



# AnoMed

Kompetenzcluster Anonymisierung für  
medizinische Anwendungen

## Deliverable 4.9.8

# Strategies to manage the risk of re-identification

©ULD 2025



Funded by  Federal Ministry  
of Research, Technology  
and Space



Funded by  
the European Union  
NextGenerationEU



<https://anomed.de>

UAP	4.9.8
Date	2025
Version	1.0
Status	Final
Distribution (only final version)	PU
Lead Contributors (© by affiliation)	Bud P. Bruegger (ULD)
Additional Contributors (© by affiliation)	
Reviewers	Harald Zwingelberg (ULD)
License	CC-BY 4.0

# Table of Contents

- 1 Introduction..... 5
- 2 Understanding Anonymity ..... 5
- 3 Types of Identity-Reduced Data ..... 6
- 4 Re-identification risk over time ..... 10
- 5 Legal consequences of successful re-identification of identity-reduced data ..... 11
  - 5.1 Possible sanctions..... 11
  - 5.2 Notification of re-identification event..... 12
  - 5.3 Actions to reach compliance after a re-identification event ..... 13
  - 5.4 Redressing the past violation ..... 13
- 6 Monitoring of the re-identification risk and responsible disclosure ..... 14
  - 6.1 The information security ecosystem as an analogy ..... 14
  - 6.2 Proposal of an ecosystem to monitor the re-identification risk ..... 15
- 7 Role of the proposed monitoring ecosystem..... 17
- 8 Incubation of a monitoring ecosystem..... 19
  - 8.1 Establishment of the authority and its infrastructure..... 19
  - 8.2 Obligations for anonymizers, Dataset Users, and re-identification researchers ..... 19
  - 8.3 Incentives for re-identification and anonymity researchers..... 20
- 9 Legal requirements for attempted re-identification of anonymous data..... 22
- 10 Conclusions..... 23

# 1 Introduction

In the large-scale sharing of data in data spaces, anonymization techniques assume an important role since they promise to significantly facilitate free sharing compared to the sharing of personal data. The present report discusses strategies to responsibly manage “anonymized” data in data spaces.

Technical transformations called “anonymization techniques” render it more difficult and unlikely to identify data subjects after the transformation; they fail to guarantee that this is impossible, however. Hence, even careful “identity-reduction” always remains subject to a residual risk of re-identification.

No objective method exists to assess anonymity—no analytical objective assessment method is known. That re-identification is possible can be determined in an empirical approach of trying different re-identification methods, however. This report therefore proposes the use of an ecosystem to manage residual re-identification risk. In the ecosystem, empirical assessment of re-identification resistance by “re-identification” researchers is encouraged and managed by an equivalent body to a CERT/CSIRT known from information security. The main strategy then to foster ethical actors to detect re-identification vulnerabilities without endangering data subjects, before malicious actors can exploit such vulnerabilities.

## 2 Understanding Anonymity

To reason about anonymity, the role of anonymization, and the risk of re-identification, it is necessary to understand the role of these concepts in the bigger context of the fundamental right to data protection. This section describes in this context and attempts to create the necessary understanding.

Data protection is a fundamental right of individuals (Art. 8 ECFR). Processing of personal data is considered to restrict this right. Therefore, the impact of such processing has to be minimized to what is necessary for legitimate purposes. In addition, processing can create risks for the rights and freedoms (i.e., not limited to the fundamental right to data protection) of natural persons (predominantly, data subjects). The GDPR takes a **risk-based approach** that mandates controllers to (i) identify such risks and (ii) implement an adequate level of **technical and organizational measures** to **protect** data subjects against excessive risks.

The GDPR uses a simple **indicator** to determine whether risk identification and mitigation is necessary: namely it is in the case the data is considered **personal data** and it is not in the case that the data is not personal data. Data that is not personal is called **anonymous**.

In the following, the term **protection measure** is used to denote any possible measure a controller takes to protect the rights and freedoms of data subjects (and possible other natural persons). It is up to a legal discussion whether that coincided exactly with the *term technical and organizational measure (TOM)* used in the GDPR. Protective measures are therefore either a superset or equal to TOMs. A different term is used here to steer free of a possible distracting discussion whether all measures in the here presented reasoning are actually TOMs.

There are two super protective measures, or simply **super measures**, that, if effective, guarantee the total elimination of all risks. Namely, they are **data erasure** (deletion) and **anonymization**. The effective use of either of these measures means that the obligations of the GDPR, namely to identify,

minimize, and mitigate risks are no longer applicable. This consequence follows from the fact that these super measures completely eliminate any risk and thus render identification, minimization, and mitigation mute.

Evidently, the super measures have this effect of taking processing outside the GDPR only, if they are indeed effective. A good illustrative example are PDF documents, that were “redacted” by using overlaid black bars. When displayed or printed, this indeed has the effect of erasing the covered data; when looking at the file in electronic form, this erasure was instead not effective.

Consequently, if the objective was to erase personal information, the measure was not effective, the risk was not eliminated, and the reason to fall outside the GDPR becomes mute. The risk inherent in the data is obviously still present. Obligations of the GDPR, such as the mandate to keep the data confidential, thus still apply. The fact that the creators of the “redacted” document were unaware of the ineffectiveness of the measure can possibly protect them from the accusation of negligence, but cannot eliminate legal obligations.

### 3 Types of Identity-Reduced Data

After discussing the role of anonymity in a wider context, this section defines concepts that are necessary for the discussion of anonymity. It defines several concepts and discusses them in some detail.

The first term that requires a precise definition is *anonymization*. In the context of this report, it is defined as follows:

**(1) Anonymization**

*Anonymization* is a transformation that takes personal data as input and yields non-personal (i.e., truly anonymous) data as output.

Note that in common practice, the term *anonymization* is used more loosely and fails to require that the output of the transformation is indeed anonymous. Consequently, scientists have shown the possibility of re-identifying (or “de-anonymizing”) “anonymized” data.

Taking the stricter point of view that anonymous is an indicator of an absence of risk and obligations, the concept of *anonymization* used here is only applicable if it indeed results in anonymous data. The notion of *re-identifying* anonymized data is therefore an impossibility. Similarly, the term “*de-anonymize*” is a contradiction in itself.

Currently, no data transformation is known which guarantees with certainty that its output is truly anonymous<sup>1</sup>. It is even likely that true *anonymization* does not even exist. This report therefore needs a concept for the transformations that are used in practice.

**(2) Identity-reduction transformation**

An *Identity-reduction transformation* is a transformation that takes personal data as input and yields *identity-reduced data* as output.

**(3) Identity-reduced data**

*Identity-reduced data* is information about natural person where the possibility of re-identification is rendered more difficult (compared to the input data of identity-reduction).

---

<sup>1</sup> Note that even differential privacy methods provide guaranteed protection against re-identification only as long as the privacy budget is not eroded by other disclosures. Much rather, the scientific community believes that any disclosure leaks some information suitable to support re-identification.

Identity-reduction thus just impedes re-identification but remains without guarantees to render re-identification impossible. With increasing identity-reduction, the difficulty and cost of re-identification increases.

A given context may specify the means available to potential re-identification. When identity-reduction aims at a sufficient protection against re-identification in such a context, the following concept is useful to describe this situation.

**(4) Completely identity-reduced data** (relative to a given context)  
 Data is called completely *identity-reduced* in a given context, if the degree of identity-reduction is sufficient to protect against any re-identification attempt that is possible with the means available in the context.

In a given situation, it may be subject to discussion what context to apply to an assessment of identity-reduced data, i.e., which means of re-identification have to be considered. In addition, demonstrating that an identity-reduction is indeed complete for an agreed on context may be difficult. This said, the proposed terminology makes a clear distinction between disagreements on the applicable context and disagreement on completeness.

In preparation of further discussing the properties of identity-reduced data, the following concept is defined:

**(5) Undecidable absence**  
*In a context, where the presence of a certain entity can be detected only with a certain likelihood and where the number of such entities is unlimited, the absence of such entities cannot be shown.*

In particular, in a situation of undecidable absence, it is impossible to distinguish the case where at least one entity was not detected from the case where no such entities are present.

A good practical illustration of this concept stems from information security. Here, “secure” is equivalent to the absence of vulnerabilities. Vulnerabilities are hard to detect, however. Therefore, a system without known vulnerabilities cannot be considered invulnerable since it may be affected by yet unknown vulnerabilities. (Pen) testing the security of a system can only check for known vulnerabilities, but cannot guarantee that a system is secure.

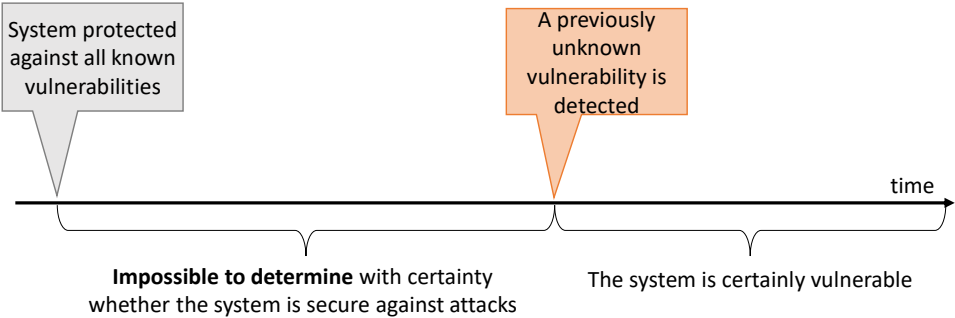


Figure 1: Undecidable absence of vulnerabilities in information security.

Figure 1 further illustrates the undecidability in information security. In a first event, a system is updated such that it is secure against all known vulnerabilities. In a second event, a previously

unknown vulnerability is discovered. Before this event, it is impossible to determine with certainty whether the system is free of vulnerabilities (although experience shows that this is highly unlikely). No procedure exists to show the absence of vulnerabilities. In contrast, the presence of vulnerabilities is detectable and can be tested for.

In consequence, information security is seen as a process much rather than a state. In particular, the process is concerned with tracking and mitigating newly discovered previously unknown, vulnerabilities.

Very similar to information security, also in the field of identity-reduction, an undecidable absence problem exists. In particular, identity-reduced data are subject to vulnerabilities relative to certain re-identification attacks. Even if an identity-reduction defends against all known attacks, it remains uncertain if a yet unknown, more sophisticated attack (even with the same means) exists that can successfully re-identify the data.

This situation is illustrated in Figure 2. Even if the means available to re-identification remain unchanged, more sophisticated methodology may unexpectedly render re-identification possible. Consider for example re-identification through linkage of individual-level data. At one point in time, the addition of noise could possibly have prevented linkage based on equality or closeness. More sophisticated noise-resistant, pattern- or structure-based linkage methods change this. Even further, AI-based linkage methods may be able to re-identify data that was previously thought to be safe.

Similarly, K-anonymity<sup>2</sup> was considered the state-of-the-art protection against re-identification. Later, it was found to be insufficient if L-diversity<sup>3</sup> wasn't considered. Even later, T-closeness<sup>4</sup> and then downcoding attacks<sup>5</sup> were added.

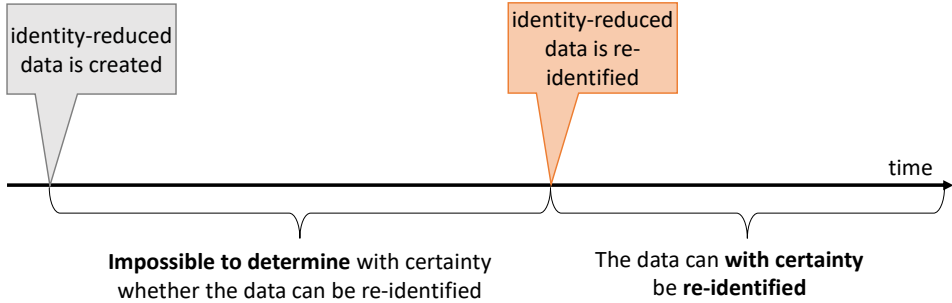


Figure 2: Undecidable absence of re-identification vulnerabilities of identity-reduced data.

The above analysis shows the impossibility to determine with certainty, whether data is completely identity-reduced for a given context, let alone whether it is anonymous in the legal sense and poses no risk at all to the rights and freedoms of natural persons and therefore falls outside the GDPR.

In absence of an objective method to determine the status of identity-reduced data, a strategy to manage this uncertainty is needed. For this purpose, two additional concepts are defined. They are

<sup>2</sup> <https://en.wikipedia.org/wiki/K-anonymity>.  
<sup>3</sup> <https://en.wikipedia.org/wiki/L-diversity>.  
<sup>4</sup> <https://en.wikipedia.org/wiki/T-closeness>.  
<sup>5</sup> Cohen, Aloni. (2022). Attacks on Deidentification's Defenses. 10.48550/arXiv.2202.13470., <https://www.usenix.org/conference/usenixsecurity22/presentation/cohen>.

illustrated in Figure 3. The upper part shows the objective, factual assessment of the identity-reduced data. Namely, it is impossible to determine whether re-identification is actually possible or not. Therefore, they are subject to a residual risk of re-identification.

The objective classification fails to provide an answer to the question whether the requirements of the GDPR apply. For example, can the data be published (as is possible for anonymous data) or is confidentiality required (as stated by the GDPR for personal data).

The solution lies in a subjective classification by the party (called *assessor* in the figure) who assesses the identity-reduced data. In data protection, controllers are required to act as assessors. Other parties such as supervisory authorities or courts could also act in this role.

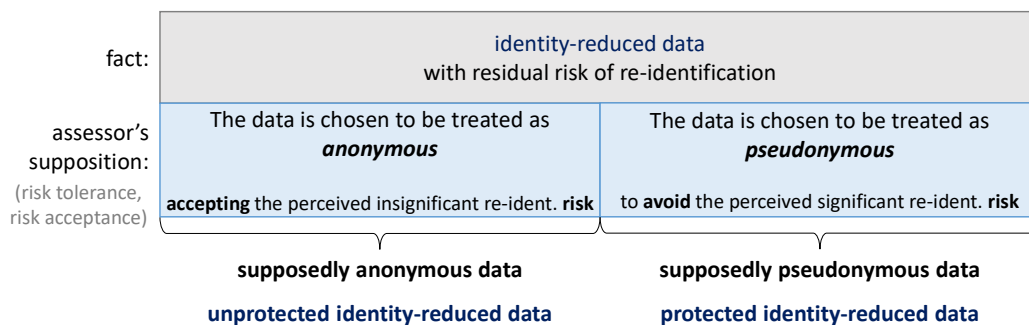


Figure 3: Possible classifications of identity-reduced data.

The assessor has to make a subjective decision of how to judge and thus handle the residual risk or re-identification.

**One option** is to consider this **risk** to be **insignificant** and consequently **accept the risk**. The assessor then treats the data as anonymous and avoids fulfilling the requirements of the GDPR for personal data. If the data unexpectedly proves to be personal after all (i.e., it can be re-identified), the assessor accepts the consequences of the erroneous decision (see Section 5). In absence of protections mandated by the GDPR for personal data, data subjects who did not accept the risk may also be exposed to negative consequences. Some of the effects of such a breach may be irreversible and unamendable.

**An alternative option** is to consider the **risk** to still be **significant** and consequently **avoid taking the risk** (i.e., play it safe). The assessor then treats the data as it was pseudonymous (i.e., personal). This means that in the case of a re-identification of the data, there are no consequences for the assessor nor for data subjects since all protections mandated by the GDPR are in place.

This analysis leads to the following two definitions:

**(6) *Supposedly anonymous data – unprotected identity-reduced data***

*Supposedly anonymous data or synonymously unprotected identity-reduced data is identity reduced data that is legally treated as anonymous data by a party who accepts the residual risk of re-identification and accepts the possible consequences.*

**(7) *Supposedly pseudonymous data – protected identity-reduced data***

*Supposedly pseudonymous data or synonymously protected identity-reduced data is identity reduced data that is legally treated as pseudonymous data by a party who protects itself from the residual risk of re-identification.*

Since there is not an objective method to relate the technical concept of identity-reduced with the legal classification of personal data vs. anonymous, these concepts inherently remain somewhat disconnected. It is important to understand that any attempt to relate these concepts is a supposition by some assessor and is necessarily subjective. Consequently, it is unavoidable that different assessors make different suppositions.

The concept of “completely identity-reduced” fails to solve this problem. By no means can the absence of known vulnerabilities guarantee that re-identification is impossible such that there is no risk for the rights and freedoms of the affected persons.

## 4 Re-identification risk over time

This section describes the temporal aspect in the legal definition of anonymity and the technical protection against re-identification over time.

Recital 26 GDPR provides guidance on how to determine whether data is personal. Among others, it states that “[...] account should be taken of all objective factors, such as [...], taking into consideration the available technology **at the time of the processing** and **technological developments**.” (Highlighting added by the author).

Clearly, this states that beyond the present situation, also **future developments** have to be taken into consideration. The Recital refers only to “technological” developments; it does not speak of other kinds of developments. Leaving a precise legal interpretation to legal scholars, the following assumption is made for the discussion: “technological developments” are contained in the sentence part that provides examples for “objective factors”. With this assumption, other kinds of developments besides just technological ones also have to be taken into account<sup>6</sup>. With this disclaimer, the following discussion looks at relevant developments from a technical perspective.

A future development is relevant to the discussion if it affects the risk of re-identification. Therefore, the following developments seem to be of interest:

- Increase of **computing power** (i.e., technological development),
- improvements in the **methodology** used for re-identification,
- increase in **additional information** suitable for identification through linkage,
- and increase of **additional disclosures** suitable for reconstruction of individual-level information (i.e., disclosures that reduce the “privacy budget”).

These kinds of future developments are illustrated in Figure 4. It shows the reasonable assumption of an increase of all types over time. At the bottom right, it also shows a summery where all these types are combined into “means reasonably likely to be used” (i.e., the wording from Recital 26 GDPR). Conceptually it seems to be sufficient to consider that the risk of re-identification increases with time. The exact reasons for this assumption may be irrelevant for the decision how to best manage this fact.

---

<sup>6</sup> These have to be taken into account at least if they are considered “objective factors”.

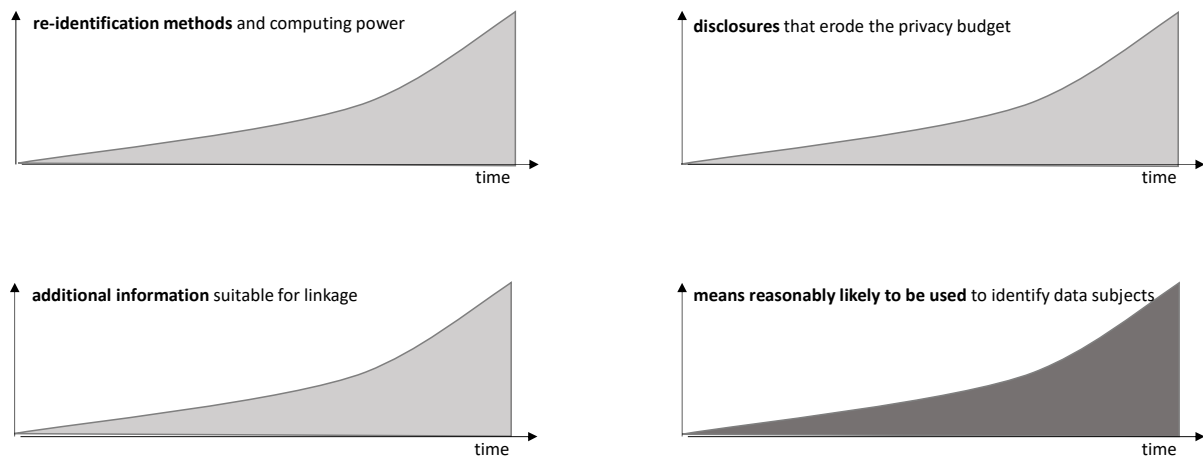


Figure 4: Types of future developments that affect re-identification.

## 5 Legal consequences of successful re-identification of identity-reduced data

When identity-reduced data is at least partially re-identified, the supposition that it was pseudonymous (i.e., personal) has been proven correct, while the supposition that it was anonymous has been proven wrong.

Controllers who treated the data as pseudonymous were aware that the requirements of the GDPR apply and have therefore already identified a valid legal basis for their processing and implemented appropriate technical and organizational measures. Consequently, the re-identification event does not change the situation.

In contrast, controllers who treated the data as anonymous supposed that the requirements of the GDPR do not apply to their processing. Thus, the re-identification event is a significant change that leads to logical consequences. These consequences are discussed in the following.

### 5.1 Possible sanctions

A re-identification event typically indicates that controllers who created the identity-reduced data and controllers who process such data fail to comply with the requirements of the GDPR.

Supervisory Authorities are in possession of corrective powers (see Art. 58(2) GDPR) including administrative fines. The question poses itself whether the incompliance of controllers can lead to differ forms of sanctions. When considering the possibility of fines, supervisory authorities have to take into account (among others) the following factors (see Art. 83(2) GDPR):

- Possible negligence,
- the manner in which the infringement became known to the supervisory authority,
- efforts to remedy the infringement, i.e. become compliant, and
- efforts to mitigate possible adverse effects of the infringement.

The first factor will be discussed in the following, the latter three are subject of the following sub-sections.

In the context of a re-identification event, negligence concerns the supposition that the data were anonymous. Whether this supposition was reasonable or negligent depends on the effectiveness of identity-reduction and re-identification methods used and the currently perceived threat landscape<sup>7</sup>.

Assume for example a controller who performed an identity-reduction with a state-of-the-art method. Unexpectedly, for example due to a previously unknown flaw in the concept or implementation of the method, a re-identification becomes possible. Further, the controller learns about this vulnerability by monitoring the state-of-the-art and takes remedial actions even before anyone could attempt to re-identify the controller's data. Such responsible behavior will likely be taken into consideration by authorities or courts when choosing a proportional corrective power.

In contrast, assume a controller who simply deletes obviously identifying attributes (such as name or e-mail address) from a dataset and then supposes that it was anonymous. A re-identification event then shows that this supposition was erroneous. Since the applied identity-reduction method is generally known to be vulnerable to re-identification, negligence of the controller regarding the choice of the method may likely be assumed.

The kind of responsible behavior that implements due care should be made explicit by some polity or guidelines. How responsible behavior can be partially guaranteed through participation in an ecosystem is described in Section 6.

## 5.2 Notification of re-identification event

The re-identification of previously supposed anonymous data likely constitutes a *personal data breach* according to Art. 4(12) GDPR and thus potentially requires notification of the competent supervisory authority according to Art. 33 GDPR and possibly communication of the breach to affected data subjects according to Art. 34 GDPR.

Art. 4(12) GDPR defines a *personal data breach* as “a breach of security leading to the accidental or unlawful destruction, loss, alteration, unauthorized disclosure of, or access to, personal data transmitted, stored or otherwise processed”.

Supposedly anonymous data is not protected against unauthorized disclosure with appropriate measures in support of confidentiality. A breach of confidentiality is considered a breach of security (see Art. 5(1)(f) GDPR). If identity-reduction is considered to be a measure comparable to a measure of access control, a re-identification event is comparable to the failure of a security measure. Hence, the GDPR's definition of *personal data breach* seems to apply.

Art. 33(1) GDPR states: “In the case of a personal data breach, the controller shall without undue delay and, where feasible, not later than 72 hours after having become aware of it, notify the personal data breach to the [competent] supervisory authority [...], unless the personal data breach is unlikely to result in a risk to the rights and freedoms of natural persons.”

A formal notification is thus only required, if the breach results in a risk for data subjects. This largely depends on the nature and availability of the re-identified data.

In the case where re-identified data was found for sale on the dark web and the data is sensitive (see Art. 9 GDPR on special categories of data), there is likely a significant risk and consequently, formal notification is required.

---

<sup>7</sup> Note that the proposed re-identification ecosystem foresees that its authority regularly publishes reports on the effectiveness of identity-reduction methods and re-identification methods).

In the case where re-identification was achieved by a re-identification researcher, the identified data was kept confidential, and the discovered vulnerability was responsibly disclosed (see Section 6), the risk resulting from the event is likely insignificant and does not require notification.

Like formal notification, also communication to the breach to data subjects according to Art. 34(1) GDPR is only necessary for cases of a significant risk to data subjects.

This illustrates the benefits of re-identification research with a well-defined protocol of responsible disclosure that aims to detect vulnerabilities while keeping the risk to data subjects insignificant (see Section 6 for details).

### 5.3 Actions to reach compliance after a re-identification event

Controllers who suppose anonymity think that their processing falls outside of the GDPR. Consequently, they fail to comply with the requirements of the GDPR. When they learn though a re-identification of their own or similar data, that the data is personal after all, they have different option on how to comply with the GDPR. These include the following:

1. Terminate the processing activity by deleting all data.
2. Use a more effective identity-reduction transformation and continue the processing activity still supposing anonymity,

Input to the more effective identity-reduction transformation are then either the original data, or the previous identity-reduced data that were found vulnerable without additional measures. An example for the latter case are k-anonymous data in which certain value intervals are merged or certain attributes are suppressed.

3. Continue the processing treating the data as personal. This includes at least the following:
  - Definition of explicit purposes of processing,
  - identifying a valid legal bases, and
  - implement appropriate measures, most notably in support of confidentiality

Any of these options aim to guarantee GDPR-compliance in the present and future. How to address effects of the past is discussed in the following sub-section.

### 5.4 Redressing the past violation

Changing the past is not possible. Therefore, redressing the past consists of limiting adverse consequences in the present and future. This concept is sometimes called “damage control”.

Two kinds of damage control are possible after a re-identification event:

- Notify all recipients of the previously supposed anonymous data that the data was found to be personal. This enables these recipients to take the appropriate actions to become compliant (see previous subsection).
- Where necessary and possible, notify data subjects to enable them to protect themselves as much as is still possible against negative consequences.

How to notify recipients of the compromised dataset and actors processing similar dataset that are likely also have become vulnerable is the motivation for a proposed ecosystem that can facilitate such notifications. This is described in the following section.

## 6 Monitoring of the re-identification risk and responsible disclosure

This section addresses the question of how a controller can monitor the re-identification risk its identity-reduced data. To gain insight in this question, it first draws an analogy to the ecosystem used in information security. On this basis, it then proposes a similar ecosystem that could be created in a given context, such as a data space, for the monitoring of the re-identification risk.

### 6.1 The information security ecosystem as an analogy

In information security, the monitoring of vulnerabilities is based on a well-established ecosystem. This ecosystem is illustrated in Figure 5.

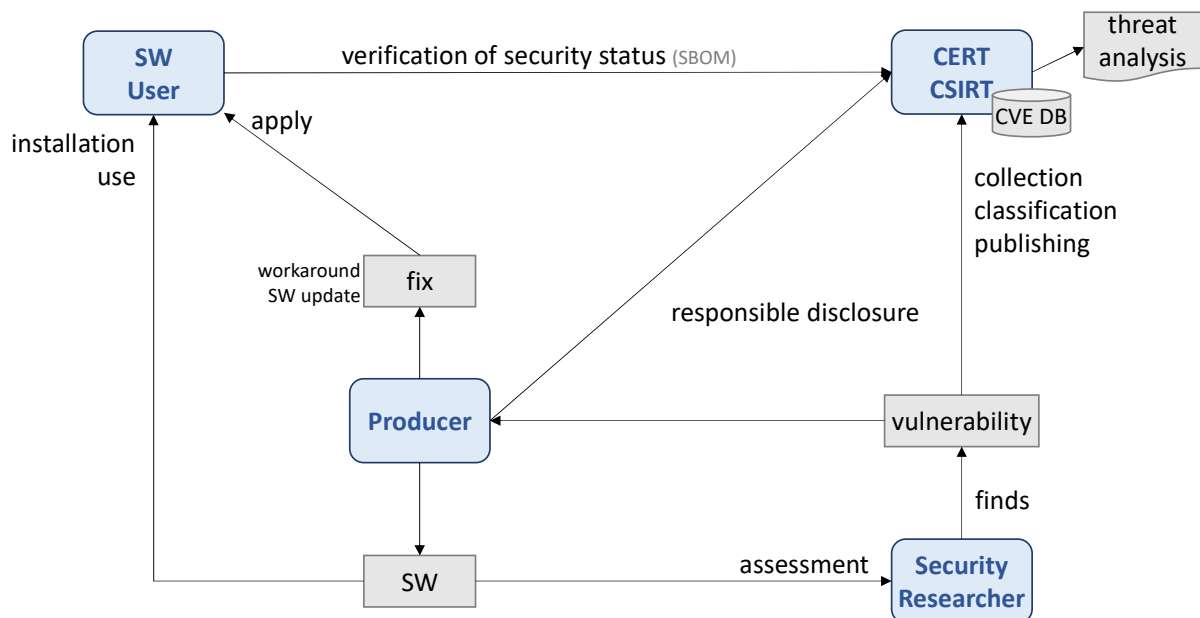


Figure 5: Monitoring of vulnerabilities in the information security ecosystem.

The ecosystem brings together a multitude of players who interact for the purpose of identifying and fixing security vulnerabilities. The types of players are the following:

- Software **producers** who create a piece of software (**SW**) that may be subject to vulnerabilities. The producer is authoritative for issuing **fixes** for the vulnerability.
- The software is then installed and used by a multitude of **Software Users**. These must be able to:
  - **assess the security status** of their system, namely by finding out whether they are affected by **known vulnerabilities**, and
  - obtain and **apply fixes** for these vulnerabilities.

- Detecting vulnerabilities is a difficult task that typically requires the expertise of **Security Researchers**. These assess a given piece of software and may be able to **identify vulnerabilities**. They then participate in a protocol of **responsible disclosure** that aims to get fixes rolled out to users before the vulnerability becomes known to (and potentially exploited by) malicious parties. A part of responsible disclosure is to make the software producer aware of the vulnerability to enable the creation of a fix. Protocols typically also require interaction with a CERT, CSIRT.
- **CERTs** (cyber emergency response teams) and **CSIRTs** (cyber security incident response teams) are typically established as part of a policy action in a certain domain (such as a nation) and context (such as the EU NIS 2 Directive). Their responsibility includes to collect, classify, and disseminate information about relevant known vulnerabilities and their fixes. Vulnerabilities are organized as CVEs (Common Vulnerabilities and Exposures, possibly with support of CVE Numbering Authorities) and assigned a unique identifier (the CVE Number). Software users can then query CERTs/CSIRTs to assess the security status of their systems. In addition, CERTs/CSIRTs can use the collected information to create generalized reports on the current threat landscape.

This well-established ecosystem is effective due to the distribution of responsibilities to different roles. No individual participant could likely assume all necessary roles.

## 6.2 Proposal of an ecosystem to monitor the re-identification risk

Based on the analogy to the information security ecosystem, this subsection proposes an ecosystem suited to monitor the re-identification risk. It can be created by establishing a specific authority and issuing regulation that mandates appropriate participation by the necessary players.

The proposed ecosystem is illustrated in Figure 6. It shows the various players and their interactions. The latter are represented by arrows and visualize the information that is exchanged. These interactions are numbered for easy reference in the text.

The similarity between the ecosystems is apparent from the figures. There is a strong similarity between actors and artefacts. The main difference lies in the relationship between the producer/anonymizer player and the User: While the relation is direct and strong in the case of information security, it cannot be relied on to exist in the proposed ecosystem. This is the case because an anonymizer can publish supposedly anonymous data. Here, anonymizers may not know the users who either obtain the data directly but without registering their identity or obtain the data indirectly. Further, while in information security the strong connection is created by a fix, in the proposed ecosystem there is no equivalent to a fix.

A second difference is motivated from the objective to prevent the possibility that malicious parties preventively retain supposedly data in order to receive help from the ecosystem to know which data can be re-identified.

Apart from these differences, the similarity is apparent in the equivalence of roles. This is illustrated by the following comparison of actors:

- While producers issue software, anonymizers create identity-reduced data sets.
- These are obtained by users of software or datasets, respectively.
- In both cases, researchers with relevant expertise identify vulnerabilities.
- In both cases, an authority enables “holds together” the ecosystem. While CERTs/CSIRTs are well-established, the proposed authorities are a new concept and instrumental in the creation of the ecosystem.

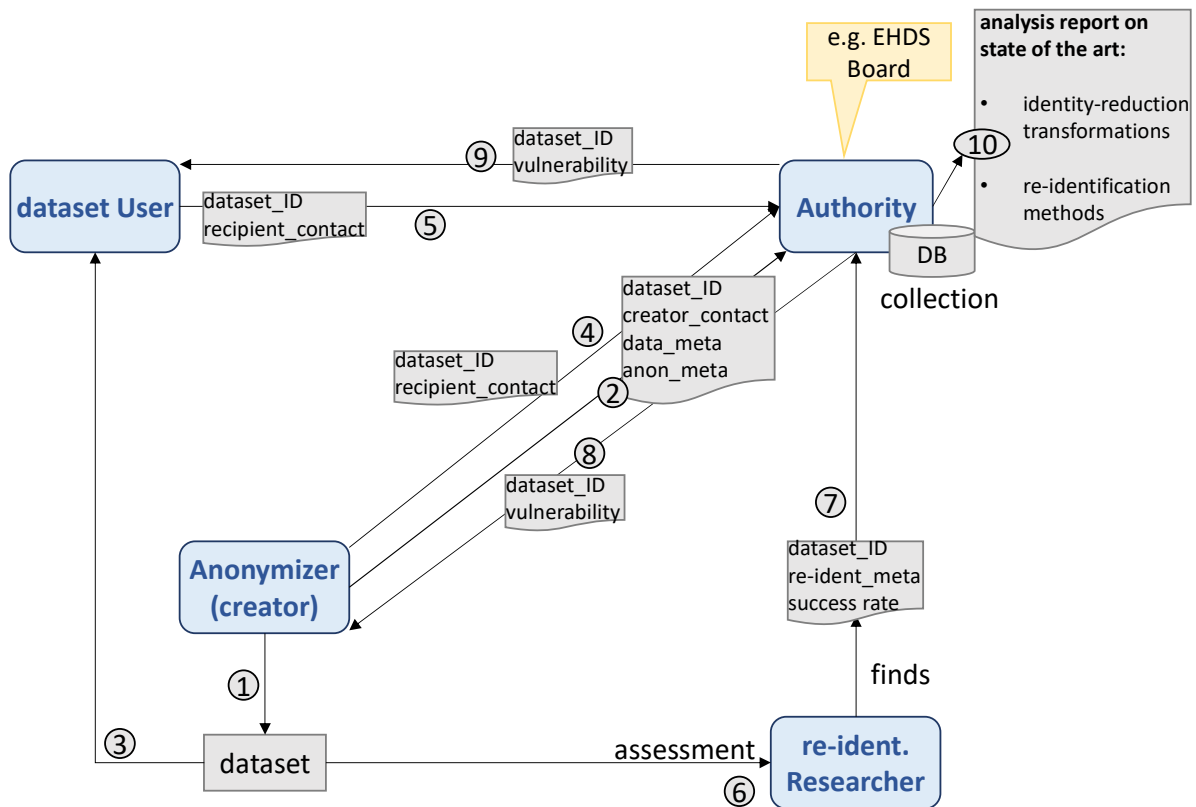


Figure 6: Proposed ecosystem for monitoring re-identification risk. .

The following describes the interaction shown in the figure in more detail:

1. An Anonymizer anonymizes a personal data set resulting in an identity-reduced dataset.
2. The Anonymizer communicates this fact to the Authority, in particular, it sends:
  - a. An ID of the resulting dataset (e.g., a digest),
  - b. An address where the Anonymizer can be contacted (e.g., a dedicated e-mail address),
  - c. Metadata about the dataset,
  - d. Metadata about the applied information-reduction transformation<sup>8</sup>. This information is collected by the Authority in a database (DB).
3. The dataset is then obtained by a Dataset User. This can either happen with the knowledge of the Anonymizer (e.g., based on controlled transfer) or without the Anonymizers knowledge (particularly when the dataset is published without keeping track of downloaders).
4. Where the Anonymizer knows about the recipient, the transfer is reported to the authority who registers it in its DB. In particular, this requires the following information:
  - a. An ID of the dataset,

<sup>8</sup> The proposed taxonomy of possible ways to claim anonymity may be part of such metadata. See AnoMed Deliverable D4.9.5 for detail.

- b. An address where the recipient can be contacted.
5. The recipient also registers its role as Dataset User with the Authority. This redundant registration covers cases where Dataset Users obtain the dataset (possibly indirectly) without knowledge of the Anonymizer. The redundancy in the registration aims to increase the registration rate.
6. A special kind of dataset Users are re-identification Researchers. They assess the protection of a dataset against re-identification. If they register as Dataset Users, they get informed by the Authority of possible future successful re-identification attempts of the data, even in the case that they were unsuccessful.
7. Re-identification Researchers can report the result of this assessment to the Authority. In particular, they report the following information:
  - a. The ID of the dataset,
  - b. metadata about the re-identification method used,
  - c. The success rate of the re-identification attempt. Note that also reports of unsuccessful re-identification are of value in the ecosystem. Also note that re-identification can be only partial, since only a certain percentage of data subjects was could be successfully identified.
8. When receiving a report of a successful re-identification, the Authority notifies the original Anonymizer of the fact. The Anonymizer, now informed that the data that was supposed to be anonymous is indeed personal, can take the necessary steps to conform with the GDPR. Most importantly, if the data was published, it has to be taken down. Anonymizers may also learn, beyond the single affected dataset, that the identity-reduction transformation they used no longer provides sufficient protection against re-identification and that they have to use stronger methods from now on. While this communication is primarily a push from the Authority to the Anonymizer, registered parties should also be able to pull the re-identification status information. A restriction to solely registered parties prevents potential malicious parties from retaining supposedly anonymous datasets and receiving help in identifying the vulnerable ones.
9. The Authority also notifies all registered Dataset Users of the discovered vulnerability. In the case that all Users are indeed registered and react lawfully to the vulnerability notification, it is possible to prevent further dissemination of the vulnerable data set to other parties.
10. In addition to managing the vulnerability status of individual data sets, the Authority can use the collected information to regularly issue an analysis on the effectiveness of various identity-reduction transformations and re-identification methods. Such reports can define the state of the art and inform policies that regulate a certain context, such as a data space.

## 7 Role of the proposed monitoring ecosystem

To further detail the proposal of the monitoring ecosystem in the last section, the following discussion illustrates how it supports the obligations of actors and how it compares to strategies without such an ecosystem.

Controllers of re-identified, previously suppose anonymous data are subject to obligations described in Section 5. Most importantly for this discussions, controllers must notify all recipients of the re-identification event. In addition, controllers, after an identity-reduction, have an obligation of reviewing the effectiveness of this measure.

How notification of recipients is supported by the proposed ecosystem was already discussed in the previous section; the following therefore focuses on the obligation of reviewing the effectiveness of the identity-reduction.

The main obligations for controllers is stated in Art. 24(1) GDPR and consist in the implementation of appropriate technical and organizational **measures** to ensure and demonstrate compliance. It further states that “Those measures shall be reviewed and updated where necessary.”

Assuming that identity-reduction is such a measure, the GDPR expresses an obligation for the **time after** the effected identity-reduction to **monitor** the effectiveness of this measure and **react** accordingly should the **effectiveness** of the measure be **deemed insufficient**. This temporally open-ended obligation ends only when there is no longer a risk (to the rights and freedoms of natural persons).

Note that this obligation stands in **stark contrast to** the “forget” in a strategy of “**anonymize, publish, and forget**”.

This poses the question, what exactly a controller needs to do to review the effectiveness of an identity-reduction measure.

Since most controllers lack the necessary specialized expertise to determine the effectiveness of identity-reduction, a collaboration with experts is necessary. In the proposed ecosystem, this is achieved by including and incentivizing re-identification researchers. The ecosystem promises a more realistic approach than a direct relation between controllers and such researchers.

Controllers also need to be concerned beyond their own identity-reduced data whether similar data has become vulnerable to re-identification. Again, they typically lack the necessary expertise to determine which identity-reduction methods are vulnerable and which are still save. It may also be practically difficult to monitor all re-identification event that occur. Therefore, they require assistance of a party that monitors re-identification events at a large scale and translates them to concrete recommendations on which identity-reduction transformations are considered safe. For this purpose, the ecosystem contains an authority who can cover these tasks.

By supporting early detection of vulnerabilities and a rapid mitigation of resulting risks, the proposed ecosystem can be considered a middle ground between an “‘anonymize’, publish, and forget” approach, and treating all data as pseudonymous from the beginning.

This middle ground position is illustrated in Table 1. In particular, the “‘anonymize’, publish, and forget” strategy is shown on the left and the supposition of pseudonymity from the start on the right.

What the middle column adds is infrastructure to detect and react to re-identification events. When comparing implementation effort, the middle ground solution adds the very moderate effort of participating in the proposed ecosystem. This is much less onerous than implementing all necessary measures for GDPR compliance.

A major benefit over an “‘anonymize’, publish, and forget” strategy is the ability of the ecosystem to detect relevant re-identification events. While the former strategy fails to provide any protection

even after such an event, the ecosystem facilitates that appropriate protection can be implemented “just in time”.

Since recipients of identity-reduced data register with the ecosystem authority, the data can still be easily shared. Even in the more careful approach of disclosing data only after receiving a contact address, the restriction to data sharing is minimal and cannot be compared to that of sharing pseudonymous data, likely only based on legal agreements.

Table 1: Middle ground between supposed anonymity and supposed pseudonymity.

strategy	always supposed anonymous (no monitoring)	start with supposed anonymous, detect and react to change	always supposed pseudonymous
implementation effort	none: “anonymize” – publish – forget	participation in re-identification ecosystem	implement full mitigation measures
breach detection	unlikely	likely	unnecessary
protection	no protection	protection as reaction to breach	protection always present
usability	unrestricted use	unrestricted use up to breach	use possibly hindered by measures (e.g. confidentiality)

## 8 Incubation of a monitoring ecosystem

In contrast to the information security ecosystem, an ecosystem for re-identification risk does not yet exist. If deemed beneficial, one or several (e.g., per dataspace) of these ecosystems have to be incubated. Objective of such an incubation is to bring all necessary players to participate in the ecosystem and to render the ecosystem sustainable.

### 8.1 Establishment of the authority and its infrastructure

The key player of the ecosystem that does not currently exist is its authority. The following action is required:

- Establishment of the authority and defining its tasks and competences. This is typically achieved through a legal act or regulation.
- Implementation of the necessary IT infrastructure to be operated by the authority, namely a data base system and the necessary messaging.
- Continued and staffing and financing of the authority.

### 8.2 Obligations for anonymizers, Dataset Users, and re-identification researchers

The participation of anonymizers and dataset users can be guaranteed through the creation of according obligations.

In particular, the following rules pursue this objective:

- Anonymizers shall:
  - Register the creation of identity-reduced data sets with the authority,
  - guarantee that a contact address of recipients is collected when the identity-reduced dataset is disclosed (transferred) to them,

- notify the authority of the recipient's address during such disclosures,
  - take appropriate remedial action when notified of re-identification events by the authority (see Section 5), and
  - monitor the authority's reports on the state-of-the-art and follow its recommendations in current and future identity-reduction activities.
- Dataset Users shall:
    - Register the reception of an identity-reduced dataset with the authority,
    - take appropriate remedial action when notified of re-identification events by the authority (see Section 5), and
    - monitor the authority's reports on the state-of-the-art and follow its recommendations in current and future identity-reduction activities.
  - Re-Identification Researchers shall:
    - Register with the authority as re-identification researcher,
    - register concrete attempts to re-identify a given dataset,
    - implement appropriate data protection safeguards for the attempt (see Section 9),
    - inform the authority about the outcome of the attempt following a standardized reporting format.

There are different options to bind the different players to these obligations. One is to integrate them into some regulation, for example an act that establishes a certain data space. Another option is to convey these obligations through contracts. For example, the disclosure of an identity-reduced dataset by the anonymizer can be executed on the basis of a contractual agreement in which the recipient assumes the necessary obligations.

### 8.3 Incentives for re-identification and anonymity researchers

Re-identification researchers fill an important role in the ecosystem. Without their activity, the ecosystem would fall apart. Since research activities cannot usually be mandated, appropriate incentives have to be created. The following discusses the importance of avoiding disincentives and discusses possible incentives.

There seems to be a line of thought that any re-identification attempt is malicious and has to be prohibited. Evidently, a possible prohibition represents a powerful disincentive for any well-intended, responsible researcher. Such prohibitions therefore need to be avoided.

The situation bears some similarities with that of ethical hacking in which researchers attempt to find vulnerabilities of information security. Alabi Samuel has compiled a summary of the legality of penetration testing in various countries<sup>9</sup>. The scientific services of the German parliament have also compiled a report about the criminalization of hacking<sup>10</sup> (in German). A presentation on the same

---

<sup>9</sup> <https://medium.com/my-identity-pay/countries-where-penetration-testing-is-illegal-and-legal-3f9f8c853aea>

<sup>10</sup> <https://www.bundestag.de/resource/blob/1005444/ed435cb1a5311bb688385a81f295c8a3/WD-7-104-23-pdf.pdf>

topic is available from Latvia's cert<sup>11</sup>. A possible amendment to the EU Cybersecurity Act addresses the issue under the key word of "managed security services<sup>12</sup>". Ethical hacking has become legal in Belgium based on their whistleblower law<sup>13</sup>. The regulation on "Coordinated Vulnerability Disclosure Policy (CVDP) and Vulnerability Detection Reward Program"<sup>14</sup> provides detailed procedures for ethical hacking.

Some prohibitions of re-identification can be found in current legislation. In particular, Article 44(3) EHDS proposal states that "[d]ata users shall not re-identify the electronic health data provided to them in pseudonymised format." This clause is limited to pseudonymous data and does not restrict the possibilities of re-identification researchers in the proposed ecosystem.

Paragraph 7(2) of the German "Gesundheitsdatennutzungsgesetz"<sup>15</sup> seems to impose a wider prohibition of re-identification. It states (in German) "Bereitgestellte Daten dürfen nicht zum Zwecke der Herstellung eines Personenbezugs oder [...] verarbeitet werden." This wording may well render the desired activities by re-identification researchers an offense.

Besides avoiding disincentives, actual incentives can be provided in particular by the ecosystem authority. These include the following:

- Guidance to re-identification researchers to identify and obtain suitable identity-reduced datasets. A variation of this incentive would be a competition to re-identify a given (possibly synthetic) dataset. (This is similar to the well-known "bounty hunting").
- Organizing a scientific conference where re-identification researchers can exchange their experiences and in whose proceedings the current state-of-the-art of re-identification methods can be documented. Participation and publication in such a conference, similarly to publication in renowned journals, can motivate researchers through providing reputation. Note that a similar conference is already organized yearly by the French supervisory authority CNIL<sup>16</sup>.
- Reward prizes for the best scientific work on re-identification. These can incentivize work through reputation and possibly financial reward.
- Positioning of re-identification and anonymization research as high priority in research funding programs such as "Horizon Europe".

The above listed incentives have a predominantly scientific focus and target the scientific community. They could be complemented with incentives that foster technology transfer from the scientific community to the practice of anonymizers. Possible incentive include the following:

- High reputation publishing and possible prizes for educational material and practical instruction on the state-of-the-art of anonymization and re-identification techniques.
- Classification of new anonymization and re-identification techniques according to their Technology Readiness Level (TRL)

---

<sup>11</sup> [https://cert.lv/uploads/pasakumi/Nathalie\\_Falot.pdf](https://cert.lv/uploads/pasakumi/Nathalie_Falot.pdf).

<sup>12</sup> [https://www.europarl.europa.eu/RegData/etudes/BRIE/2023/754556/EPRS\\_BRI\(2023\)754556\\_EN.pdf](https://www.europarl.europa.eu/RegData/etudes/BRIE/2023/754556/EPRS_BRI(2023)754556_EN.pdf)

<sup>13</sup> <https://www.law.kuleuven.be/citip/blog/belgium-legalises-ethical-hacking-a-threat-or-an-opportunity-for-cybersecurity/>.

<sup>14</sup> <https://ccb.belgium.be/regulation/cvdp>

<sup>15</sup> <https://www.recht.bund.de/bgbl/1/2024/102/VO.html>

<sup>16</sup> <https://www.cnil.fr/en/call-for-papers-privacy-research-day-2025>

- Publication of success stories of applying the state-of-the-art anonymization and re-identification techniques in practice.
- Indorsement of suitable open source implementations of state-of-the-art techniques.

## 9 Legal requirements for attempted re-identification of anonymous data

Re-identification researchers expect to be successful with a significant probability. If they expected an insignificant chance of success, the motivation for a re-identification attempt would likely be missing. Consequently, the processing of personal data is planned and the requirements of the GDPR apply. This section discusses the most relevant requirements in more detail.

The situation is similar to that of “data protection prophylactics” (a concept coined in German as “Datenschutz-Vorsorge”) that has been discussed in the literature<sup>17</sup>. Instead of re-identification that leads to personal data, data protection prophylactics is concerned with ethical hacking that results in obtaining personal data with some likelihood.

The difference of re-identification and data protection prophylactics is that the type (and thus sensitivity) and volume of potentially created personal data us much better known. This renders it much easier to identify and then mitigate possible data protection risks.

Some of the main data protection requirements as they pertain to re-identification research are described in the following:

- The **purposes** of the processing have to be clearly **specified** and processing for other purposes is prohibited. In particular, the purpose is **limited** to “identifying possible re-identification vulnerabilities in datasets”. This excludes for example researchers looking up persons they know in re-identified data. The limitation of activities to the specified purposes must be made clear to all parties involved (for example to participating students).
- The processing necessary for a re-identification attempt requires the identification of a suitable **legal basis** according to Art. 6(1) GDPR. The detailed discussion of suitable legal bases is left to legal scholars. Since the processing activity is desired for the benefit of society, it is assumed here that a suitable legal basis can be identified without excessive difficulty.
- The third requirement is to implement appropriate technical and organizational measures in support of the data protection principle. These include the following:
  - **Data minimization** mandates that only as much data is used as required for the purposes. In a re-identification attempt based on linkage, for example, it may either take the form of prior **deletion of unnecessary attributes**, or limiting the study to only a **sub-sample of data subjects**.
  - Only the **necessary level of identification** shall be re-constructed. For example, to assess the effectiveness of a given re-identification method, it may be sufficient to

---

<sup>17</sup> Boll, Alina & Stummer, Sarah & Selzer, Annika. (2024). Datenschutz-Vorsorge. Datenschutz und Datensicherheit - DuD. 48. 172-176. 10.1007/s11623-023-1902-x, [https://www.researchgate.net/publication/379080835\\_Datenschutz-Vorsorge](https://www.researchgate.net/publication/379080835_Datenschutz-Vorsorge) (in German).

limit the work to a first step of reconstructing individual-level data. The second step of fully identifying this data through linkage is possibly unnecessary for this purpose. In this approach, the re-identification is only partial, producing **solely pseudonymous**, not fully identified data. This kind of limitation of the processing evidently reduces the risk for data subjects.

- The disclosure of the actually processed data has then to be minimized to parties that actually require access to the data in order to reach the specified purposes. Disclosure to additional parties raises the risk that the data is used for other purposes. The measures to ensure this are those of **confidentiality** and **access control**. An appropriate implementation of this measure affect decisions such as on the possibility to employ private laptops for this work, take data home for processing over the weekend, or to send data as unencrypted e-mail attachments to colleagues.

## 10 Conclusions

This report has proposed an empirical, ecosystem-based approach to managing the residual re-identification risk in data spaces where identity-reduced data is freely shared. This ecosystem mirrors the well-established information security ecosystem that is for example established by the Cyber Resilience Act. The proposal represents a middleground between sharing supposedly anonymous data without protection and treating all data as personal (i.e., pseudonymous). The proposed ecosystem then provides protection of data subjects through early detection and reaction by ethical parties to newly discovered re-identification vulnerabilities. It further implements a continuous observatory of what level of identity-reduction is necessary in order to avoid exposing data subjects to excessive risk.