

# Pseudonimisatieservices voor wetenschappelijk onderzoek

Houten, 5 november 2019

# Welkom!

## Even voorstellen...

**Simone van Wijngaarden**

Adviseur, CIPP/e, Stichting ZorgTTP

### **Dienstverlening Stichting ZorgTTP**

- Onomkeerbare pseudonimisatie d.m.v. ons pseudonimisatieplatform
- Omkeerbare pseudonimisatie d.m.v. onze encryptiedienst TRES (Trusted Reversible Encryption Service)
- Data Protection Impact Assessment (DPIA) / gegevensbeschermingseffectbeoordeling

### **Maar ook:**

- Integratie van meerdere diensten
- Nieuw product/andere invulling opdracht met bestaande dienst als basis
  - Zoals de pseudonimisatievoorziening voor wetenschappelijk onderzoek waar we het vandaag over gaan hebben

# Welkom!

# Even voorstellen...

**Jan Lucas van der Ploeg**

Data Engineer, afdeling Informatiemanagement Onderzoek, UMCG

**Robert Griffioen**

Projectleider de-identificatie project, SURF

# Waar gaan we het over hebben?

## **Pseudonimisatieservices voor wetenschappelijk onderzoek**

- Introductie pseudonimisatievoorziening
- Noodzaak tot pseudonimisatievoorziening
- Toelichting op de eerste vier use cases
- Release 1
- Opschaling naar SURF
- Vragen en afronding

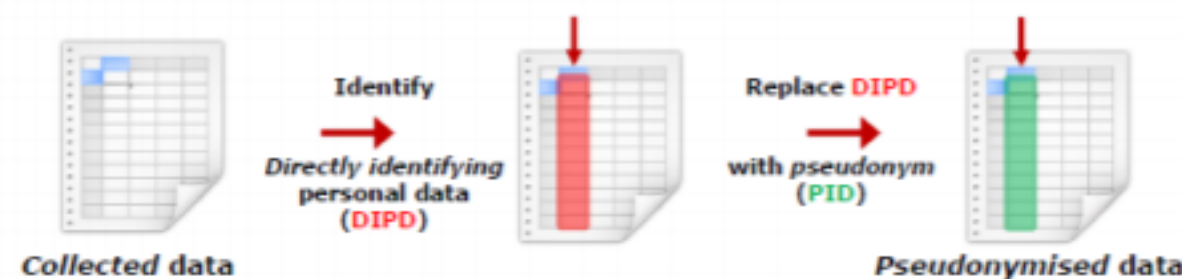
*Vragen tijdens de presentatie? Stel ze gerust!*

## Pseudonymisation Service for Research

Jan Lucas van der Ploeg and Francisco Romero Pastrana

Data Federation Hub - University Medical Center Groningen- University of Groningen

### Pseudonymisation



**Pseudonymisation** means that direct identifiable identifiers such as a name, date of birth and address, are replaced with a pseudonym.

Pseudonymisation involves separating directly identifying personal data from *substantive* data, optionally maintaining a link through an arbitrary key. The GDPR explicitly mentions pseudonymisation as one approach for GDPR requirements compliance, increasing the privacy and security of personal data processing.

#### Pseudonymisation ≠ Anonymization

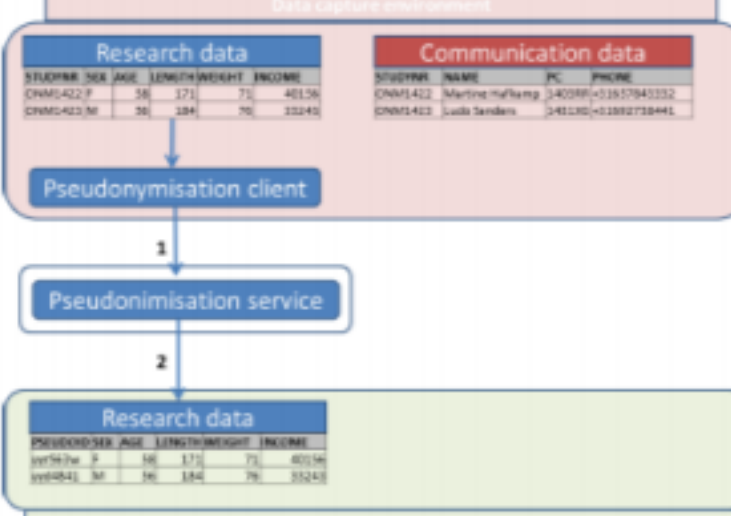
**Pseudonymisation** is one step in data process for research. Pseudonymised data is not necessarily anonymous. Re-identification is possible, because of indirect identifiers.

#### Pseudonymisation Service for research:

Pseudonymisation is not a trivial process. The UMCG and UG developed a pseudonymisation service to support researchers with pseudonymisation and linking datasets. The service provides software and a support desk for pseudonymisation, linking and anonymization.

### Pseudonymisation: No data linkage possible

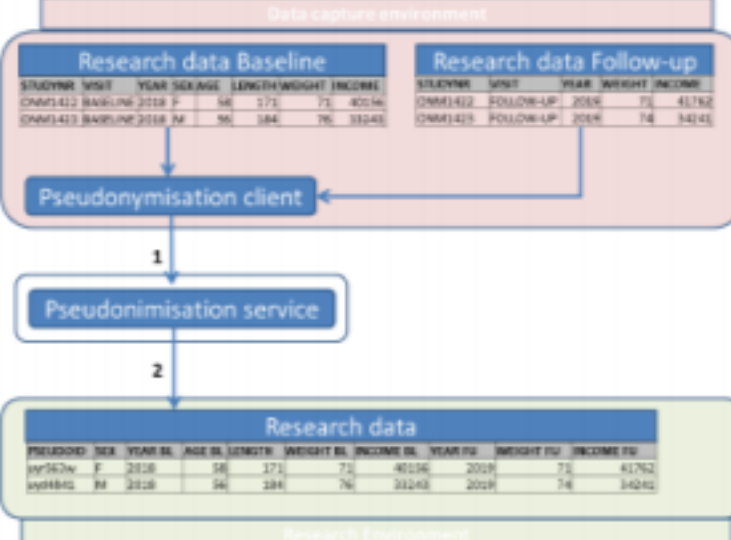
#### For researchers, self service



- Small studies
- Irreversible
- Different pseudonyms for different sessions
- Not possible to link published datasets
- Creating research pseudonyms

### Pseudonymisation: Data linkage possible

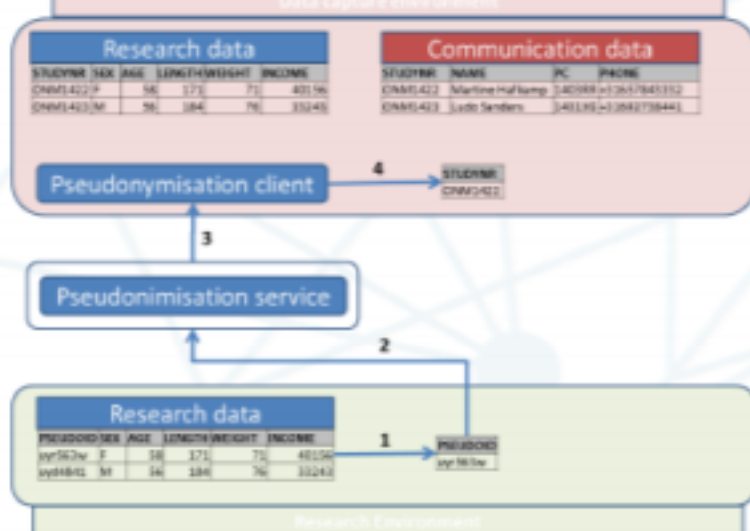
#### For data stewards/advanced data managers, no self service



- Longitudinal studies, large cohort studies and biobanks
- Subjects are traceable over time
- Pseudonym for subject is stable in project over different sessions

### Process for reversible pseudonymisation

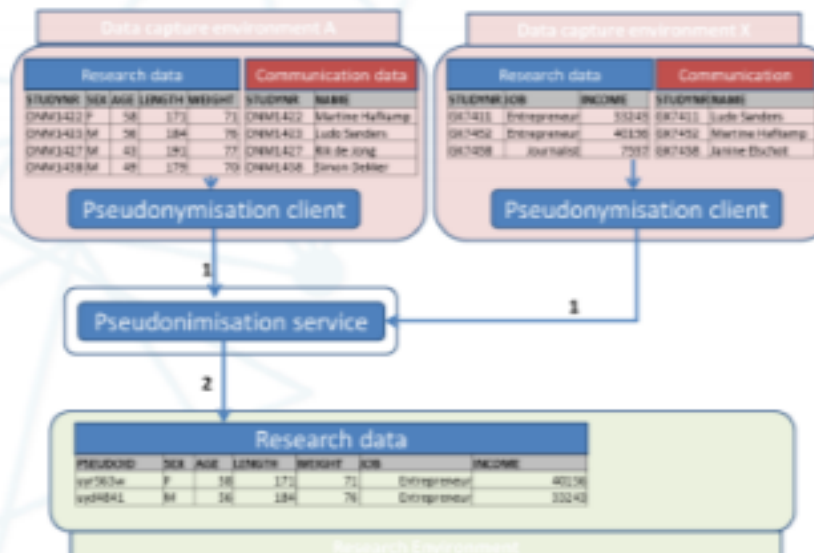
#### For data stewards/advanced data managers, no self service



- Only when necessary
- Not the default option, extra procedures and contracts when needed
- Incidental findings

#### Pre matching by Trusted Third Party

#### By Trusted Third Party (TTP), no self service



- Linking data sets from different sources
- Pre matching and linkage by TTP
- Only researcher knows overlap between sets

### Pseudonymisation service for Research

#### First version of software available soon, with support desk for pseudonymisation, data linkage and anonymization

The service makes it as simple as possible for researchers to pseudonymise their data in a secure manner. The pseudonymisation service contributes to high-quality linking of health data for research that complies with legislation.

- Only direct identifiers are pseudonymised.
- Only for quantitative data
- For next releases:
  - Pseudonymisation of subjects in research data set (e.g. names of doctors or nurses)
  - Pseudonymisation of keys (maintaining data integrity)
  - Other types of data (qualitative, video, audio, genetic)
  - Anonymization

#### The Data Federation Hub / Human Data

The University Medical Center Groningen and the University of Groningen joined efforts to set up an integrated research support platform: the Data Federation Hub/Human Data. The DFH/Human Data support the whole research data lifecycle while ensuring data security and efficiently protecting the privacy of participants.

# Pseudonimisatieservices voor wetenschappelijk onderzoek

2016 – 2018

Uitvoerige gesprekken over use cases RUG-UMCG

Midden 2018

Eerste contact UMCG-ZorgTTP

Q3-Q4 2018

Vooronderzoek

Eind 2019

Demo voor eerste use case

- Pseudonimisatieservice voor kleine studies

Begin 2019

Start requirements en ontwikkeling tweede use case

- Pseudonimisatieservice voor longitudinale studies met EPD data

Midden 2019

Vervolg van de oplossingen via SURF dankzij gewonnen aanbestedingstraject

Eind september 2019

Oplevering eerste release

- Pseudonimisatieservice voor longitudinale studies

Q4 2019 – Q1 2020

Start requirements en ontwikkeling volgende use case

- Pseudonimisatieservice voor kleine studies

# Pseudonimisatieservices en Use Cases

Jan Lucas van der Ploeg

Data Engineer, afdeling Informatiemanagement Onderzoek, UMCG

# Pseudonimisatieservice Groningen

## **Samenwerking UMCG en RUG**

- Onderdeel project Human Subject Research van UMCG en RUG

## **Onderscheid maken (in proces/stap en techniek)**

- Pseudonimiseren
- Anonimiseren
- Koppelen

## **Zowel techniek als organisatie/proces**

- Service met software plus support/consultancy

# Scope en fasering

## Scope bepalen

- Welke data pseudonimiseren?
- Direct identifiers, ook indirect identifiers, quasi identifiers, keys?
- Alleen research subjecten of ook andere personen (zorgverleners)?
- Kwantitatief, ook kwalitatief?
- Gestructureerde data, images, video, audio, geo-data, genetische data?

## Gefaseerd

- In releases



# Context en prioritering

## **Kleine studies**

- Veel ( $\pm$  1000 per jaar)
- Geen budget
- Weinig kennis pseudonimiseren & anonimiseren
- Scheiden directe persoonsgegevens van onderzoeksdata
- Grootste risico

## **Grote studies/projecten**

- Wel budget
- Kennis en personeel aanwezig
- Maatwerkoplossingen beschikbaar
- Vaak koppelen aan andere bronnen

# Eerste release van de service

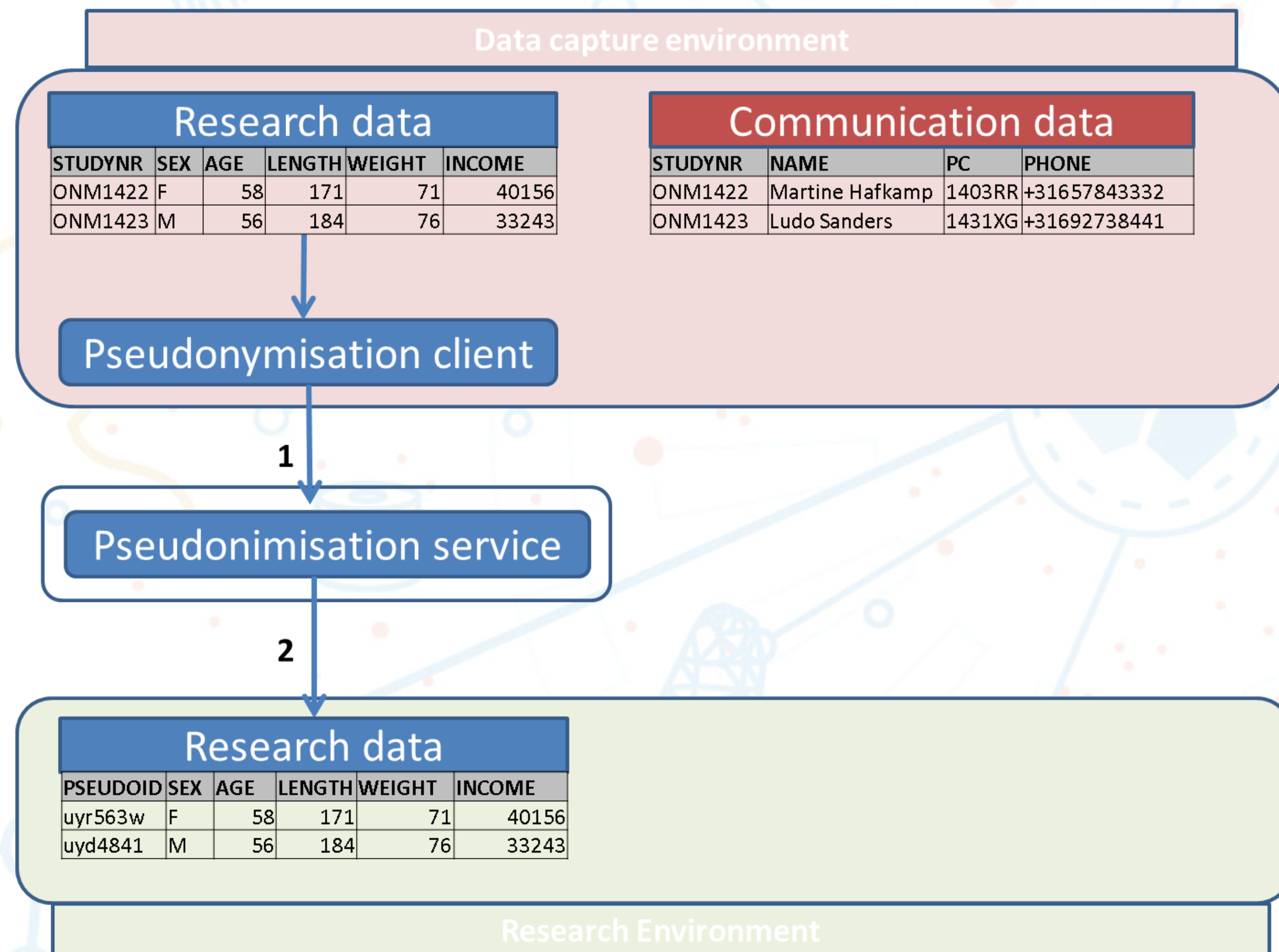
## Scope

- Software for pseudonymisation
  - Only direct identifiers are pseudonymised
  - Only for quantitative data
- Support for other data and for anonymization and data linkage

## 4 Use cases

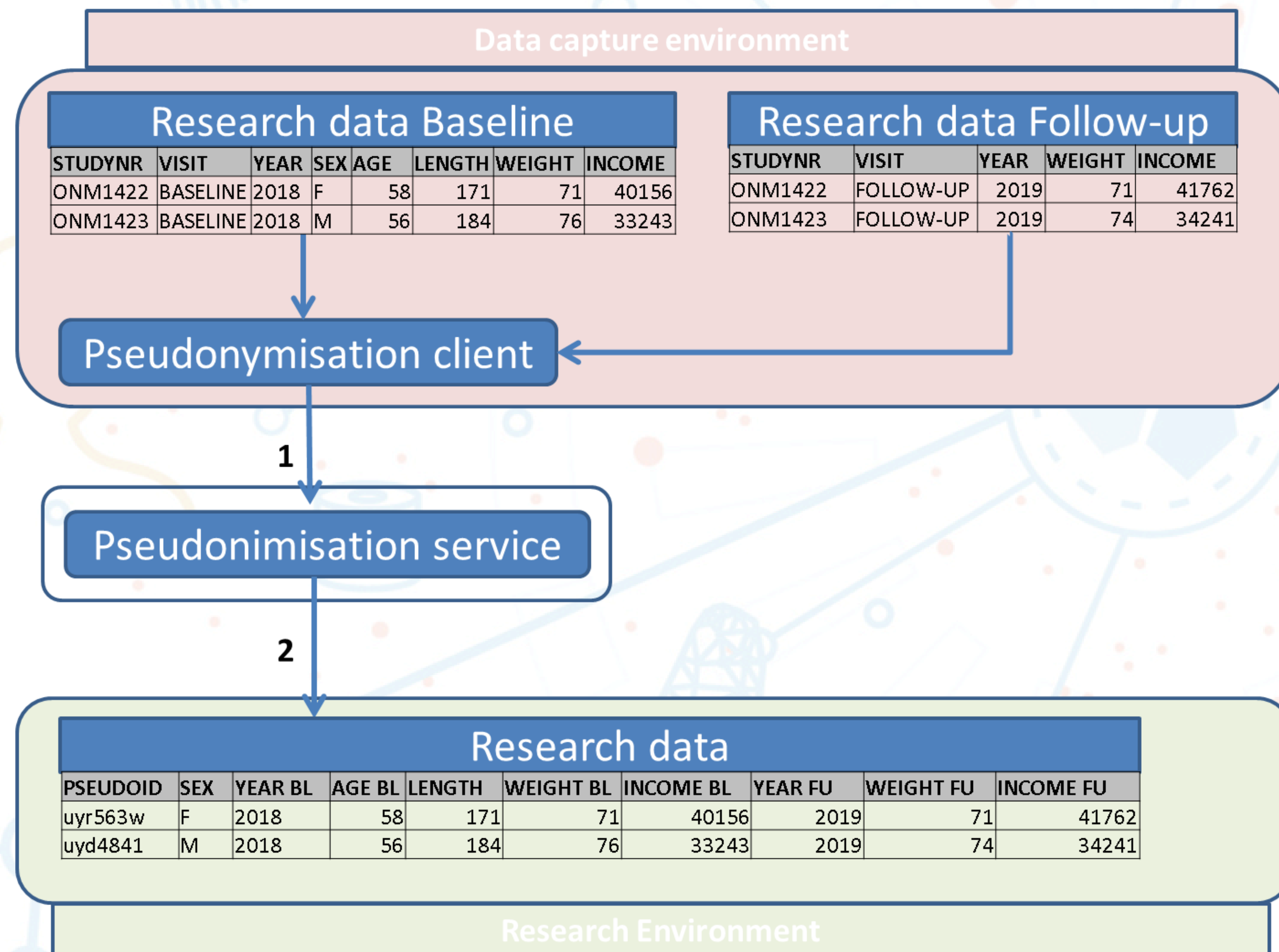
- Pseudonymisation: No data linkage possible
- Pseudonymisation: Data linkage possible
- Process for reversible pseudonymisation
- Pre-matching (and linking) by Trusted Third Party

# Pseudonimisatie: geen data linkage



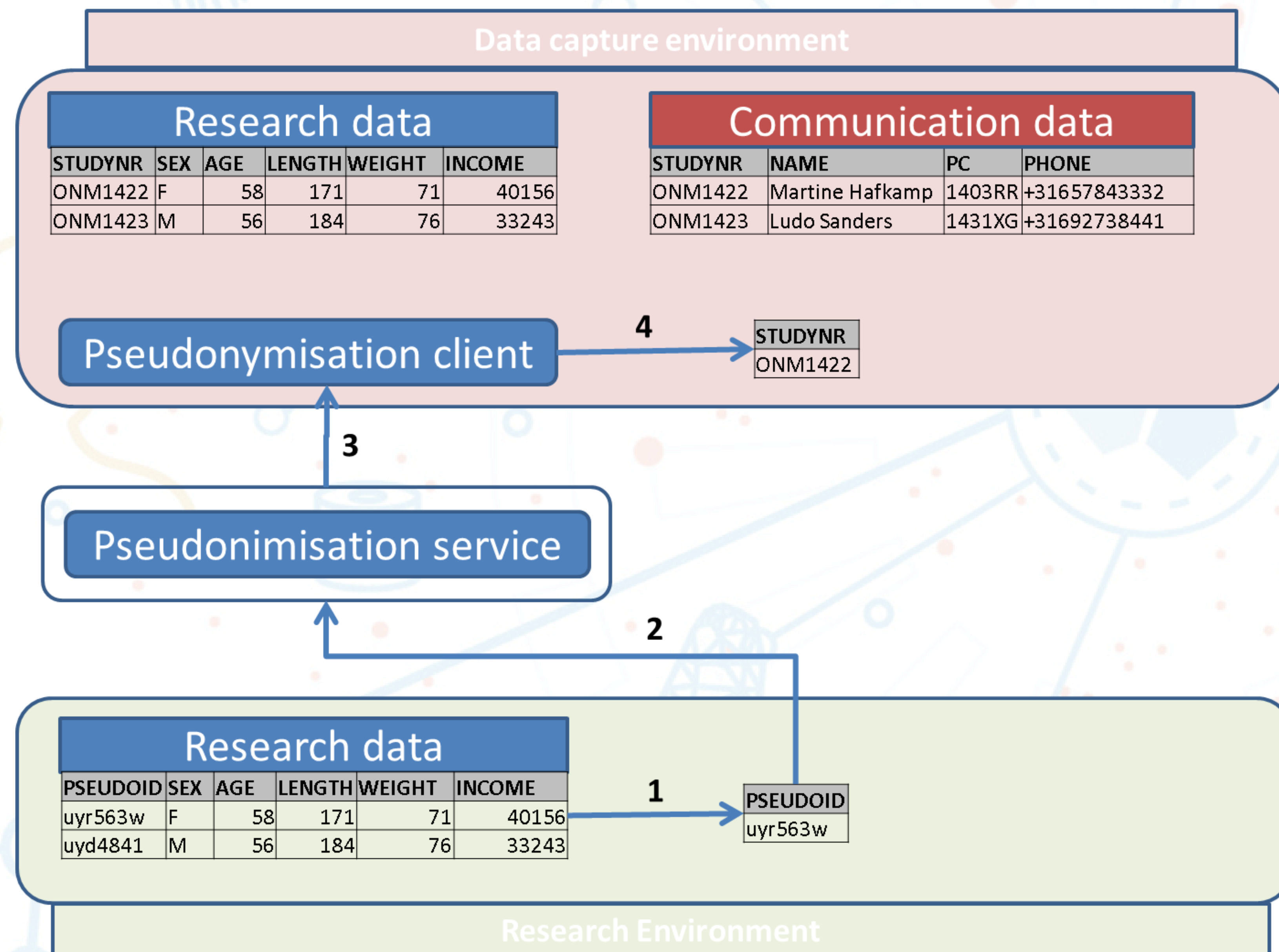
- For researchers, self service
- Small studies
- Irreversible (default)
- Different pseudonyms for different sessions (default)
- Not possible to link published datasets
- Creating research pseudonyms

# Pseudonimisatie: data linkage mogelijk



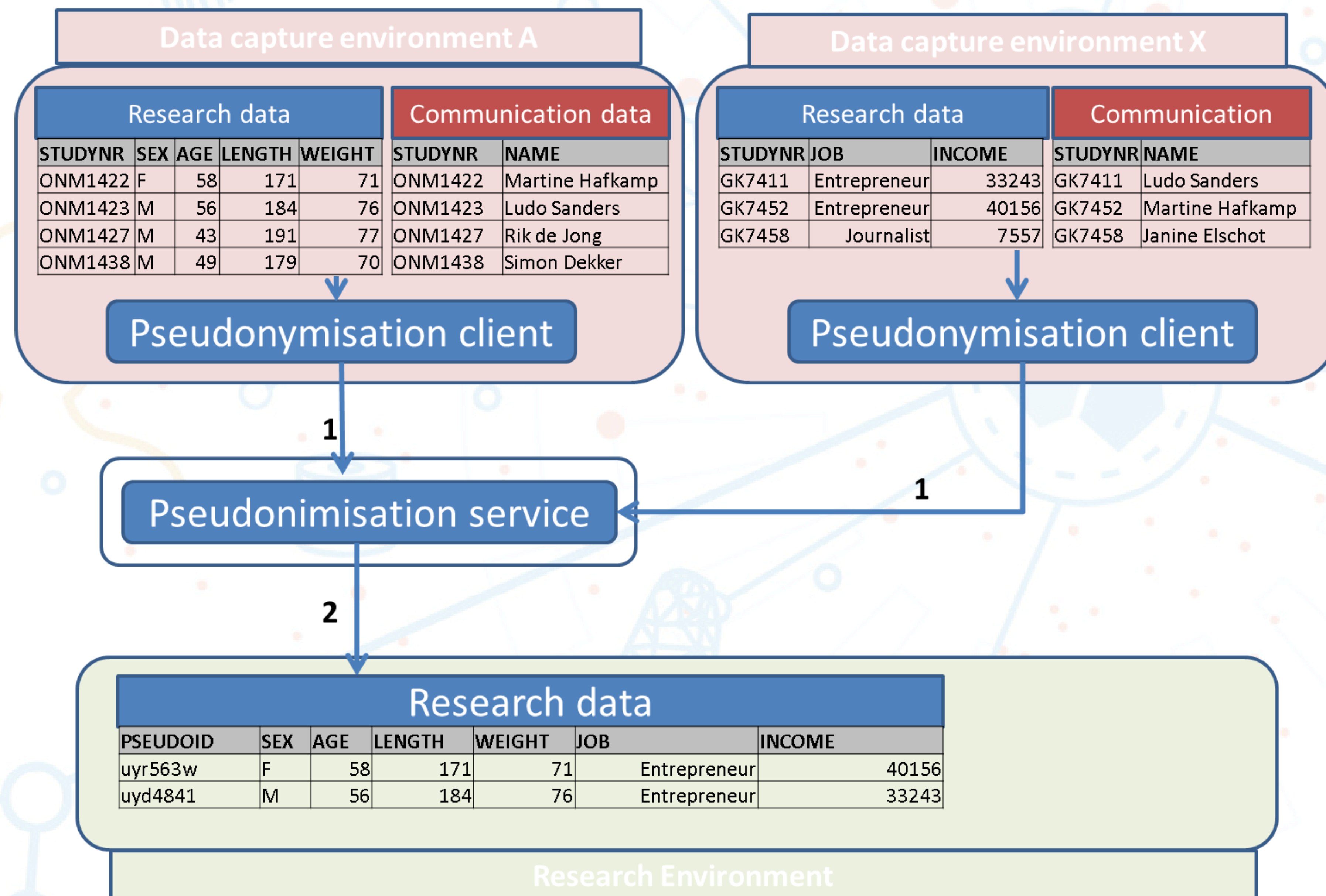
- For data stewards/advanced data managers, no self service
- Longitudinal studies, large cohort studies and biobanks
- Subjects are traceable over time
- Pseudonym for subject is stable in project over different sessions

# Proces voor omkeerbare pseudonimisatie



- For data stewards/advanced data managers, no self service
- **Only when necessary**
- Not the default option, extra procedures and contracts when needed
- Incidental findings

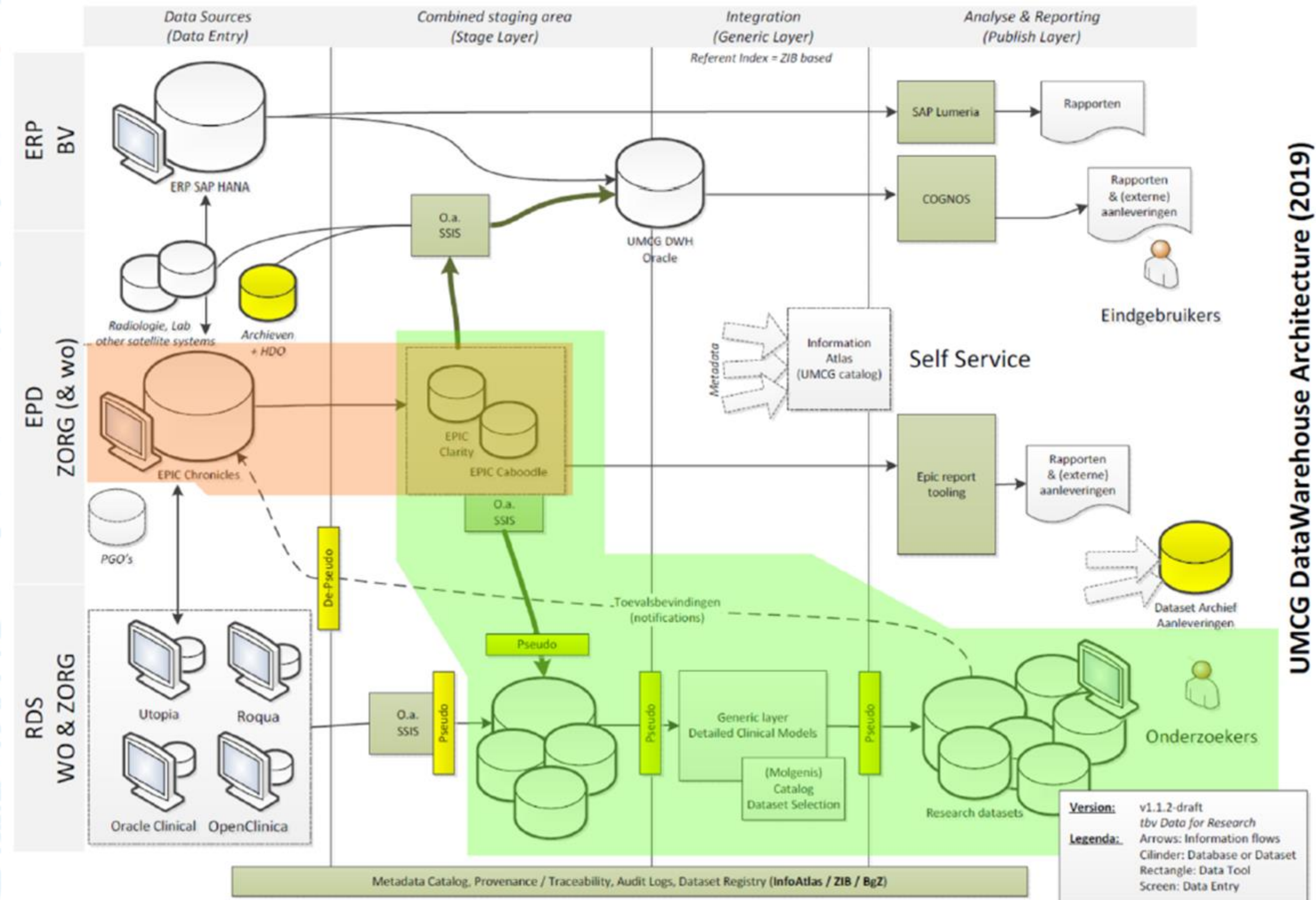
# Prematching en linkage door Trusted Third Party



- By TTP, no self service
- Linking data sets from different sources
- Pre matching and linkage by TTP
- Only researcher knows overlap between sets

# 2019: Prioriteit EPD data voor onderzoek

- Pseudonimisatie service is voorwaarde voor EPD data voor onderzoek
- Prio 1<sup>ste</sup> release software
  - 2<sup>e</sup> use case; longitudinale studie met stabiele pseudoniemen



UMCG Data Warehouse Architecture (2019)

# **Release 1: Pseudonimisatieservice voor longitudinale studies**

Simone van Wijngaarden  
Projectleider, Stichting ZorgTTP



**ZorgTTP**  
Privacy & vertrouwen



# Vereisten pseudonimisatieservice longitudinale studies (use case 2)

- REST API
- FHIR bundle
- JSON
- Self service
- Onomkeerbare pseudoniemen
- Pseudoniemen
  - Zelfde input, andere studie = ander pseudoniem
  - Stabiele pseudoniemen binnen studie
  - Unieke pseudoniemen
- Limieten
- Responstijd
- Geteste oplevering
- Eindgebruikersdocumentatie



## Pseudonymisation Service for Research

Jan Lucas van der Ploeg and Francisco Romero Pastrana  
Data Federation Hub - University Medical Center Groningen - University of Groningen

### Pseudonymisation



**Collected data**  
Pseudonymisation means that direct identifiable identifiers such as a name, date of birth and address, are replaced with a pseudonym.  
Pseudonymisation involves separating directly identifying personal data from substantive data, optionally maintaining a link through an arbitrary key. The GDPR explicitly mentions pseudonymisation as one approach for GDPR requirements compliance, increasing the privacy and security of personal data processing.

### Pseudonymisation ≠ Anonymization

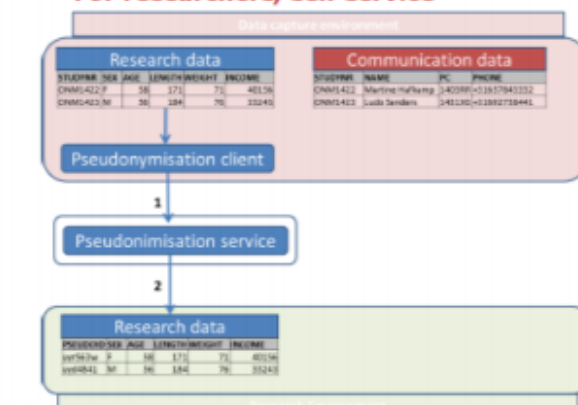
**Pseudonymisation** is one step in data process for research. Pseudonymised data is not necessarily anonymous. Re-identification is possible, because of indirect identifiers.

### Pseudonymisation Service for research:

Pseudonymisation is not a trivial process. The UMCG and UG developed a pseudonymisation service to support researchers with pseudonymisation and linking datasets. The service provides software and a support desk for pseudonymisation, linking and anonymization.

### Pseudonymisation: No data linkage possible

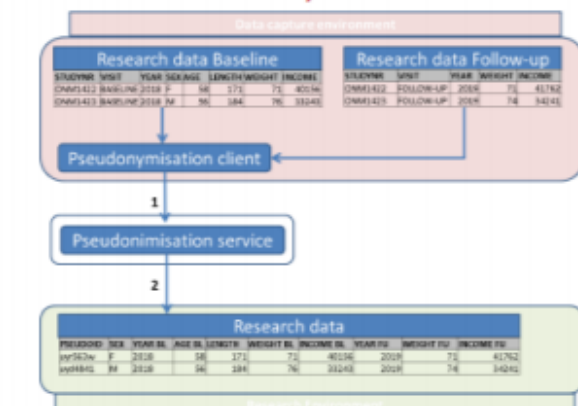
#### For researchers, self service



- Small studies
- Irreversible
- Different pseudonyms for different sessions
- Not possible to link published datasets
- Creating research pseudonyms

### Pseudonymisation: Data linkage possible

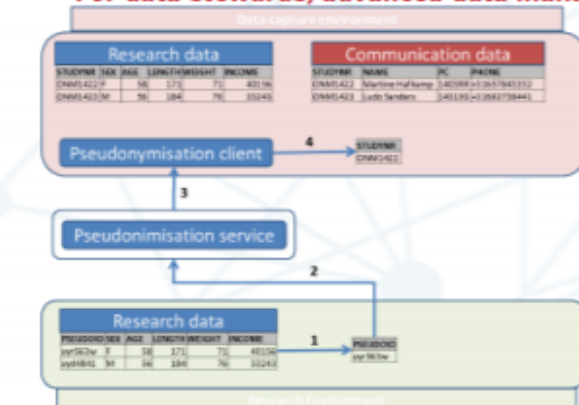
#### For data stewards/advanced data managers, no self service



- Longitudinal studies, large cohort studies and biobanks
- Subjects are traceable over time
- Pseudonym for subject is stable in project over different sessions

### Process for reversible pseudonymisation

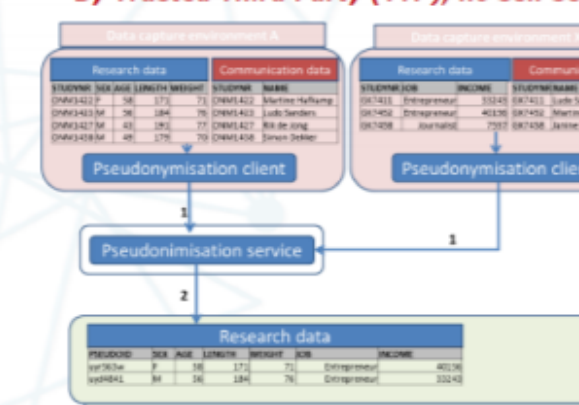
#### For data stewards/advanced data managers, no self service



- Only when necessary
- Not the default option, extra procedures and contracts when needed
- Incidental findings

### Pre matching by Trusted Third Party

#### By Trusted Third Party (TTP), no self service



- Linking data sets from different sources
- Pre matching and linkage by TTP
- Only researcher knows overlap between sets

### Pseudonymisation service for Research

#### First version of software available soon, with support desk for support for pseudonymisation, data linkage and anonymization

The service makes it as simple as possible for researchers to pseudonymise their data in a secure manner. The pseudonymisation service contributes to high-quality linking of health data for research that complies with legislation.

- Only direct identifiers are pseudonymised.
- Only for quantitative data
- For next releases:
  - Pseudonymisation of subjects in research data set (e.g. names of doctors or nurses)
  - Pseudonymisation of keys (maintaining data integrity)
  - Other types of data (qualitative, video, audio, genetic)
  - Anonymization

### The Data Federation Hub / Human Data

The University Medical Center Groningen and the University of Groningen joined efforts to set up an integrated research support platform: the Data Federation Hub/Human Data. The DFH/Human Data support the whole research data lifecycle while ensuring data security and efficiently protecting the privacy of participants.

[www.rug.nl/dfh](http://www.rug.nl/dfh)

[dfh@rug.nl](mailto:dfh@rug.nl)

[f.romero.pastrana@rug.nl](mailto:f.romero.pastrana@rug.nl)  
[j.l.van.der.ploeg@umcg.nl](mailto:j.l.van.der.ploeg@umcg.nl)

# Opschaling pseudonimisatieservices

Robert Griffioen  
Projectleider de-identificatie project, SURF

# SURF coöperatie voor onderzoek en onderwijs

- Coöperatie van Nederlandse instellingen hoger onderwijs en onderzoek (>160), bijv. universiteiten, academische medische centra, hoge scholen en mbo's.
- Leverancier diensten of makelaar van diensten
- Innovatie partner en coördinator van dienstontwikkeling
- Kenniscentrum
- We worden 1 SURF per 1 juli 2020: verschillende doelgroepen en culturen
- SURFsara van oudsher gericht op onderzoek



# SURFsara: voorbeeld projecten

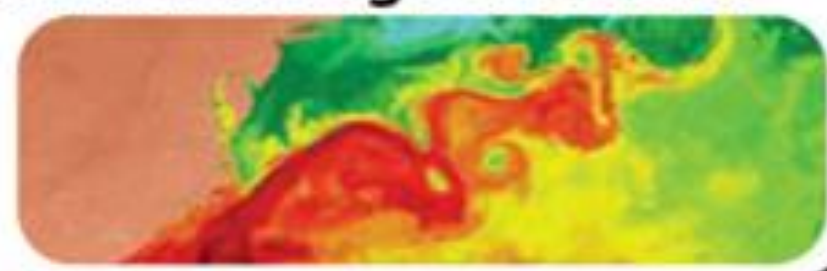
- 100 PFlop/s
- 100 gigabit/s
- 100 petabytes
- dedicated GPUs
- collaboration & access management

How to process 1 million news articles per day




Prof. Vossen  
VU

How to calculate the chance of increasing sea levels?



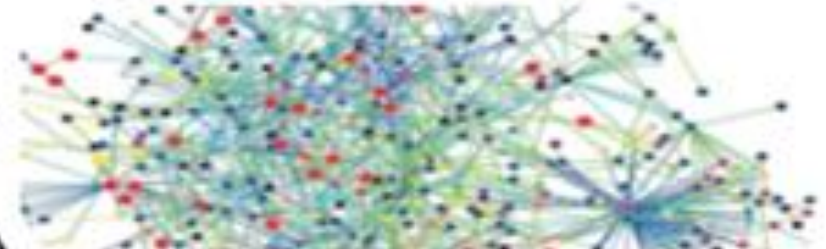
Prof. Dijkstra  
UU

Researching structures and properties of proteins



Prof. Bonvin  
UU

How to process and transfer proteomics data




Prof. Swertz  
UMCG/BBMRI

Researching the universe with thousands antennas



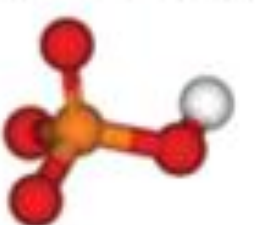
Dr. Holties  
LOFAR

Billions of tweets of research data



Prof. vd Bosch  
RUN

Computational catalysis for sustainable energy



Prof. Hensen  
TUE

# TTP-pseudonimisatie project

- Coördinerende SURF Contactpersonen – WO project
- Brede scope: koppelen van data centraal
- Analyse fase
- Meervoudige onderhandse aanbesteding: Hierbij kwamen project ZorgTTP-Groningen en dit project samen

- 1 Opslag of archivering van anonieme persoon-herleidbare data of persoonsgegevens
- 2 Pseudonimisatie van persoonsgegevens
- 3 Sleutelmanagement (sleutelgeneratie, -opslag, -beheer en –vernietiging)
- 4 Beheer van versleutelde data
- 5 Data koppeling op basis van versleutelde data
- 6 Supportloket ter ondersteuning van onderzoekers
- 7 Informed consent

# De-identificatie pilot

- Nadruk van de pilot ligt op de-identificeren, pseudonimiseren of anonimiseren, van persoonsdata.
- Nut/Noodzaak: allerlei lokale oplossingen binnen instellingen waarvan men op centraal bestuurlijk niveau niet op de hoogte is, laat staan over de kwaliteit. Vanuit governance oogpunt is het daarom wenselijk een uniforme wijze van de-identificeren te introduceren die voldoet aan de privacy en security richtlijnen van de verschillende instituten.
- Tool is vereiste, maar voor een complete oplossing zijn bijvoorbeeld ook organisatorische aspecten een eis.

# De pilot: werkwijze

- Deelnemers: Rijksuniversiteit Groningen, UMC Groningen, Universiteit Leiden en Vrije Universiteit
- Doel: ontwikkeling van pilot- tot Nationale SURF productiedienst/*SURF*-dienst. Verder uitwerken van reeds ontwikkelde diensten van 'Groningen' en ZorgTTP (zie presentatie Jan Lucas).
  - Use case 1: Kleinschalige studies (onomkeerbare pseudonimisatie)
  - Use case 2: Studies die periodiek worden bijgewerkt.
- Agile ontwikkeling: ontwikkeling en evaluatie tussenproducten
- Requirements en evaluatie door:
  - Onderzoekers
  - Data Steward
  - Functionaris Gegevensbescherming
  - IT

# De pilot: verder..

- Voordelen:
  - Werken met een ervaren en erkende partij op privacy gebied
  - Inrichten van organisatie samen met partner instellingen
- Deelname aan de productiedienst door instellingen
  - Reeds kenbaar maken van interesse om op de hoogte te worden gehouden:  
Simone van Wijngaarden ([simone.van.wijngaarden@zorgttp.nl](mailto:simone.van.wijngaarden@zorgttp.nl)) en  
Robert Griffioen ([robert.griffioen@surfsara.nl](mailto:robert.griffioen@surfsara.nl)).



# De toekomst: ontwikkeling nieuwe privacy diensten; van pilot tot dienst

- Ontwikkeling nieuwe privacy diensten:
  - Use case 3: Omkeerbare pseudonimisatie
  - Use case 4: Anonimiseren van data
  - Koppelen van data
  - Beelddata, etc.
- Binnen dit project, maar mogelijk ook: andere deelnemers, andere leverancier, ander project.

# Vragen?

# Dank!

**Jan Lucas van der Ploeg**

Data Engineer, Afdeling Informatiemanagement Onderzoek, UMCG

**Robert Griffioen**

Projectleider de-identificatie project, SURF

**Simone van Wijngaarden**

Adviseur, CIPP/e, Stichting ZorgTTP

[Simone.van.wijngaarden@zorgttp.nl](mailto:Simone.van.wijngaarden@zorgttp.nl)

06-13631500



**ZorgTTP**  
Privacy & vertrouwen