

Research data beveiligen en delen met behulp van Polymorfe Encryptie en Pseudonimisatie (PEP)

Jean Popma

Interfaculty Hub for Security, Privacy and Data Governance

J.Popma@ru.nl 29-10-2019



About myself

- Technical support manager 1987-2000 at Radboud University
- Managing director of ICT-service centre from 2000 to 2013 at Radboud University
- Corporate Information Security Officer and acting Privacy Officer from 2013-2016 at Radboud University
- Chairman of CERT-RU 2001-2017
- SURF CyberSaveYourself steering committee 2012-2016 <https://www.cybersaveyourself.nl/>
- Member of the board of the Privacy by design Foundation (working on IRMA, <https://privacybydesign.foundation/>)
- Currently Project Manager Applied Security Research working on PEP <https://pep.cs.ru.nl/> at iHub, Radboud University <https://www.ru.nl/iHub>



Biomedical research & big data ... What's the big deal?

Biomedical data are **sensitive** personal data.

- In case of a **data breach**:
 - harmful to participants/patients and
 - reputation damage for researchers
 - undermining patient's willingness to participate in future research is a **show-stopper** for medical research
- **Legal requirements** are strict
 - More so after GDPR has become fully operational (May 2018)
 - high fines / repercussions
- **Ethical restrictions** apply
- Professional **cooperation** between specialists from multiple disciplines (computer science, legal, ethics) and medical scientists is essential



Data Breach – As-A-Symptom

- | | |
|---|--|
| <ul style="list-style-type: none"> • Authorizations to broad • Orphaned data • Stolen or lost devices • Copies traveling around • Combination of different sources • Identity theft • Ransomware • Lack of maintenance • Hidden tracking/profiling • Changing context | <ul style="list-style-type: none"> • Too many data collected • Data reused for unintended purposes • Fraud • Revenge • Spionage • Wrongly addressed mail • Forgotten or lost print-output • Active hackers • Stupidity |
|---|--|



Privacy By Design?

- Data Design Strategies
 - Minimise
 - Separate
 - Aggregate
 - Hide
- Proces Design Strategies
 - Inform
 - Control
 - Enforce
 - Demonstrate

source: Jaap Henk Hoegman @ <https://arxiv.org/pdf/1210.6621.pdf>



Personalised Parkinson's Project



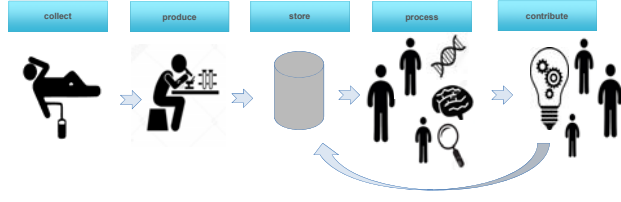

- 850 participants suffering from Parkinson's
- Followed over a period of 2 years
- Data will be available to scientists at dutch UMC's and other research institutes (worldwide)
 - Clinical data (questionnaires, tests)
 - Biophysiological data (from blood, plasma, CSF)
 - Wearable sensor data (2 years)
 - ECG
 - fMRI
 - Genome
 - Microbiome

Study Total: > 0.5 PB of data





Research Data Repository

2 main functions : share data, and ensuring scientific integrity

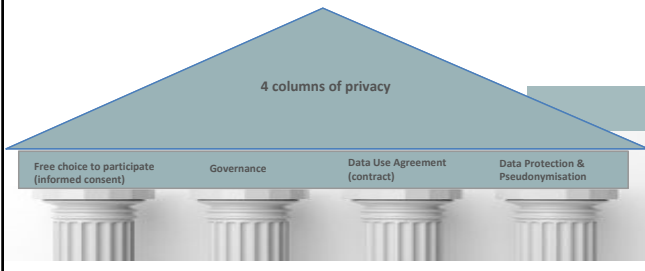

Challenges

- Protect the privacy and personal data of participants
- Enable the use of data collected for legitimate research
- If data are shared, receiving party obtains controllership over the data
- What about subject's rights?
- What about the obligations of the primary controller?




Building Trust


4 columns of privacy

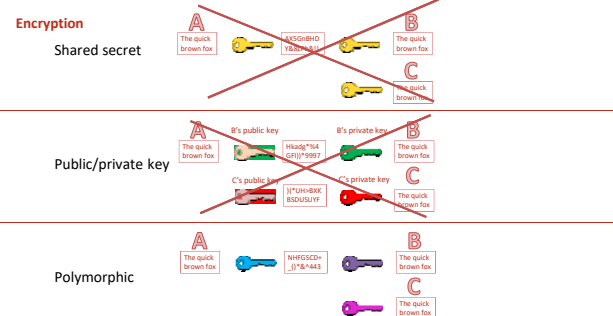

Data Protection & Pseudonymisation



- Research data should be encrypted and pseudonymised
- Encryption means that data cannot be read unless you have a proper key to unlock them.
- Pseudonymisation means that all directly identifying information is removed from the data and replaced by meaningless identifiers.
- Special system built as a data repository for these data: PEP = Polymorphic Encryption en Pseudonymisation (<https://pdp.cs.ru.nl>). PEP infrastructure takes care of the key- and data management



Encryption

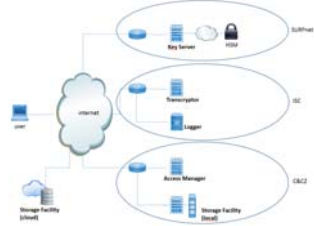
The key to encryption is the key










Distributed management of cryptographic keys






- User has client Software (encryption/decryption, up/download)
- 3 trusted parties take care of key management. Developers have no knowledge of keys.
- Keys in Hardware Security Modules (tamperproof)
- Storage Facility can be anything (public cloud)

Pseudonymisation




Data Collection

Subject	DOB	Gender	ECG	HeartRate	HRV	Stress	Stress 2
1000000001	1980-01-01	M	120	60	0.5	10	15
1000000002	1985-02-15	F	110	55	0.4	8	12
1000000003	1975-03-20	M	130	65	0.6	12	18
1000000004	1990-04-10	F	105	50	0.3	7	10
1000000005	1982-05-25	M	115	58	0.45	9	14
1000000006	1978-06-30	F	125	62	0.55	11	16
1000000007	1988-07-15	M	108	52	0.35	8	11
1000000008	1972-08-20	F	135	68	0.65	13	19
1000000009	1992-09-10	M	102	48	0.3	7	10
1000000010	1980-10-25	F	118	56	0.42	9	13

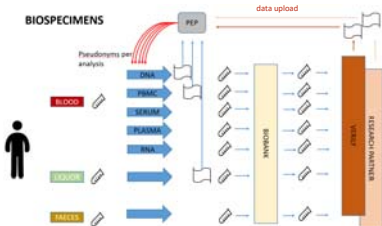




Polymorphic Encryption and Pseudonymisation (PEP) in a Data Repository




- Polymorphic encryption close to the data source
- Unique keys for different data and different users can be generated a posteriori, at the time access is granted.
- Decryption in the target processing environment
- All data streams are based on unique pseudonyms
- Unique and persistent pseudonyms are cryptographically generated for each user (user group) obtaining data from the repository

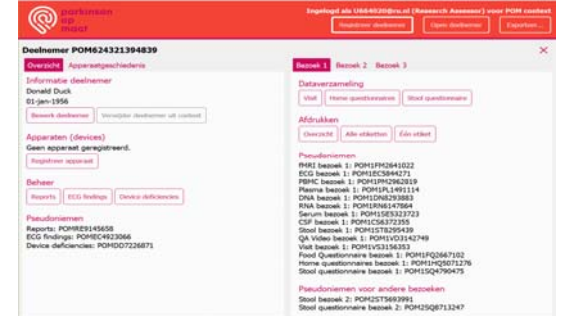





Pseudonymisation during data collection / analysis phase




- All datastreams based on unique pseudonyms
- When uploaded, these pseudonyms are translated to a polymorphic pseudonym, so data from different sources can be linked.
- Same idea for data from other sources (wearables, MRI, ECG etc.)













Pseudonymisation for data use by researchers




1. Research project **requests** data of certain types regarding subjects that meet certain criteria, based on a research proposal
2. Request is **evaluated** by a review board.
3. In case of positive **decision**:
 - Required columns for specific subjects are **authorized** to User
 - User **requests** data from PEP
4. User gets:
 - **Access** (download) to only these authorized data
 - **Pseudonyms** personalized to user(group).
 - A unique **key** is constructed to decrypt these (and only these) data
5. In the user's authorized processing environment (like DRE):
 - Data are **decrypted** for further processing and analysis
 - Derived data can be encrypted and **uploaded** back to the repository


	DNA	RNA	Biome	ECG	Wearable/IMS	Form1	Form2
Page 1							
Page 2							
Page 3							
Page 4							
Page 5							
Page 6							
Page 7							
Page 8							
Page 9							
Page 10							
Page 11							
Page 12							
Page 13							
Page 14							
Page 15							
Page 16							
Page 17							
Page 18							
Page 19							
Page 20							
Page 21							
Page 22							
Page 23							
Page 24							
Page 25							
Page 26							
Page 27							
Page 28							
Page 29							
Page 30							
Page 31							
Page 32							
Page 33							
Page 34							
Page 35							
Page 36							
Page 37							
Page 38							
Page 39							
Page 40							
Page 41							
Page 42							
Page 43							
Page 44							
Page 45							
Page 46							
Page 47							
Page 48							
Page 49							
Page 50							




Scientific Integrity




- Data can not be erased from the repository
- Data can become invalid at a certain point in time
- Participants can withdraw consent
- Shareable data may change over time
- Historic queries (for reproduction of earlier studies) must be possible



Results



- During the data collection phase samples cannot be linked to the same subject by different labs involved in production of the data.
- During transport (to and from repository) and storage data are protected by encryption at all times
- Key management is distributed over multiple parties: no single party can get access to the data (storage can be anywhere)
- Authorization is fine grained, and pseudonyms are personalized for individual researchers (research groups)
- Historical queries to facilitate replication studies are possible



Privacy By Design revisited




- Data Design Strategies
 - Minimise ✓
 - Separate ✓
 - Aggregate ✗
 - Hide ✓
- Process Design Strategies
 - Inform ✓
 - Control ✓
 - Enforce ✓
 - Demonstrate ✓



From: <https://www.oxfordjournals.org/doi/pdf/10.1093/oxfordjournals.issp.a011111>




Summary



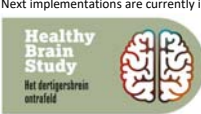
- Privacy risks are not reduced to 0 ...
- Data can be self-identifying (video/photo material, genome data etc.)
- Conspiring scientists can link data based on content.
- Enforcement and control of data use agreements is crucial but very difficult to do.
- Technology is not snake oil.
- Big Data, Cloud & Privacy can be combined if Privacy & Security are part of the design, at legal, organizational and technical levels **combined**.


Where do we go from here?



- First PEP pilot should demonstrate useability of the technology. It is a proof of concept (not for the Parkinson's project though)
- All software will be published as open source for reasons of transparency and public value.
- Next implementations are currently in progress





Healthy Brain Study
Het dertigjarig hersin onderzoek



Chronische Pijn
Langdurige pijn bij de volgende aandoeningen:
- Migraine
- Hoofdpijn
- Spierpijn
- Rugpijn
- Artritis
- Fibromyalgie
- Chronische pijn van onbekende oorsprong

- After pilot phases, a more general PEP system will be built and released
- Transfer of technology to [Privacy by Design Foundation](#) for further development and support





A CRYPTO NERD'S IMAGINATION:
HIS LAPTOP'S ENCRYPTED. LETS BUILD A MILLION-DOLLAR CLUSTER TO CRACK IT.
NO GOOD! IT'S 4096-BIT RSA!
BLAST! OUR EVIL PLAN IS FOILED!

WHAT WOULD ACTUALLY HAPPEN:
HIS LAPTOP'S ENCRYPTED. DRUG HIM AND HIT HIM WITH THIS \$5 WRENCH UNTIL HE TELLS US THE PASSWORD.
GOT IT.

More? <http://pep.cs.ru.nl>

