

Review

Data Protection Using Polymorphic Pseudonymisation in a Large-Scale Parkinson's Disease Study

Bernard E. van Gastel*, Bart Jacobs and Jean Popma

Interdisciplinary Hub for Security, Privacy and Data Governance, Radboud University, Nijmegen, the Netherlands

Accepted 11 May 2021

Pre-press 1 June 2021

Abstract. This paper describes an advanced form of pseudonymisation in a large cohort study on Parkinson's disease, called Personalized Parkinson Project (PPP). The study collects various forms of biomedical data of study participants, including data from wearable devices with multiple sensors. The participants are all from the Netherlands, but the data will be usable by research groups worldwide on the basis of a suitable data use agreement. The data are pseudonymised, as required by Europe's General Data Protection Regulation (GDPR). The form of pseudonymisation that is used in this Parkinson project is based on cryptographic techniques and is 'polymorphic': it gives each participating research group its own 'local' pseudonyms. Still, the system is globally consistent, in the sense that if one research group adds data to PPP under its own local pseudonyms, the data become available for other groups under their pseudonyms. The paper gives an overview how this works, without going into the cryptographic details.

Keywords: Computer security, privacy, big data, data protection

INTRODUCTION

A doctor is not supposed to treat patients as numbers. A medical researcher on the other hand should only see numbers (pseudonyms), not individuals. This is a big difference, which requires that the same person—a doctor also doing research—acts and thinks differently in different roles. It is a legal, data protection requirement to hide the identity of participants in scientific studies. Additionally, people are in general only willing to participate in medical research if their identity remains hidden. Hence, it is in the interest of researchers themselves to thoroughly protect the data and privacy of their study

participants, in order to provide sufficient comfort and trust to participate, now and in the future. With the increasing number of large-scale data gathering studies, high quality protections need to be in place, mandated by both legal requirements and ethical considerations.

We study how pseudonymisation can be applied to actual large scale data gathering projects to protect the data of participants. Pseudonymisation is one of many data protection techniques. The aim of pseudonymisation is to decrease the risk of reidentification and to decrease the risk of data being linked to data concerning the same participant without approval. This is difficult, for several reasons.

- Many other sources, outside the research dataset, such as social media, provide publicly available information that facilitates re-identification.

*Correspondence to: Bernard E. van Gastel, Radboud University (iHub), Erasmusplein 1, 6525 HT Nijmegen, the Netherlands. Tel.: +31 24 3652632; E-mail: b.vangastel@cs.ru.nl.

- 50 • Medical datasets are often very rich with many
51 identifying characteristics, for instance with
52 DNA or MRI data that match only one individ-
53 ual. Thus, the data themselves form persistent
54 pseudonyms.
- 55 • Continuous input from wearable devices pro-
56 vides identifying behavioral data or patterns.

57 The context of this paper is formed by a large
58 cohort study on Parkinson’s disease called Personal-
59 ized Parkinson’s project (PPP, see [2] for an elaborate
60 description of the study). This project aims to over-
61 come limitations of earlier cohort studies by creating
62 a large body of rich and high-quality data enabling
63 detailed phenotyping of patients. The PPP aims to
64 identify biomarkers that can assist in predicting dif-
65 ferences in prognosis and treatment response between
66 patients. It is clear that collecting such data is an
67 elaborate and costly undertaking. Maximizing the
68 scientific benefits of these data by sharing them for
69 scientific research worldwide is therefore important,
70 and one of the explicit goals of the PPP project.
71 Sharing sensitive biomedical and behavioural data is
72 a challenge in terms of legal, ethical and research-
73 technical constraints.

74 To enable responsible ways of data sharing, a novel
75 approach has been designed and implemented in the
76 form of a data repository for managing and shar-
77 ing of data for the PPP project. This design involves
78 so-called Polymorphic Encryption and Pseudonymi-
79 sation (PEP, see [8, 9]). The PEP system improves
80 over the current best practices in pseudonymisation
81 techniques as described in [1] by using a stronger
82 form of pseudonymisation based on asymmetric
83 encryption, in such a way that enables sharing of data
84 amongst different researcher groups, while not rely-
85 ing on a third party for pseudonymisation. Sharing
86 of data amongst different research groups requires
87 an easy but secure way to translate pseudonyms
88 as used by one research group to pseudonyms as
89 used by another research group. This implies that
90 pseudonymisation should be reversible, which is not
91 the case for some of the methods described in the
92 aforementioned overview of best practices. Using a
93 third party introduces a large amount of trust into
94 one single party. If that party acts in bad faith, data
95 protections can be circumvented. The infrastructure
96 hosting the PPP data repository is referred to as the
97 PEP-system, managed by the authors. Design and
98 development of the system was done in close coop-
99 eration with the PPP team. Of course this approach
100 is not restricted to the PPP study, but this study is

101 the first implementation and thus provides an exam-
102 ple for use in future studies, also outside the field of
103 Parkinson’s disease.

104 The working of our PEP system can be laid out after
105 describing the legal requirements to pseudonymisa-
106 tion and constraints for pseudonymisation. We will
107 briefly describe the PEP system, at a functional level,
108 with emphasis on polymorphic pseudonymisation—
109 and not on encryption, even though pseudonymi-
110 sation and encryption are tightly linked in PEP.
111 Pseudonymisation has become a scientific topic
112 in itself, but this article focuses on the practical
113 usage of the relatively new technique of polymor-
114 phic pseudonyms, whereby, roughly, each research
115 group participating in the study gets a separate set of
116 pseudonyms.

117 PSEUDONYMISATION AND THE GDPR

118 In health care, it is important to establish the iden-
119 tity of patients with certainty, for various reasons.
120 First, the right patient should get the right diagnosis
121 and treatment, but for instance also the right bill. In
122 addition, the patient’s identity is important for pri-
123 vacy/data protection, so that each patient gets online
124 access to only his/her own files and so that care
125 providers discuss personal medical details only with
126 the right individuals. In hospitals patients are fre-
127 quently asked what their date of birth is, not out of
128 personal interest, but only in order to prevent identity
129 confusion.

130 In contrast, in medical research identities need
131 not and should not be known, in principle. Certain
132 personal attributes, like gender, age, etc., are use-
133 ful in certain studies, but the identity itself is not
134 relevant. Indeed, if such attributes are used, they
135 should not be identifying. Treatment and research
136 are two completely different contexts [4], each with
137 their own legal framework and terminology. We
138 shall distinguish ‘patients’ and ‘health care profes-
139 sionals/providers’ in care, and ‘study participants’
140 and ‘researchers’ in research. In practice the same
141 person can be active both in health care and in
142 research and switch roles frequently. The awareness
143 of this context difference is therefore an important job
144 requirement.

145 In general, one can hide an identity via either
146 anonymisation or pseudonymisation. After apply-
147 ing pseudonymisation there is still a way back
148 to the original data set when combined with
149 additional information. With anonymisation this is
150

impossible. Europe’s General Data Protection Regulation (GDPR) uses the following descriptions.

- In Art. 4, pseudonymisation is defined as “the processing of personal data in such a way that the data can no longer be attributed to a specific data subject without the use of additional information, as long as such additional information is kept separately and subject to technical and organizational measures to ensure non-attribution to an identified or identifiable individual”.
- Anonymisation is described in Recital 26 as “... information which does not relate to an identified or identifiable natural person or to personal data rendered anonymous in such a manner that the data subject is not or no longer identifiable”.

The difference is legally highly relevant, since the GDPR does not apply to anonymous data. But it does apply to pseudonymous data.

In larger studies, such as PPP, participants are monitored during a longer period of time and are (re-)invited for multiple interviews and examinations. In such a situation a random number (a ‘pseudonym’) is assigned to each participant and this number is used in each contact with the participant during data collection. The GDPR thus applies to such studies.

Proper pseudonymisation is a topic in itself (see, e.g., [6]) since re-identification is always a danger [7]. Just gluing together the postal code and date of birth of a participant and using the result as pseudonym is not acceptable: it does not fit the above requirement: “... subject to technical and organizational measures to ensure non-attribution”. Pseudonymisation can in principle be done via tables of random numbers, often relying on trusted third parties to perform the translation. However, modern approaches use cryptographic techniques, like encryption and keyed hashing, see e.g., [5], which generate pseudonyms—the entries in such random tables. On the one hand such cryptographic techniques form a challenge, because they require special expertise, but on the other hand, they offer new possibilities, especially in combining pseudonymisation, access control and encryption. Researchers retrieving data from PEP need to authenticate (prove who they are) before they can get encrypted data that they are entitled to, together with matching local pseudonyms and decryption keys. There are commercial pseudonymisation service providers, but outsourcing pseudonymisation introduces additional dependencies and costs and makes such integrated pseudonymisation difficult.

This paper describes how so-called polymorphic pseudonymisation protects participants—and indirectly also researchers.

CONSTRAINTS

In practice, there are a number of constraints on how pseudonymisation can be applied. There are a number of sources these constraints originate from: biological data properties, from standard practices such as how to handle incidental findings, and how bio samples are handled. These constraints need to be taken into account into a system for data management. We can classify these constraints in the following types: re-identifying due to the nature of the data, re-identifying in combination with additional outside sources, regulation-based constraints, and practical constraints.

Even when pseudonyms are generated and used properly, digital biomedical data itself may lead to re-identification. This may happen in two ways.

1. Such biomedical data often contain patient identifiers. The reason is that devices for generating such biomedical data, like MRI-scanners, are also used for health care, in which it is important to ensure that data are linked to the right patient. Such devices may thus automatically embed patient identifiers in the data.
2. The biomedical data itself may be so rich that it allows for re-identification. This may be the case with MRI-scans of the brain which contain part of the face, or with DNA-data from which identifying characteristics can be deduced, or which can be compared to other DNA-databases

Hence, whatever the (cryptographic) pseudonymisation technique is, technical measures must be in place to remove such identifying characteristics from the data, especially of the first type.

An early experience in the PPP study illustrates the point raised above: through some error, a couple of MRI-scans got detached from their (internal) pseudonyms in the PEP-system. This is a disaster because it means that the scans have become useless. But the MRI-operator was not concerned at all and said: next year, when study participants return for their next visit (and MRI-scan), I can easily reconnect the few lost scans with matching new ones! Such matches do not involve the identities of the study participants but work simply because such an MRI scan is uniquely identifying any participant.

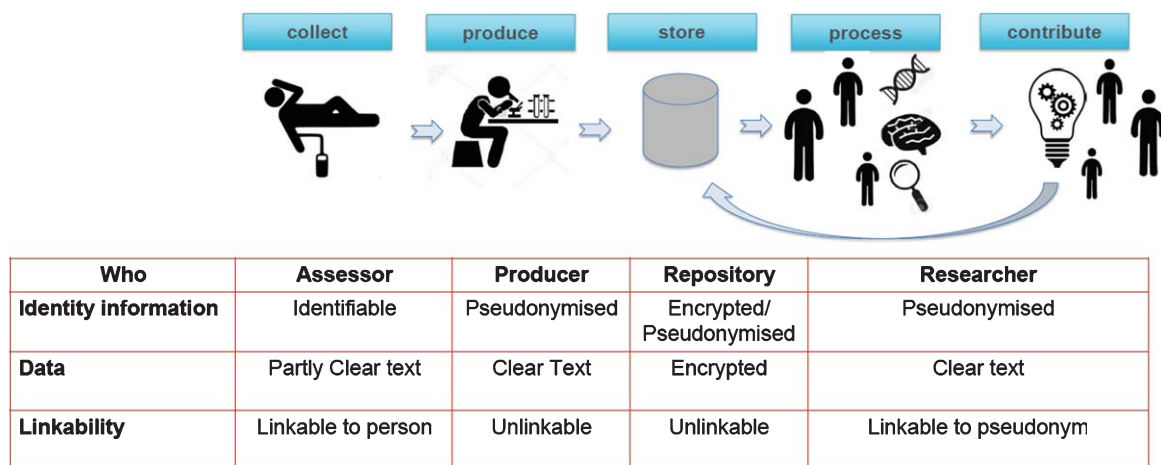


Fig. 1. Phases in research data management.

249 Combining the data with outside sources could also
 250 lead to re-identification. Besides the data itself, meta-
 251 data such as timestamps can be identifying, especially
 252 if combined with external additional data sources.
 253 Not storing metadata such as time stamps is not
 254 always enough: if there is frequent data synchron-
 255 ization, it is easy to determine when certain data
 256 was first made available. This is with high proba-
 257 bility the day a participant came in for tests. The
 258 effect of combining external information sources
 259 is not always clear for everybody involved, as can
 260 be seen in another experience of our team. Find-
 261 ing study participants is always a challenge, and so
 262 the staff of the PPP study suggested to people who
 263 had already signed up that they could enthusiastically
 264 report their PPP-involvement on social media, in
 265 order to further increase participation. We were upset
 266 when we heard about this and had to explain how such
 267 disclosures on social media undermine our careful
 268 pseudonymisation efforts. Data protection, including
 269 pseudonymisation, is not solely a technical matter.
 270 It can only work with a substantial number of study
 271 participants whose participation is not disclosed, at
 272 least not in a way that allows linking to actual data in
 273 the system, like for instance date and time of visits.
 274 The smaller the population of participants, the eas-
 275 ier it is to single them out based on just very little
 276 information.

277 There are also practical constraints. The pseudo-
 278 nyms used need to be human-readable and short
 279 enough to fit on labels. However, the used polymor-
 280 phic pseudonyms are based on non-trivial cryptogra-
 281 phic properties, namely the malleability of El Gamal
 282 encryption. The handling of these pseudonyms within

283 the PEP-system happens automatically and is not a
 284 burden for the researchers that interact with the sys-
 285 tem for storing and retrieving study data. Internally,
 286 these pseudonyms are very large numbers (65 charac-
 287 ters long, see below). In fact, they are so large that they
 288 are not suitable for external usage by humans. For
 289 instance, these internal pseudonyms do not fit on regu-
 290 lar labels on test tubes. As a result, shorter external
 291 representations of these polymorphic pseudonyms
 292 are used for input—and also output—of the PEP-
 293 system.

294 There are regulation-based exceptions. In scienti-
 295 fic research on medical data, researchers may come
 296 across an ‘incidental finding’ about the health con-
 297 dition of the study participant. Under such special
 298 circumstances it may be needed to de-pseudonymise
 299 deliberately, and to turn the study participant into a
 300 patient. In some studies, like in PPP, separate, excep-
 301 tional procedures must be in place for such cases.

302 PEP IN FIVE PHASES

303 The five phases of research data management via
 304 the PEP-system are summarized in Fig. 1, with spe-
 305 cial emphasis on the identity of the study participant.

306 These phases will be discussed separately below.
 307 They suggest a certain temporal order, but in practice
 308 these phases exist in parallel, for instance because
 309 study participants are followed during a longer time
 310 (two years, in principle) and are monitored (1) during
 311 repeated visits at discrete intervals, and (2) also con-
 312 tinuously, via a special sensor based wearable device
 313 in the form of watches, provided by Verily.

Plasma visit 1 POM1PL0610058

fMRI visit 1 POM1FM2641022

ECG visit 1 POM1EC5844271

PBMC visit 1 POM1PM2962819

DNA visit 1 POM1DN8293883

RNA visit 1 POM1RN6147864



Fig. 2. Examples of short pseudonyms, also printed on a test tube. Such a short pseudonym like POM1PL0610058 consists of five random numbers (here 10058), a prefix (POMPL) identifying the project (POM, Dutch for PPP) and type of data (1PL for plasma taken during the first visit), and a checksum (58) such as used in an International Bank Account Number (IBAN).

314 The PEP-system makes crucial use of so-called
 315 polymorphic pseudonyms for study participants.
 316 These are large numbers that typically look as fol-
 317 lows—65 characters in hexadecimal form: 0EAD7
 318 CB2D70D85FE1295817FA188D22433C2237D946
 319 964A6B4C063E6274C7D7D.

320 Such numbers are not meant for human consump-
 321 tion. They have an internal cryptographic structure
 322 so that they can be transformed to another such
 323 number, called a local pseudonym, to be used by a
 324 particular researcher (or group of researchers) of the
 325 PEP-system.

326 This transformation of pseudonyms is done by
 327 several, independently managed components of the
 328 PEP-system, working together to produce the final
 329 pseudonymisation result. This is done “blindly”,
 330 which means that the components of the system per-
 331 form their tasks without knowing the identity of the
 332 study participant involved¹: internally the system is
 333 able to produce a persistent local pseudonym for each
 334 research-group.

335 In this way each international research group that
 336 joins the PPP research effort gets its own local
 337 pseudonyms for all the study participants. If for some
 338 reason data get compromised or re-identified in such
 339 a research group, the effect is local, in principle,
 340 and does not affect the whole system. There are fur-
 341 ther organizational and legal safeguards, implying for
 342 instance that participating research groups get access
 343 to only the data that is relevant for their research ques-
 344 tions, and that the data may only be used for specific

and well-defined purposes, but that is a different
 topic.

We now discuss the five phases in Fig. 1. The col-
 lection phase consists of two parts, with repeated
 site-visits and with continuous monitoring via sensor-
 based wearable devices. To the participant they have
 the appearance and function of a watch, hence the
 term ‘study watch’ is used in practice. These site-
 visits typically take a whole day and involve several
 medical examinations, tests, and interviews. During
 such a visit, a dedicated assessor accompanies each
 study participant. During the first visit the study par-
 ticipant receives a study watch, provided by Verily
 Life Sciences, which is permanently active—except
 during daily charging and transmission of collected
 data. The watch contains a serial number, which is
 linked to a local polymorphic pseudonym associated
 with Verily. Verily receives the combination <serial-
 number, local-pseudonym>, but does not learn which
 study participant gets which watch. Verily collects
 the data from the watches into its own system, via the
 serial number, and uploads the sanitized data into the
 PEP-system via the associated local pseudonym.

For each visit of a study participant the assessor
 gets an automatically prepared table of short external
 pseudonyms, connected to the internal pseudonym
 that is associated with the study participant. At each
 lab for biospecimens or measurement (say for MRI),
 the associated short pseudonym is put on the relevant
 tube or file, see Fig. 2.

The production phase is for cleaning-up and for
 producing data in such a format that it can be
 uploaded into the PEP-system. This may involve
 sanitization like for the study watch data, transform-
 ing raw data into data that can be used in further

¹ As an aside: since these components also perform logging, it is cryptographically ensured that nothing can happen in the PEP-system without producing an audit log.

380 analysis, or performing measurements on biospecimens.
381 The measurement devices in the PPP study
382 are typically also used in a health care setting, for
383 patient diagnosis. This means that the output files
384 often contain identifiers that become persistent if
385 not removed before uploading data to the repository.
386 They have to be removed from the data, before
387 upload to the PEP-system, to prevent unwanted re
388 identification. Similarly, MRI-data must be de-faced,
389 so that the study participants cannot be recognised
390 visually or by face-recognition algorithms. During
391 the storage phase, the sanitised, appropriately formatted
392 data is encrypted, uploaded and stored. This
393 encryption is also done in ‘polymorphic’ manner but
394 how that works is out of scope here (see [3] for
395 some more information, and [9] for cryptographic
396 details). Actual storage happens in a data centre of
397 Google in Europe. Since all stored data are encrypted
398 and encryption keys are handled internally in the
399 PEP-system, the actual storage provider is not so
400 relevant, since the data is protected from the storage
401 service provider. What matters most is that the
402 encrypted data is available only to legitimate users of
403 the PEP-system, when needed, and is not corrupted.
404 The processing phases exists for research-groups that
405 have been admitted to the PPP project, after submitting
406 and approval of a research plan, and after signing
407 a data use agreement (see [3]). Such a user-group then
408 gets its own local pseudonyms and a decryption key
409 for locally decrypting a download from the PEP-data
410 repository containing the research dataset it is entitled
411 to. Typically, this dataset is necessary and minimised
412 for the approved research plan. Such a research-group
413 thus has access to unencrypted (clear text) medical
414 research data, but only in a locally pseudonymised
415 form. Data may only be processed in a secured
416 processing environment. There is an additional
417 contribution phase in which such a research-group can
418 return certain derived results to the PEP-systems. This
419 group uses its own local pseudonyms for the upload.
420 Once uploaded, these additions become available for
421 all other participating research-groups, under their
422 own local pseudonyms. The PEP-system ensures
423 consistency of these pseudonyms, via its blind translation
424 mechanism.

425 The table in Fig. 1 provides an overview of the
426 identity and encryption status of each of the phases
427 and of the names/roles involved. Only during the
428 collection phase, the identity of the study participant
429 is (necessarily) known, for instance to the assessors,
430 for personal contact and logistic purposes. Personal,
431 identifying data like name and contact information of

study participants are stored within the PEP-system,
but these data are never handed out, except via a very
special procedure, involving a designated supervisor,
for emergency reporting of incidental findings.

FINAL REMARKS

We have described the main lines of the PEP-system that is designed and built for privacy-friendly management of research data, focused on medical studies.

The PEP methodology combines advanced encryption with distributed pseudonymisation, and distribution of trusted data with fine-grained access management, allowing access to be restricted to subsets of participants and subsets of different data types. It thus allows cooperation of public and private research organizations on a global level. PEP provides secure access to minimized datasets with local pseudonymisation, in a global setting, including (selected) contributions and research results via those local pseudonyms. Although this article concentrates on pseudonymisation, within the PEP-system pseudonymisation is tightly integrated with access control, audit-logging and encryption of the research data. PEP is currently being used in the large-scale PPP study, which will encompass approximately one terabyte of data per participant, for a total of 650 participants. This final section contains some additional remarks on specific points and comes to a conclusion.

Under the GDPR, every ‘data subject’ has a right of access, that is, a right to see which data a particular organisation holds (on oneself). This is particularly challenging for a research setting with pseudonymisation. The PPP study supports this right of access by telling study participants that they can come and visit and then get access to their data in the repository, via the special data administrator (supervisor). But the PPP does not give participants online access to their data; that would require dedicated client-side software together with a reliable form of authentication that is coupled to polymorphic pseudonyms as used in the PEP system. Such a coupling is a challenge and does not exist at this stage. Another practical reason is that the PEP-system contains mostly raw biomedical data which can only be interpreted by specialists.

Within the GDPR subjects can always withdraw their consent and ask that data on them is removed. This right clashes with the obligation that scientific research must be reproducible. For reasons of archival and scientific research, this is not possible however.

481 In the PPP project, a compromise is implemented in
482 the PEP system, that after withdrawing consent the
483 existing data is not removed but is not used any more
484 for new studies. In this way, the data can still be used
485 for reproducing results of already published articles
486 (e.g., in case of doubts or even fraud allegations).

487 As a final remark, a dedicated software team of
488 4-5 people has been developing the PEP-system for
489 roughly four years (2017-2020). It is now in stable
490 form and other research teams are starting to use
491 PEP as well. The software has become open source in
492 late 2020 in order to provide maximal transparency
493 about how it works². No special hardware or licences
494 are required to run PEP. However, running PEP does
495 require some guidance, which the PEP-team is plan-
496 ning to provide for the coming years.

497 CONFLICT OF INTEREST

498 The authors have no conflict of interest to report.

499 REFERENCES

- 500 [1] Agrafiotis I, Bourka A, Drogkaris P (2019) *Pseudonymisa-*
501 *tion techniques and best practices*. European Union Agency
502 for Cybersecurity, Greece. Available at: [https://www.enisa.](https://www.enisa.europa.eu/publications/pseudonymisation-techniques-and-best-practices)
503 [europa.eu/publications/pseudonymisation-techniques-and-](https://www.enisa.europa.eu/publications/pseudonymisation-techniques-and-best-practices)
504 [best-practices](https://www.enisa.europa.eu/publications/pseudonymisation-techniques-and-best-practices).
- [2] Bloem BR, Marks WJ Jr, Silva de Lima AL, Kuijf ML, 505
van Laar T, Jacobs BPF, Verbeek MM, Helmich RC, van 506
de Warrenburg BP, Evers LJW, intHout J, van de Zande T, 507
Snyder TM, Kapur R, Meinders MJ (2019) The personalized 508
Parkinson project: Examining disease progression through 509
broad biomarkers in early Parkinson's disease. *BMC Neurol* 510
19, 160. 511
- [3] Jacobs B, Popma J (2019) Medical research, big data and the 512
need for privacy by design. *Big Data Soc* **6**, 1-5. 513
- [4] Nissenbaum H (2009) *Privacy in context. technology, policy,* 514
and the integrity of social life. Stanford University Press, 515
Redwood City. 516
- [5] European Commission: Article 29 Data Protection Working 517
Party (2014) *Opinion 05/2014 on anonymisation techniques*. 518
Available at: [https://ec.europa.eu/justice/article-29/documen-](https://ec.europa.eu/justice/article-29/documentation/opinion-recommendation/files/2014/wp216_en.pdf) 519
[tation/opinion-recommendation/files/2014/wp216_en.pdf](https://ec.europa.eu/justice/article-29/documentation/opinion-recommendation/files/2014/wp216_en.pdf). 520
- [6] Pfitzmann A, Hansen M (2008) Anonymity, unlinkability, 521
undetectability, unobservability, pseudonymity, and identity 522
management — a consolidated proposal for terminology. 523
Available at: [https://dud.inf.tu-dresden.de/literatur/Anon-](https://dud.inf.tu-dresden.de/literatur/Anon-Terminology_v0.31.pdf) 524
[Terminology_v0.31.pdf](https://dud.inf.tu-dresden.de/literatur/Anon-Terminology_v0.31.pdf). Accessed on May 13, 2021. 525
- [7] Sweeney L (2000) *Simple demographics often identify people* 526
uniquely. Carnegie Mellon University, Pittsburgh. Avail- 527
able at: [http://dataprivacylab.org/projects/identifiability/pa-](http://dataprivacylab.org/projects/identifiability/paper1.pdf) 528
[per1.pdf](http://dataprivacylab.org/projects/identifiability/paper1.pdf). 529
- [8] Verheul E (2015) Privacy protection in electronic education 530
based on polymorphic pseudonymization. In *IACR Cryptol-* 531
ogy ePrint Archive 2015/1228. Available at: [https://eprint.](https://eprint.iacr.org/2015/1228) 532
[iacr.org/2015/1228](https://eprint.iacr.org/2015/1228).2015. 533
- [9] Verheul E, Jacobs B (2017) Polymorphic encryption and 534
pseudonymisation in identity management and medical 535
research. *Nieuw Archief Wiskunde* **5/18**, 168-172. 536

² For details, see <https://pep.cs.ru.nl>