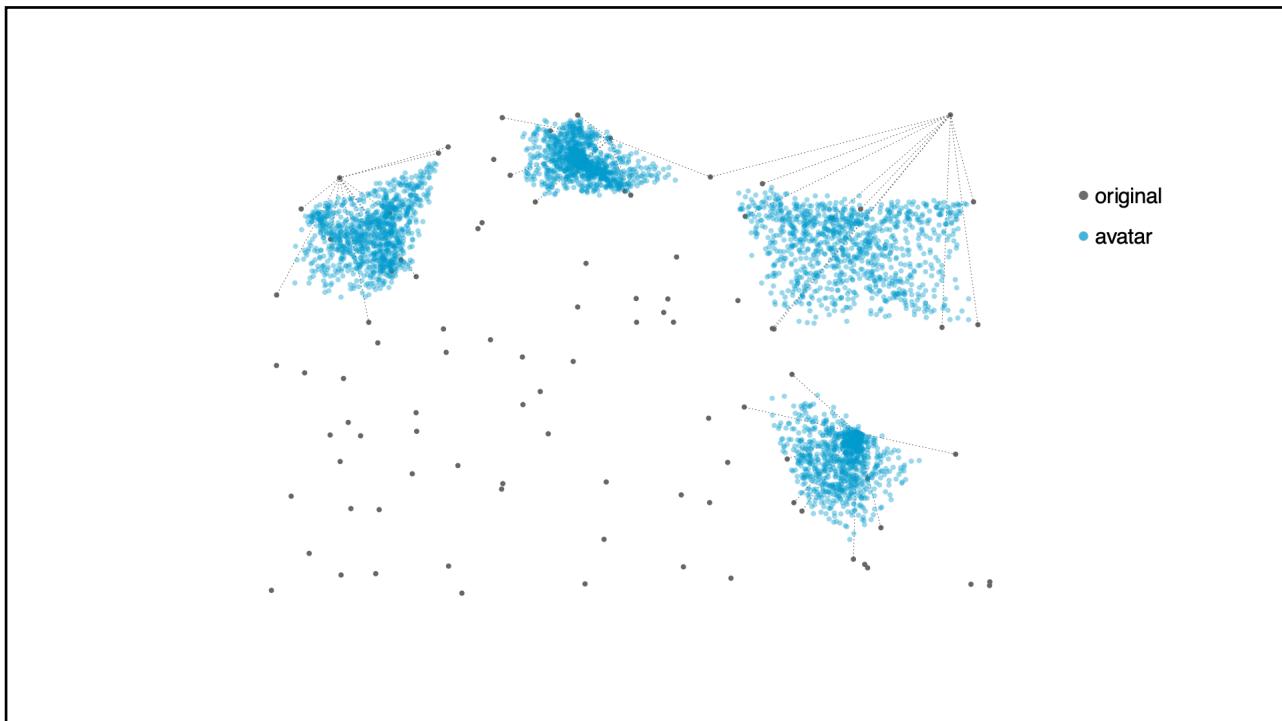
Individual	age	height	color
...
Original 167	9	132	Blue
Avatar 167 (1)	9 11	132 140	Blue Blue
Avatar 167 (2)	9 12	132 146	Blue Red
Avatar 167 (3)	9 8	132 124	Blue Blue
Avatar 167 (4)	9 10	132 129	Blue Blue
...

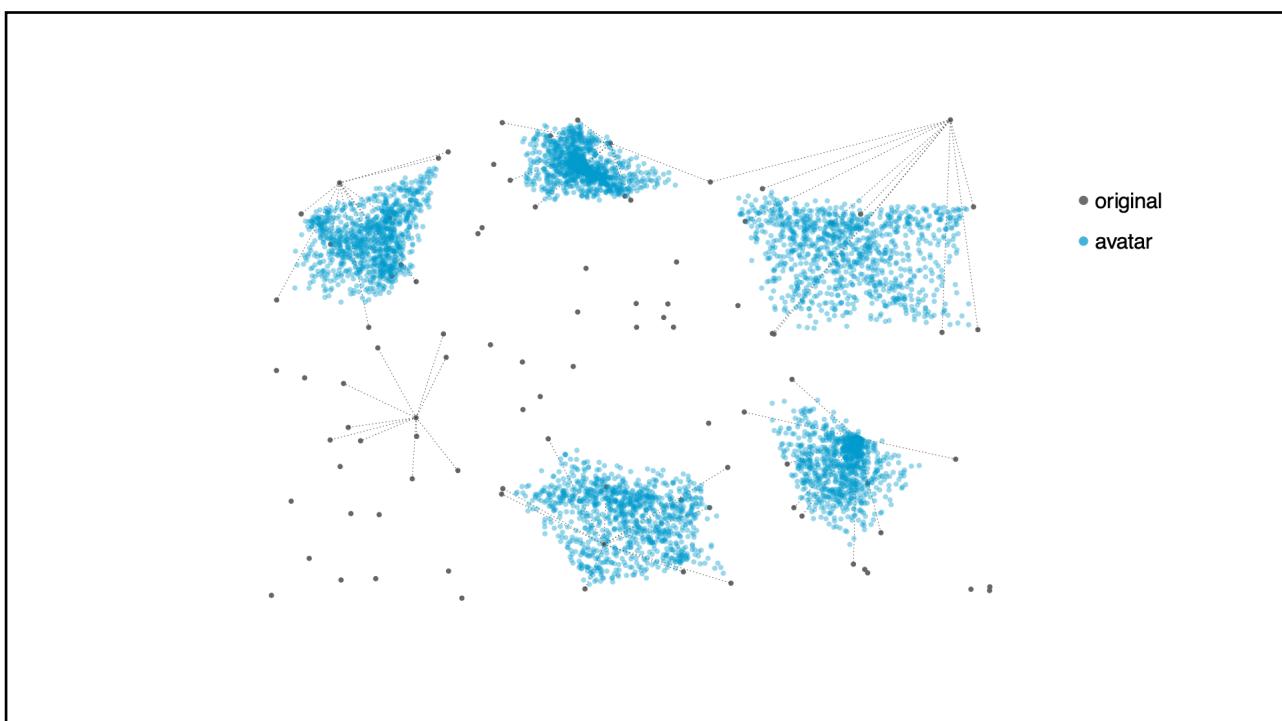
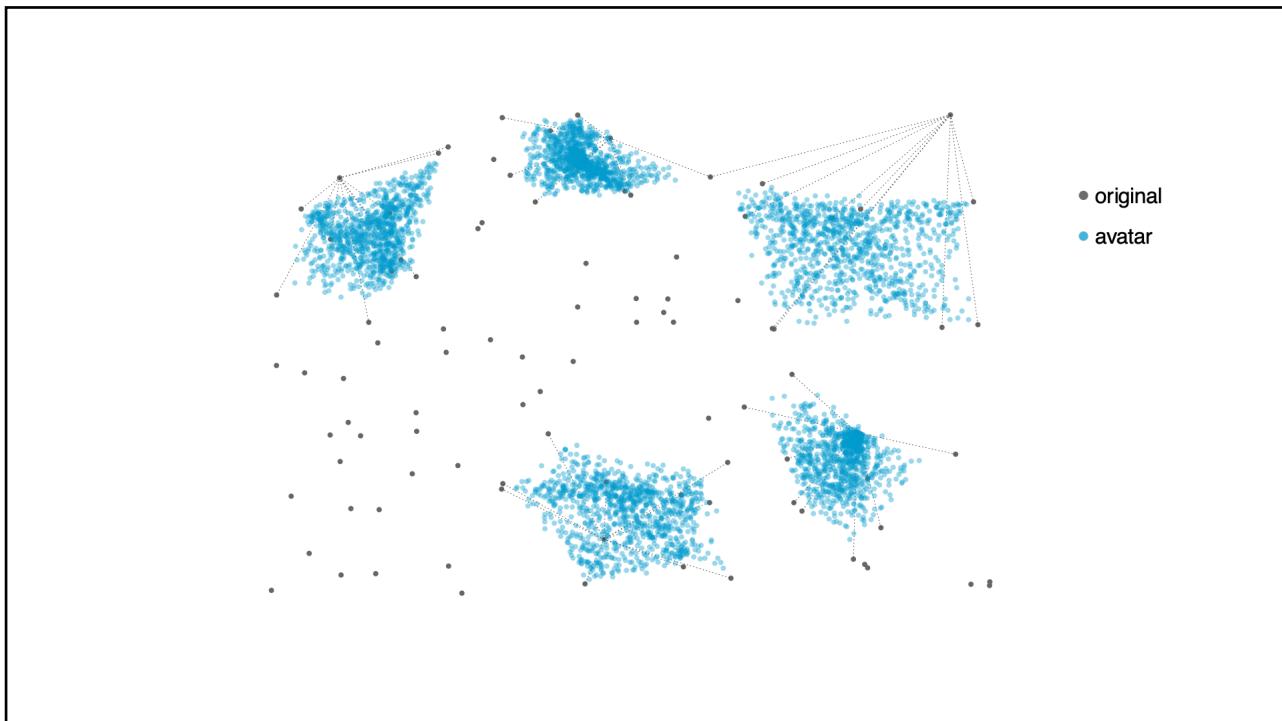


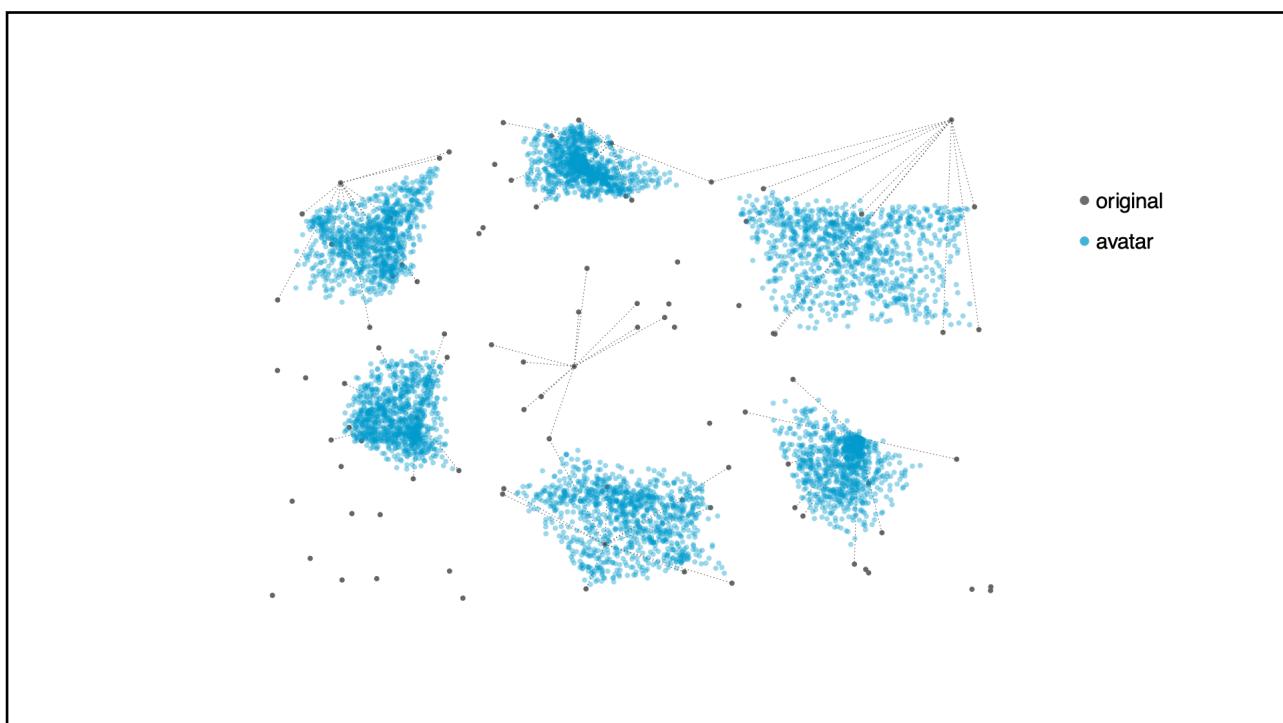
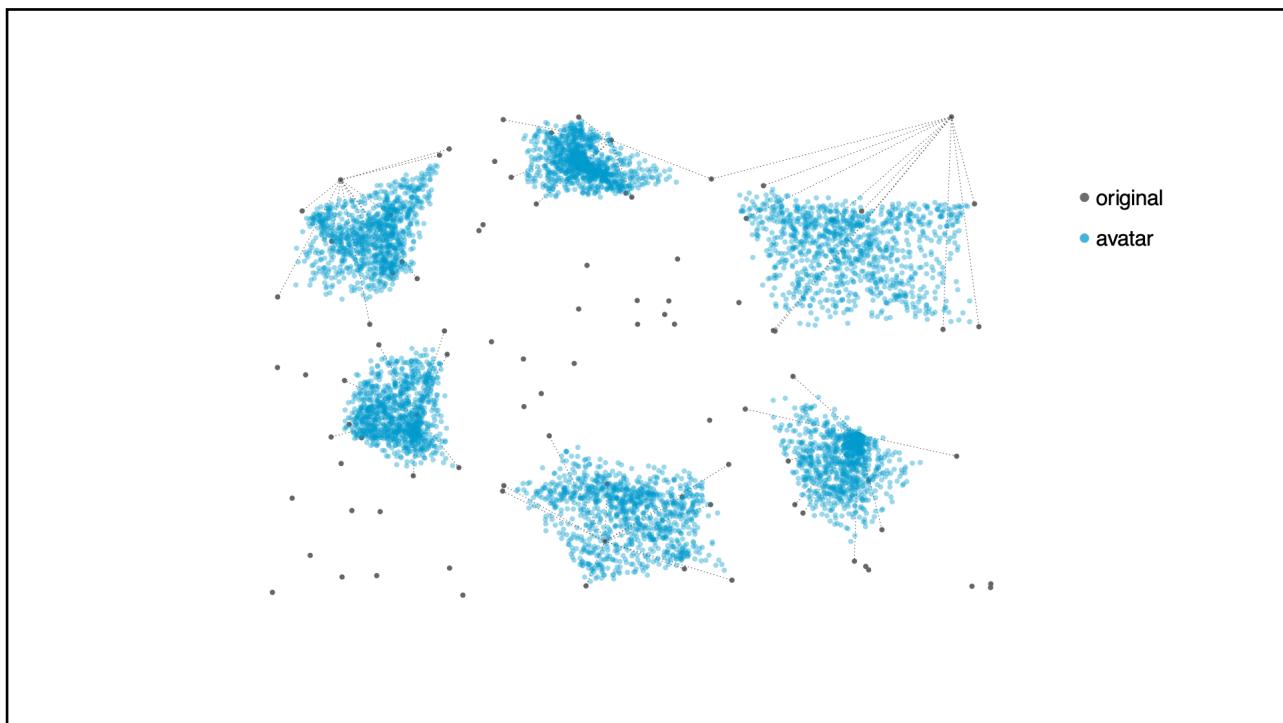


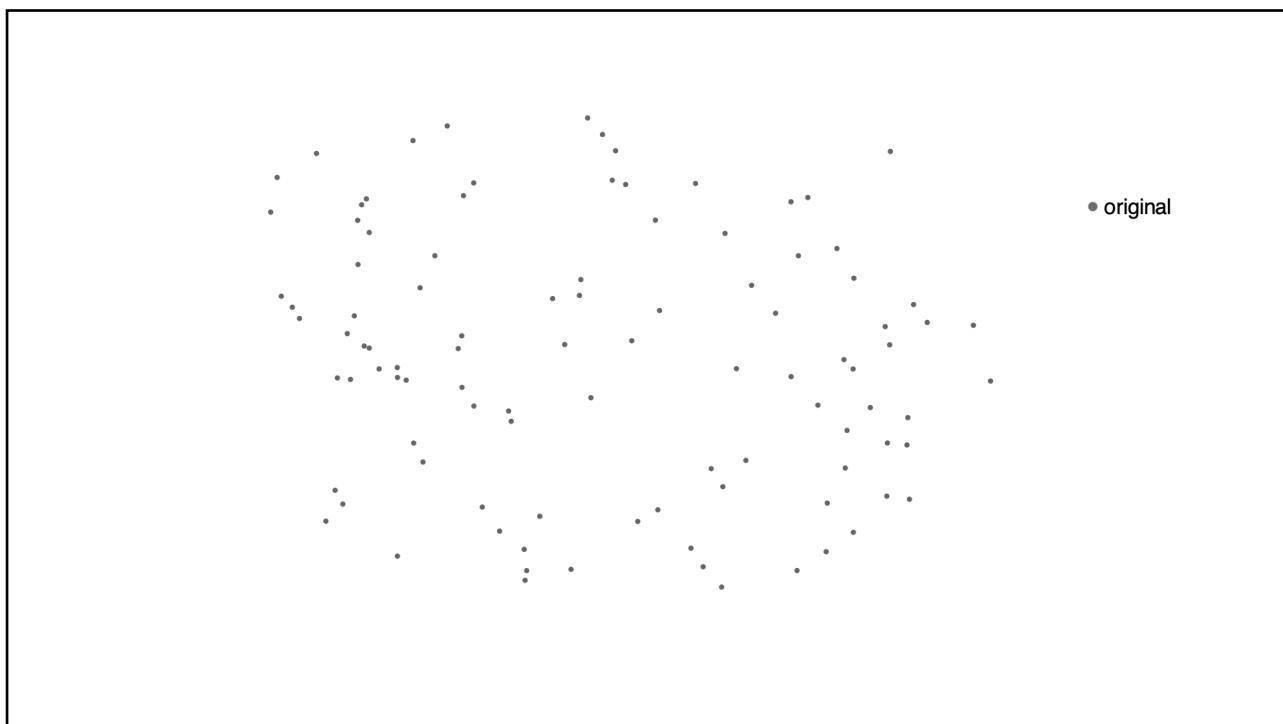
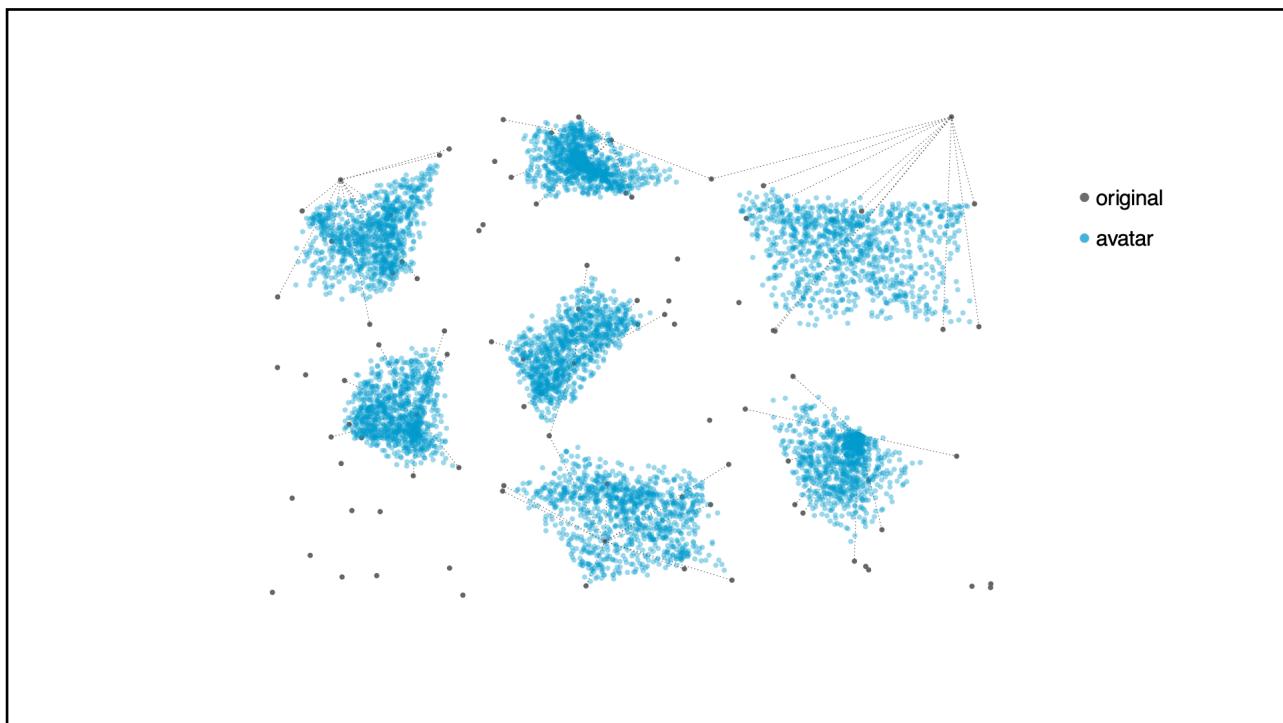


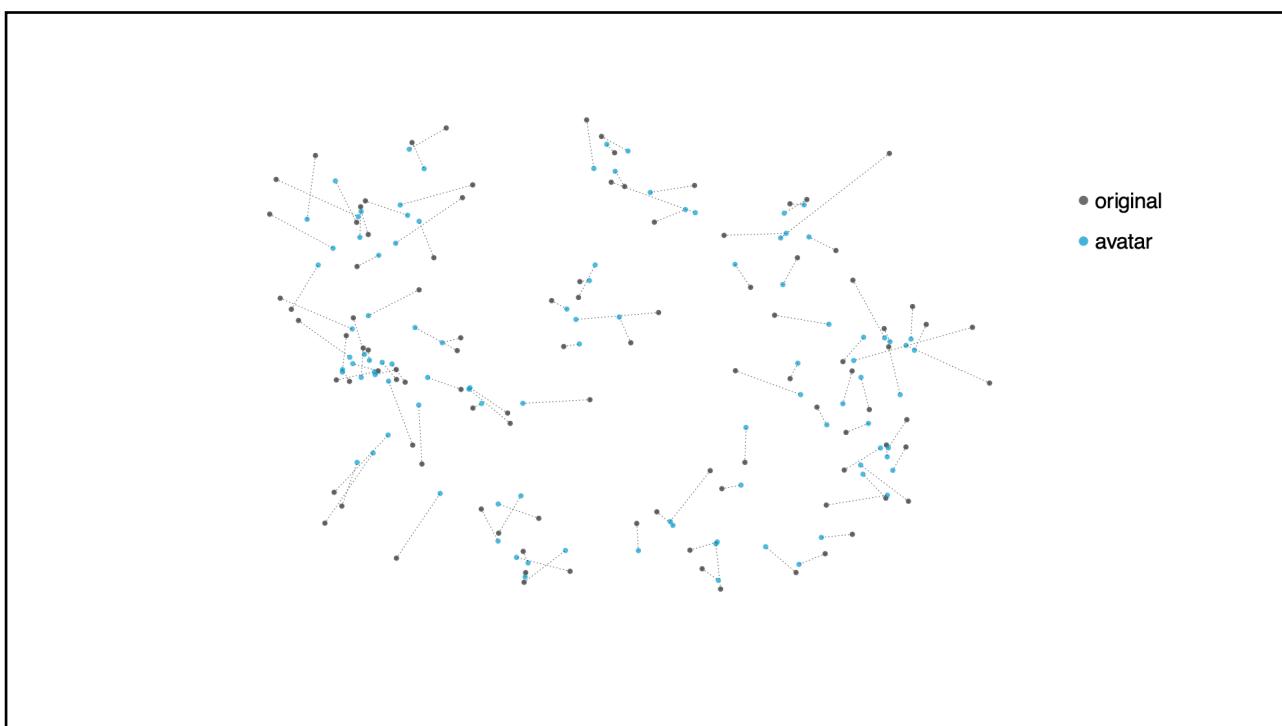
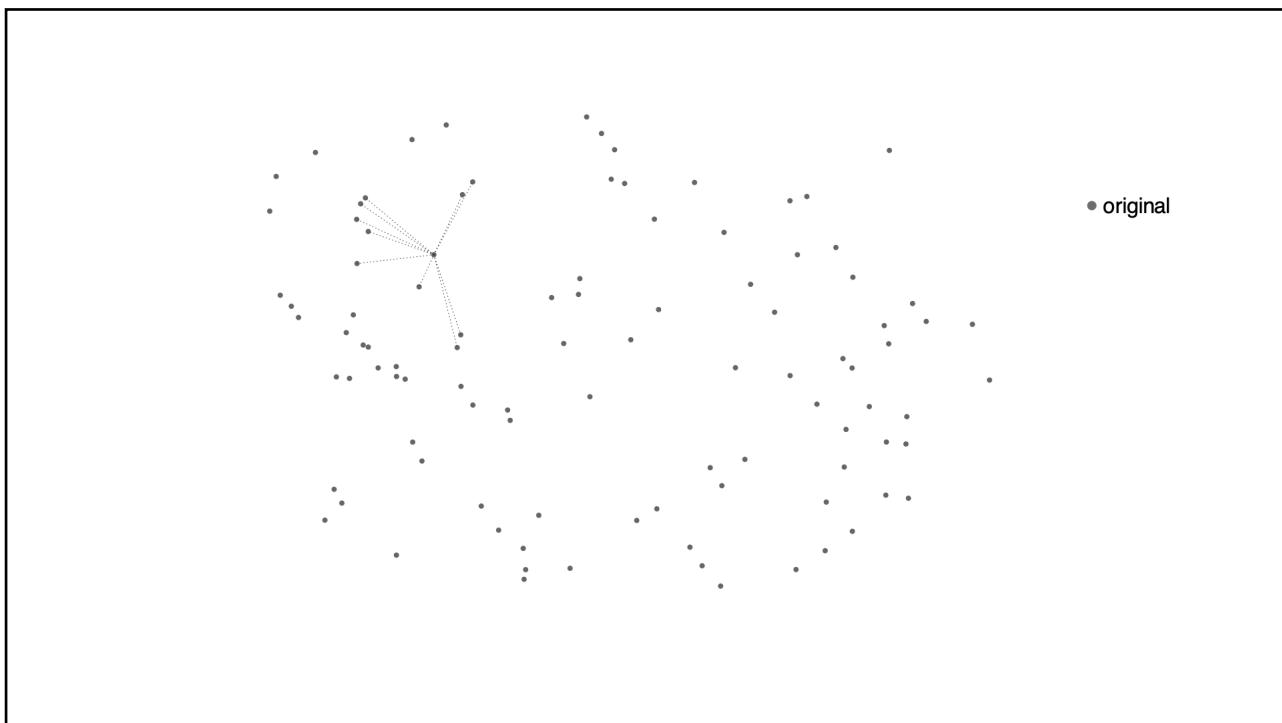


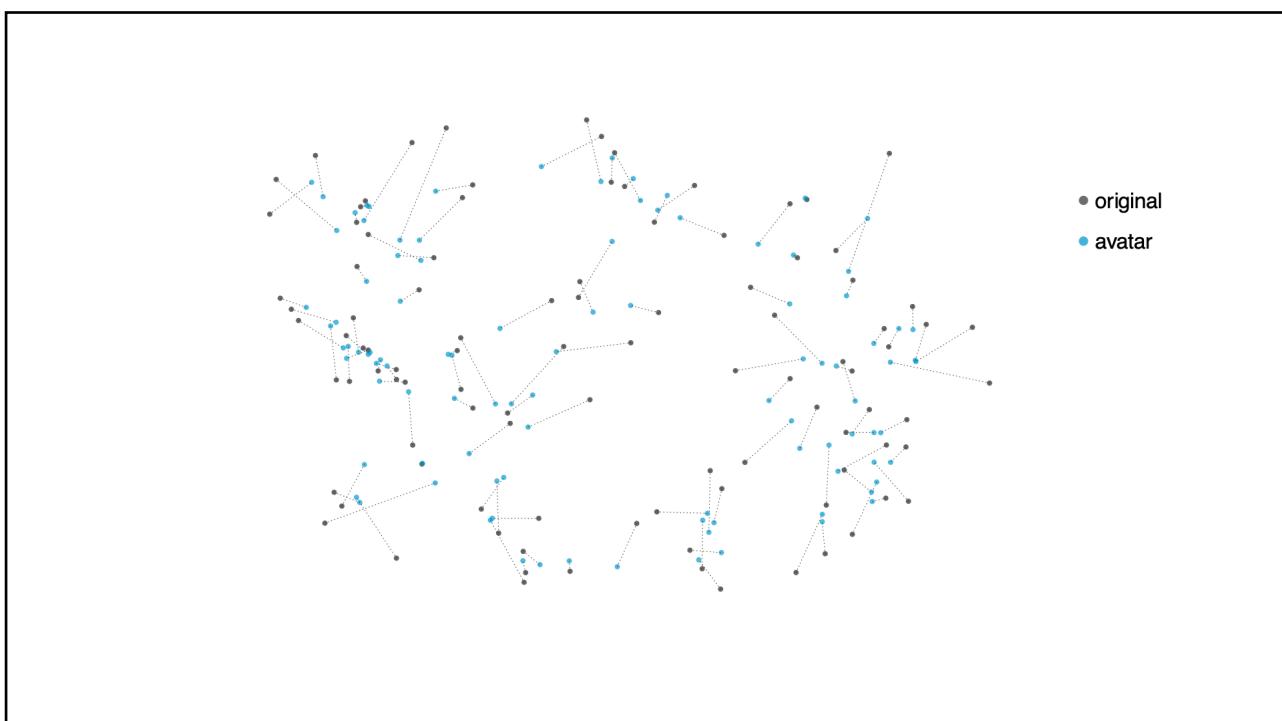
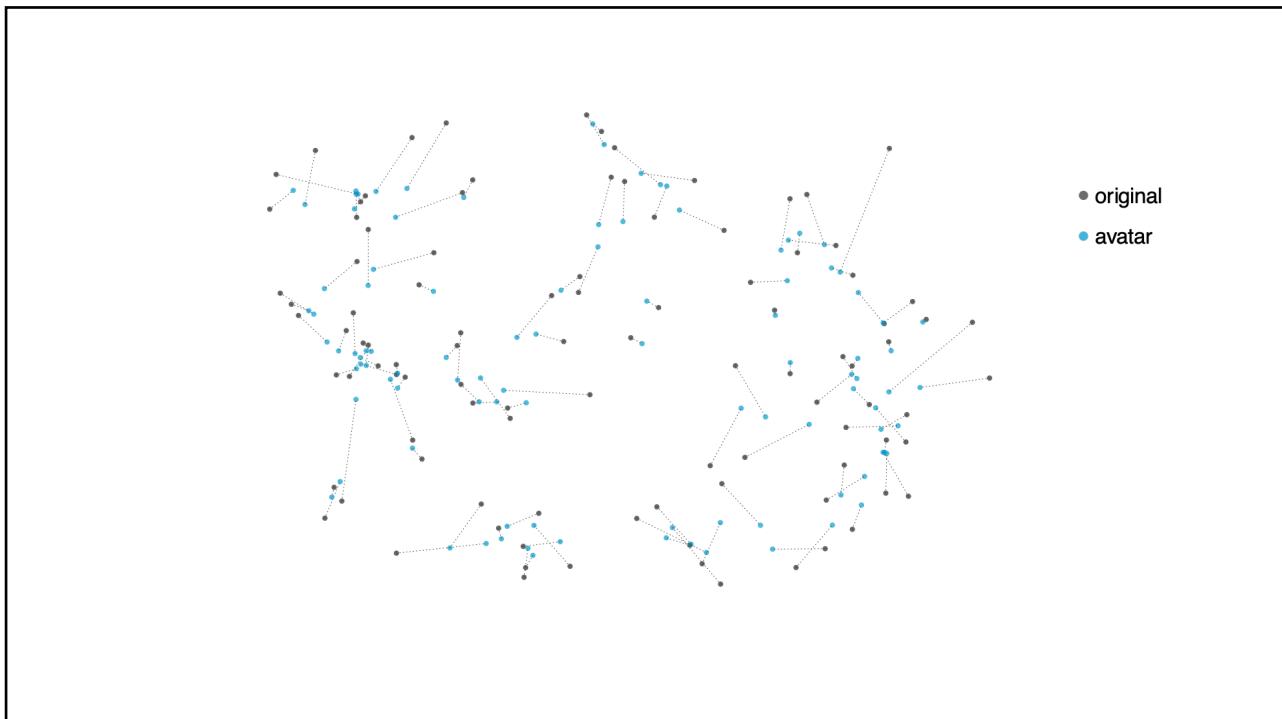


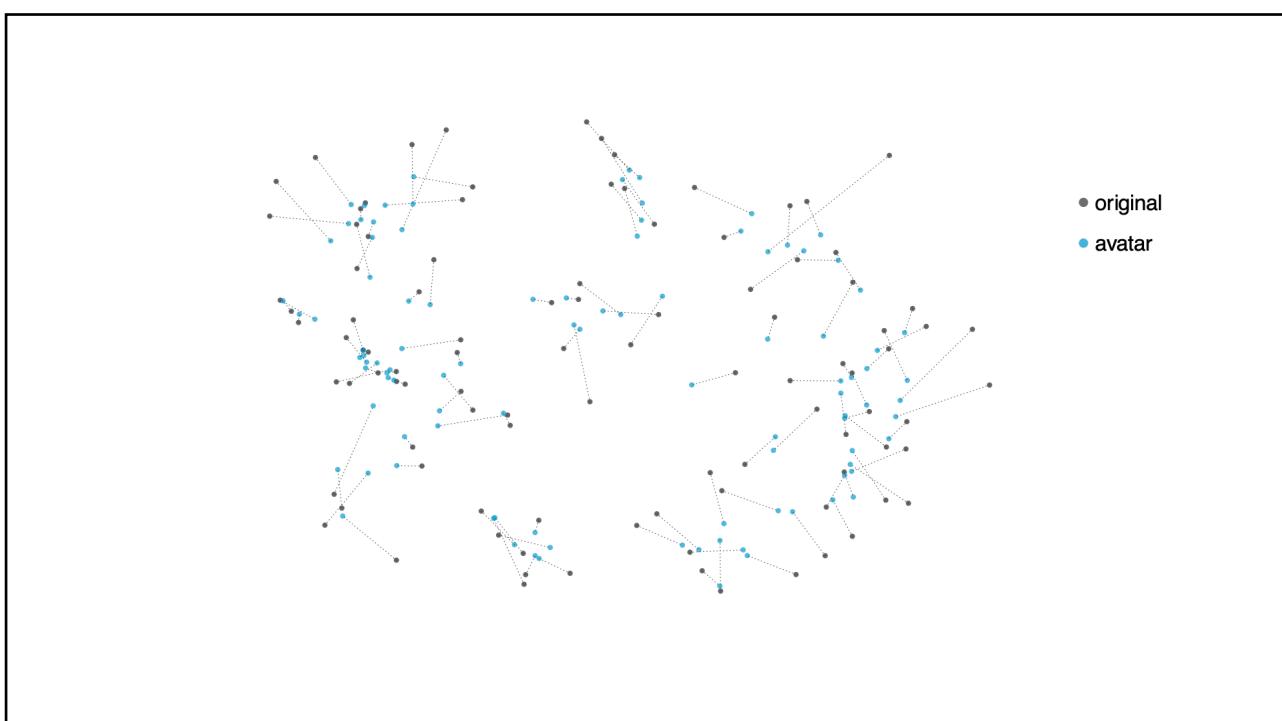
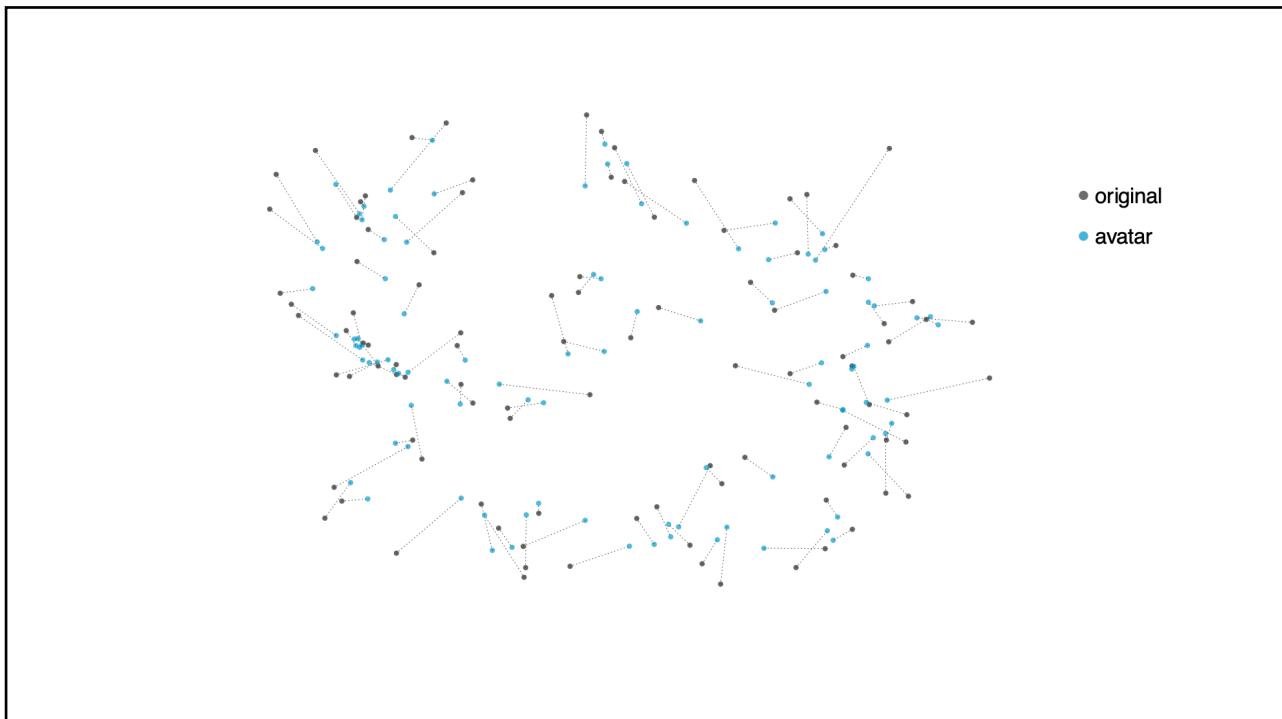


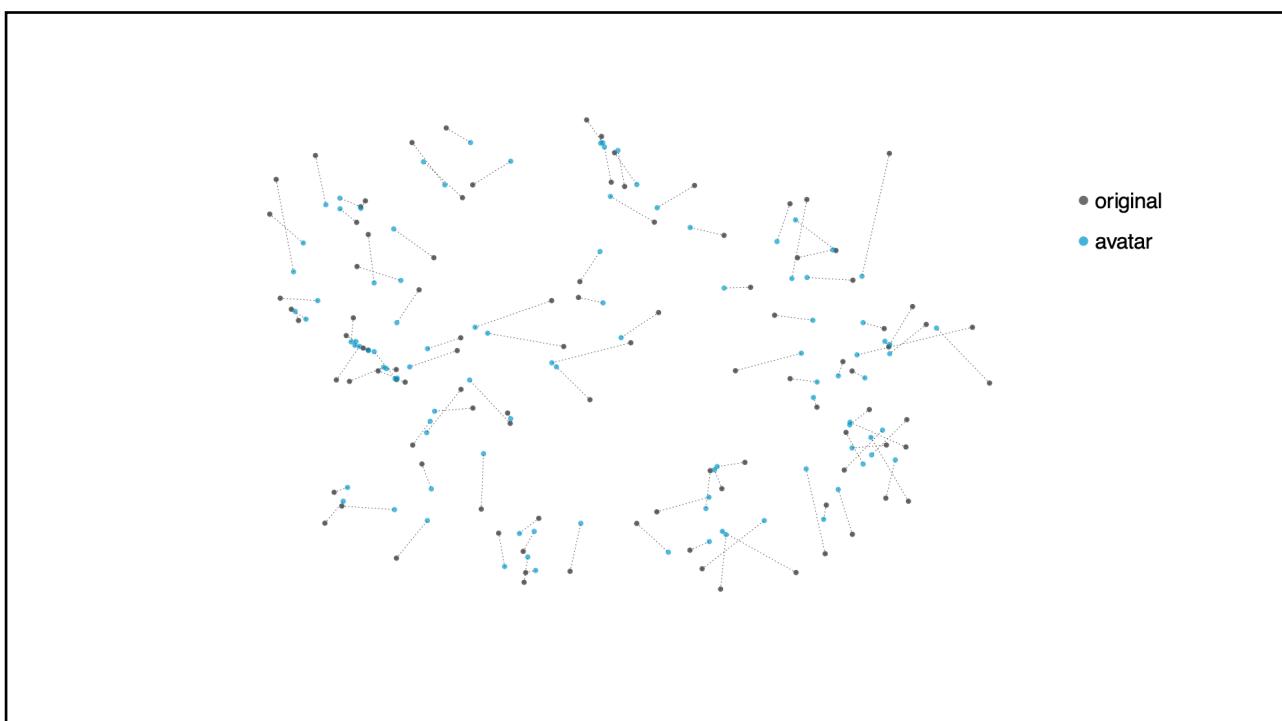
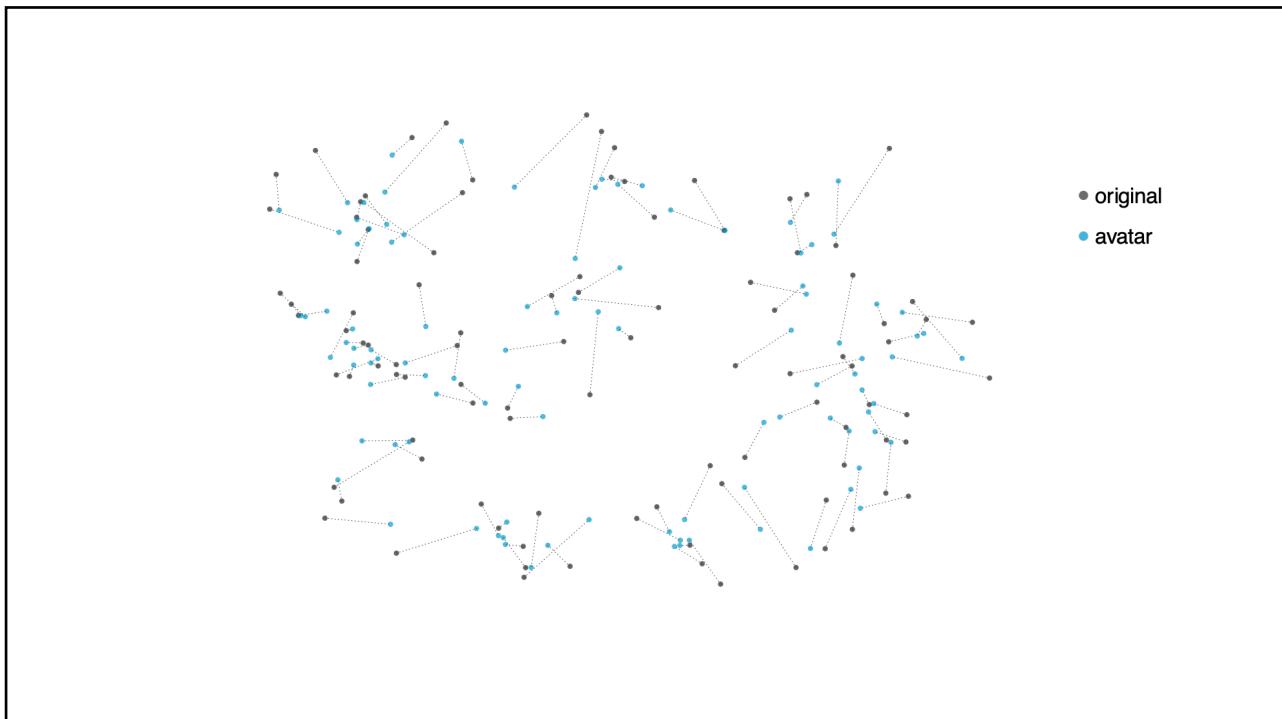












PRIVACY is the key



- Evaluated through quantitative metrics

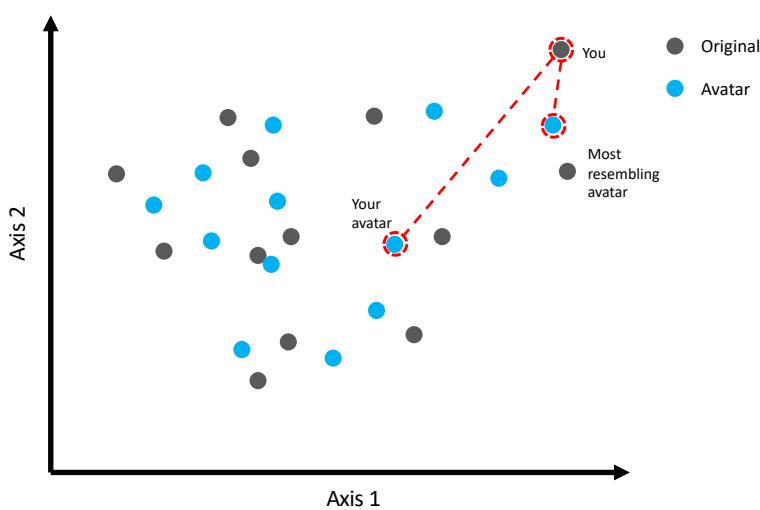
Before: Identify most at-risk records

During: define avatarization parameters

After: verify that the link between the record and the avatar is unmanageable to establish

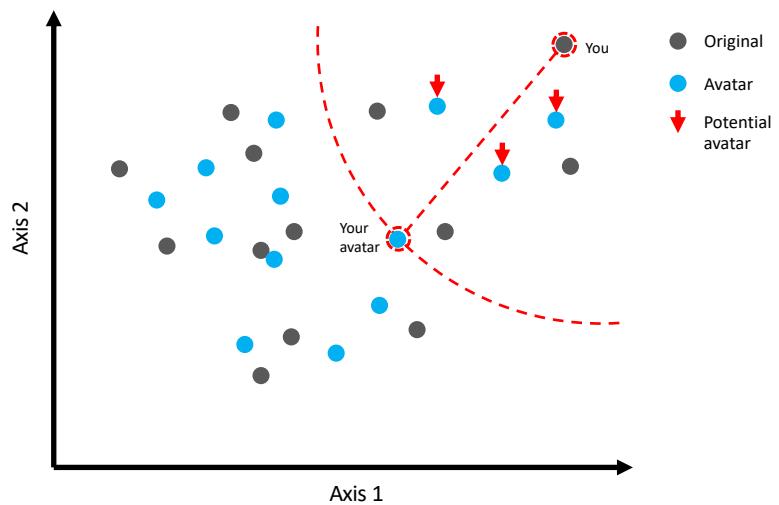
53

RISK ASSESSMENT



- Avatars to quantify the risk thanks to metrics
- These metrics can't be used with other method
- Right First Hit = protection percentage
- Proof for data owner that resembling avatar not been generated from his data

RISK ASSESSMENT



- Radius KNN = local cloaking
- The more avatars inside the radius, the stronger is the protection

AIDS
~ EXAMPLE
1 ~

BACKGROUND

Aids dataset

Trial comparing nucleoside monotherapy with three other combination therapy in HIV-infected adults

Dataset:

2139 individuals in the cohort
27 variables

57

COHORT

Aids dataset

wtkg	hemo	homo	drugs	cd40	...	cd420	days	arms
89.8128	FALSE	FALSE	FALSE	422	...	477	948	2
49.4424	FALSE	FALSE	FALSE	162	...	218	1002	3
88.452	FALSE	TRUE	TRUE	326	...	274	961	3
85.2768	FALSE	TRUE	FALSE	287	...	394	1166	3
66.6792	FALSE	TRUE	FALSE	504	...	353	1090	0

<https://clinicaltrials.gov/ct2/show/NCT00000625>

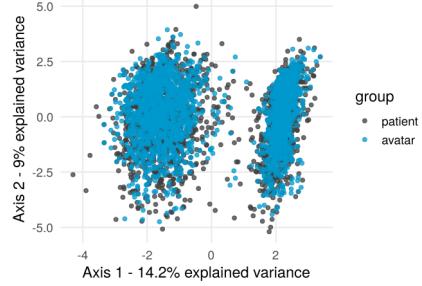
58

RESULTS

Aids dataset

Individuals with **marginal values** and by essence easy to re-identify are **protected** after avatarization

30 resembling patients ($K = 30$)



59

PRIVACY

Aids dataset

Almost **95%** of records do not produce the closest avatar

On average, a **record** is surrounded by **14** avatars:

It is therefore **unmanageable** for an attacker to establish the link between the record and its avatar

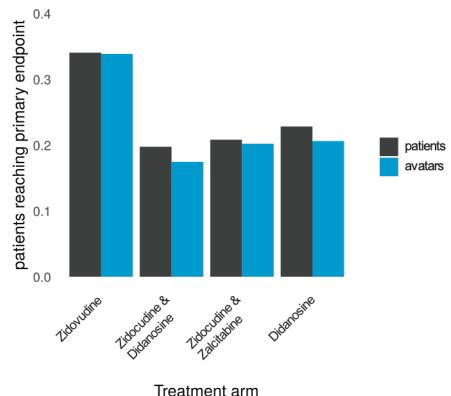
60

RESEARCH

Aids dataset

- The same conclusions can be drawn from the avatar dataset

Patients reaching primary endpoint depending on treatment



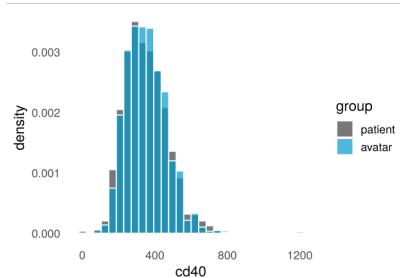
61

DISTRIBUTIONS

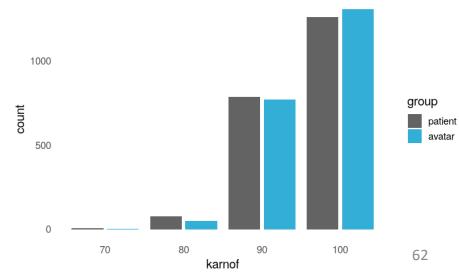
Aids dataset

- Distributions are also mainly conserved
- Marginal values are recentered

Continuous variable



Categorical variable



62

RELATIONS

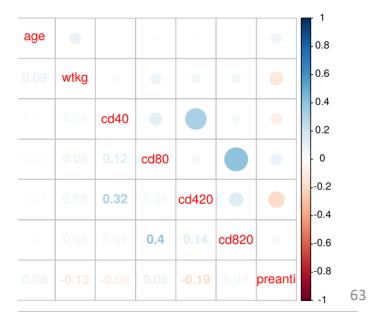
Aids dataset

- Correlations between variables remain comparable
- The avatarization seems to have **preserved signal** of the original dataset

Original dataset



Avatars dataset



BODY FAT

~EXAMPLE 2~

BACKGROUND

The objective of this study was to **evaluate** the calculation of percentage body fat from other simple body measurement

Dataset:

252 individuals in the cohort

15 variables

65

COHORT

Body fat dataset

siri	age	weight	height	adipos	...	chest	abdom	hip
22.8	31	148	70.75	21.6	...	88.5	83.5	94.5
14.9	68	179	68.25	27.2	...	100	91.6	96.4
27.2	34	205.75	71.25	27.6	...	105.2	98.8	108.3
25.3	22	176.75	70.5	25.4	...	101.8	87.2	98.5
28	52	219.15	74.25	26.1	...	107.5	104.2	107.5

Fitting Percentage of Body Fat to Simple Body Measurements, RW Johnson, Journal of Statistics Education, 1996

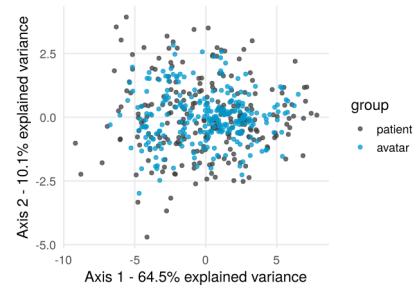
66

RESULTS

Body fat dataset

Individuals with **marginal values** and by essence easy to re-identify are **protected** after avatarization

25 resembling patients ($K = 25$)



67

PRIVACY

Body fat dataset

Almost **95%** of records do not produce the closest avatar

On average, a **record** is surrounded by **7 avatars**:

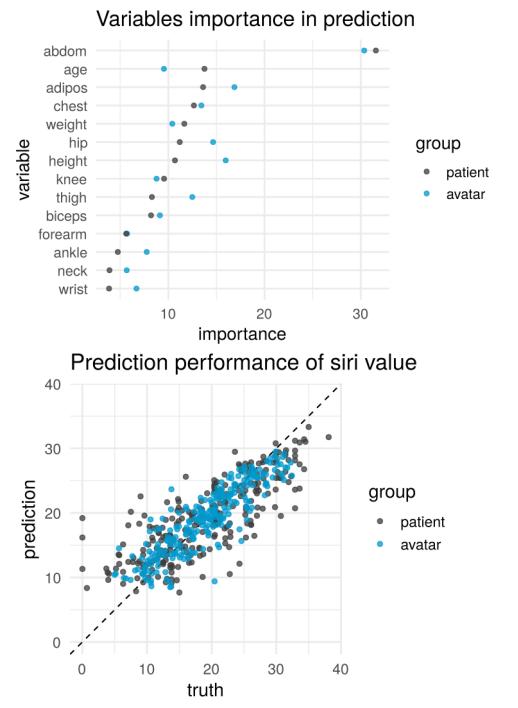
It is therefore **unmanageable** for an attacker to establish the link between the record and its avatar

68

RESEARCH

Body fat dataset

- Variables importance in prediction are comparable
- Prediction performance remains

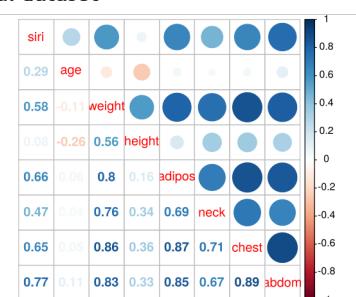


RELATIONS

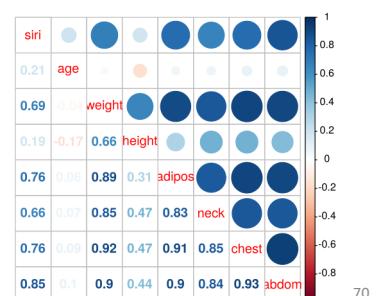
Body fat dataset

- Correlations between variables remain comparable
- The avatarization seems to have **preserved signal** of the original dataset

Original dataset



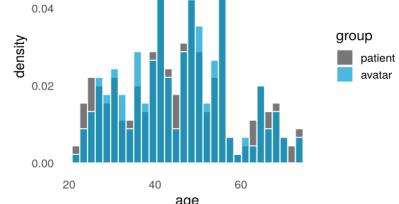
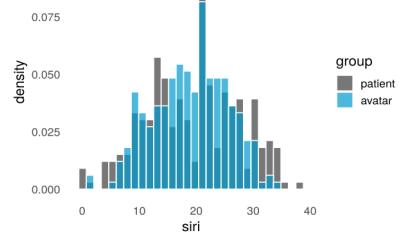
Avatars dataset



DISTRIBUTIONS

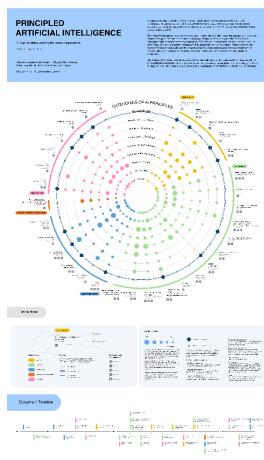
Body fat dataset

- Distributions are also mainly conserved
 - Marginal values are recentered



71

Conclusion : Bottom-up “ethics by-construction”



<https://ai-hr.cyber.harvard.edu/primp-viz.html>

- The systematic use of synthetic data (avatars) in our biomedical data warehouse is an ethical question of governance.
 - Over 70 top- down approaches on AI ethics
 - Finish line ?
 - Collective building and revision of a set of ethical principles
 - Starting point for diverse public-private community
 - Contradictory principles promote ethics by construction



Conclusion

- **Les Technologies créent des opportunités techniques**

- Tension existentielle dans le saisissement de ces outils.



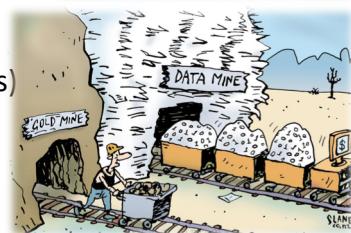
- **Un rôle pour les données de santé ?**

- Exceptionnalisme? Exemplarité ?
- Intentionalité, consentement, transparence ?
- Santéisme ?



- **Les entrepôts de données biomédicales**

- “Mine” de données (multiples jeux de mots possibles)
- Usage secondaires, non intentionnels
- Opportunité de Soli-data-rité



“La Clinique des données”

Pr Pierre-Antoine GOURRAUD

« Un check-up de l'IA pour la santé numérique, c'est grave docteur(s) ? »

Séminaire Aristote 27 février 2020,

École Polytechnique, Amphithéâtre Arago
91120 Palaiseau

ATIP-Avenir Team 5 «Translational Immunogenomics of Transplantation and Autoimmunity »,
ITUN - CRTI - UMR Inserm 1064 -CHU de Nantes

Pôle Hospitalo-Universitaire 11 : Santé Publique, Santé au Travail et Pharmacie,

Hôpital St-Jacques - CHU de Nantes

