# Anonymisation: Privitar's recommendations to the ICO

## February 2020

Guy Cohen
Marcus Grazette

# Privitar – who we are

## Privitar

Privitar is an enterprise software company headquartered in London, with a global client-base across North America, Europe, and Asia.

Privitar is leading the development and adoption of privacy engineering technology, enabling our customers to innovate and leverage data with an uncompromising approach to data privacy.

Privitar's mission is to promote and facilitate the ethical and safe use of valuable data assets. Using leading privacy engineering techniques, we help companies get maximum value from data while preserving customer's privacy.

## Privitar Research and Policy

This paper was written by Privitar's Policy Team. Assuring data privacy in an increasingly digital world is a significant and pressing challenge. To create and market effective data privacy software Privitar must have a deep understanding of all facets of the problem. We commit significant resources to our Research and Policy teams to help us do that.

The teams work with world leading academics, policy makers, regulators, and other experts to think deeply about how technology can help us to preserve privacy while utilising data. We make this investment because we're passionate about solving the privacy problem, and because we believe it's only by grappling with the hard issues that we're going to be able to build and market the solutions that our clients need. This work also enables us to produce high-quality, thought-leading content that helps companies in the selection and use of privacy enhancing technologies, and at times shapes the public debate.

# Introduction

The benefits of using data are clear, as are the risks to individuals' privacy as organisations collect and process ever increasing volumes of data. We strongly agree with HM Treasury's conclusion that "data-driven innovation holds the key to addressing some of the most significant challenges facing modern Britain"[1]. But innovation can only happen if organisations are able to use data.

There can be genuine trade-offs in protecting privacy whilst innovating with data. However, for many organisations the barriers stem from a lack of clarity, meaning that data controllers hesitate to use or share information for fear of falling foul of the rules. Advances can be made on both; privacy enhancing technologies (PETs) can improve the privacy-utility trade off, and clearer policy supported by more detailed guidance on implementation will reduce uncertainty. This paper will focus on the latter. Specifically it will lay our Privitar's recommendations for improving and clarifying anonymisation policy and guidance in the UK.

Anonymisation takes data outside of the data protection regime, opening up use cases that would not otherwise be possible. The bar to achieve anonymisation is rightly set high, as the data subject no longer benefits from the rights and protections that they would otherwise have. Anonymisation also comes at a cost to data utility and usability, meaning it will not be appropriate in all use cases.

We believe that mitigating privacy risk by applying controls to both the data and the environment, and relying on legal bases other than consent can allow processing within the scope of the data protection regime and will be a more appropriate approach for many use cases. This is particularly important given the evolving debate around consent, illustrated by the recent Joint Committee on Human Rights report on the right to privacy and the digital revolution[2], and the emergence of new data processing technologies. We therefore believe the ICO should support data controllers by providing guidance on when it is appropriate to anonymise data, and when another route may be more appropriate.

We have split our recommendations into two clusters. The first cluster brings together four areas where we recommend that the ICO should review their existing policy and provide additional guidance. First, the ICO should clarify the definition of "identifiable", a key concept in the assessment of re-identification risk. Second, the ICO should update the "motivated intruder" test. We believe that the attack based methodology for interpreting the "means reasonably likely" test in law is sensible, but that the current standard is too weak. Third, the ICO should provide guidance on what "objective factors" data controllers can and should consider when assessing re-identification risk, specifically how environmental factors should be considered. Finally, the ICO should provide guidance on applying the "motivated intruder" test in practice.

The second cluster brings together five more practical recommendations. First, the ICO should provide guidance on when anonymisation is an appropriate

---

[1] HM Treasury, The economic value of data: discussion paper, published Aug 2018
[2] UK Parliament, Joint Committee on Human Rights, The Right to Privacy (Article 8) and the Digital Revolution, published Oct 2019

**PRIVITAR**
Research & Policy

strategy, and when instead de-identification[3] and other controls are sufficient to process the data within the data protection regime. Second, the ICO should incentivise compliance through enforcement action targeting poor anonymisation practices. Third, the ICO should strengthen their health warning on the vulnerabilities of using hashing as a de-identification technique. Fourth, the ICO should establish a mechanism for active dialogue with relevant expert stakeholders, including the private sector. Finally, the ICO should lead the international debate, championing consistent interpretation across Europe and beyond.

# Cluster 1 – Core Policy and Guidance

This section includes four recommendations for the ICO to provide further guidance on the interpretation of key terms relating to anonymisation, which are currently unclear. This lack of clarity presents a risk; it could lead to an inconsistent application of the rules, missed opportunities, or unintentional non-compliance.

The extracts quoted below, with our emphasis added, frame this cluster of recommendations. Recital 26 of the GDPR defines "anonymous information" as

> *"...information which does not relate to an identified or <u>identifiable</u> natural person or to personal data rendered anonymous in such a manner that the data subject is <u>not or no longer identifiable."</u>*

Recital 26 also states that

> *"To determine whether a natural person is identifiable, <u>account should be taken of all the means reasonably likely to be used</u>, such as singling out, either by the controller or by another person to identify the natural person directly or indirectly."*

And that

> *"To ascertain whether means are reasonably likely to be used to identify the natural person, <u>account should be taken of all objective factors</u>, such as the costs of and the amount of time required for identification, taking into consideration the available technology at the time of the processing and technological developments."*

The first three recommendations below address issues with the interpretation of the concepts <u>underlined</u> in the extracts. The fourth considers how these combine to create an operational test which can be deployed by data controllers as they seek to assess re-identification risk.

---

[3]  We use the term "de-identification" to refer to any action which reduces the likelihood of an individual being identified. De-identification does not necessarily mean that the output data is anonymous.

PRIVITAR
Research & Policy

# Recommendation 1(a) - Clarify the definition of "identifiable"

The ICO should clarify their definition of "identifiable" and whether they agree with the definition provided in the Article 29 Working Party (A29WP) Opinion[4] (A29WP Opinion). If the ICO agrees with the definition (based on linkability, singling out, and inference), then the ICO should also define what constitutes a significant inference.

## Rationale

Our understanding, based on the ICO's 2012 Anonymisation Code of Practice[5] (ICO Code) and the Data Protection Act 2018[6] (DPA 2018), is that for identification to take place it is sufficient to "be able to establish a reliable connection between particular data and a known individual". The A29WP Opinion, which predates the GDPR, defines re-identification on the basis of three elements (our underlining added for emphasis):

> *"An effective anonymisation solution prevents all parties from <u>singling out</u> an individual in a dataset, from <u>linking two records within a dataset</u> (or between two separate datasets) and from <u>inferring any information in such dataset</u>."*

Other European regulators, such as the Commission Nationale de l'Lnformatique et des Libertés (CNIL)[7] in France and the Data Protection Commission[8] (DPC) in Ireland have adopted the A29WP's three-part definition in their guidance or decisions on anonymisation.

We believe that a clear definition of "identifiable" is important to enable consistent standards of anonymisation. However, the concept of inference is problematic. The A29WP Opinion defines inference as:

> *"...the possibility to deduce, with significant probability, the value of an attribute from the values of a set of other attributes."*

This definition opens further questions on the meaning of "significant". Data analysis will often include some statistical modelling, which is itself probabilistic and due to errors in data collection and entry, even the "raw" data often will also come with some uncertainty. So if most or all data is to some extent probabilistic, what level of confidence is needed for an inference to constitute identification? We have heard the concept of "plausible deniability" used as an equivalent test,

---

[4] Article 29 Data Protection Working Party, Opinion 05/2014 on Anonymisation Techniques, published April 2014
[5] ICO, Anonymisation: managing data protection risk code of practice, published 2012
[6] DPA 2018 section 171(2)(a) "personal data is de-identified it has been processed in such a manner that it can no longer be attributed, without more, to a specific data subject". Re-identification in that context is defined as "taking steps which result in the information no longer being de-identified" within the meaning of that paragraph.
[7] CNIL, L'anonymisation des données, un traitement clé pour l'open data, accessed Oct 2019
[8] Data Protection Commission, Guidance on anonymisation and pseudonymisation, published June 2019

PRIVITAR
Research & Policy

however this is also vulnerable to the same issues. We can illustrate this with an example:

**When is an inference 'significant'?**

One way to assess whether an inference is 'significant' may be to ask whether the inference has a material impact on an individual.

Imagine a fictitious health insurance company, DataSure Ltd, which gathers data about potential clients in order to assess whether a client is at high risk of developing diabetes. DataSure decides that it will not offer insurance to clients where they have a 30%+ confidence that the client is high risk.

Alice applies for an insurance policy. Based on the data available, DataSure is 10% confident that she is high risk.

Bob, a researcher, then publishes a one-off survey data from Alice's area which shows that she is in fact 20% likely to be high risk.

Charlie, a second researcher, then publishes two data releases each of which increases the confidence that Alice is at high risk of developing diabetes by 7.5%. With the first publication, DataSure becomes aware that likelihood of Alice being high risk is 27.5%, then 35% with the second publication.

Although each of Charlie's publications changed the confidence interval by the same amount (+7.5%), the cumulative impact was to push Alice over DataSure's confidence threshold, leading to her application being rejected.

Should we consider Bob's first publication anonymous, but Charlie's third personal data even though Bob revealed more information? Or is Charlie's second publication anonymous but third personal data even though they led to the same change in confidence, because the third led to a significant impact on Alice? What if they happened at the same time?

The problem can be thought of as analogous to the heap paradox, whereby grains of sand are taken from a heap, until none are left, at which point it is clearly not a heap. At what point did it cease to become a heap? In the same way, when does what a release reveals about an individual cease to be significant?

This example shows how anonymisation can be complex and what is appropriate is highly contextual. What should be consistent across use cases is the method used to assess whether data is anonymous. As such we would recommend the ICO focus on clear guidance on the methodology as opposed to the specific controls required.

**PRIVITAR**
Research & Policy

# Recommendation 1(b) - Update the "motivated intruder" test to enable a proper consideration of "all the means reasonably likely to be used"

We believe that attack based evaluation is a sensible approach to defining "means reasonably likely". We recognise that one cannot draw a general conclusion about the robustness of a system based on one attacker's failure to compromise it. However, we believe that attack based evaluations provide the best practical response currently available to applying a legal standard of reasonableness. This approach does not demonstrate that a dataset is impervious to all possible attacks, rather that it is protected against known reasonably likely threats.

We believe that the ICO's current attack model, based on a "motivated intruder", is not appropriate because it is too weak and too narrow. We recommend that the ICO re-evaluate this test and consider using an array of attacker profiles, to reflect the fact that different attackers have different motivations and access to different auxiliary information, resources, and skills, and may resort to different types of unethical or criminal behaviour.

We recommend that the ICO require those evaluating anonymisation to consider a range of attackers. A minimum set of attackers should always be considered. Additional stronger attackers should also be considered when the data is reasonably likely to be of interest to those types of attackers.

This recommendation links to 1(c) below on defining "all objective factors" to be taken into account when assessing re-identification risk, 1(d) which suggests the ICO provide guidance on applying the test in practice, and Annex A which explores the question of attackers who may break the law.

## Rationale

Penetration testing is well-established in IT security, where an authorised tester identifies vulnerabilities in a system using the same tools and techniques that an attacker might. The tester then provides a report, allowing any vulnerabilities they identify to be fixed.

A similar approach can be deployed to assess the re-identification risk of a particular dataset. A hypothetical attacker would attempt to re-identify the data, using "all the means reasonably likely to be used".  The ICO's 2012 Code suggests using an attack based evaluation, based on an attacker profile called the "motivated intruder". The motivated intruder:

*"is reasonably competent, has access to resources such as the internet, libraries, and all public documents, and would employ investigative techniques such as making enquiries of people who may have additional knowledge of the identity of the data subject or advertising for anyone with information to come forward. The 'motivated intruder' is not assumed to have any specialist knowledge such as computer hacking skills, or to have access to specialist equipment or to resort to criminality such as burglary, to gain access to data that is kept securely."*

**PRIVITAR**
Research & Policy

We believe this definition and its interpretation arose from a standard used by the Office for National Statistics (ONS) in 2002 and, while perhaps reasonable at that time, is no longer a realistic view of today's potential attackers. As such we do not believe that it reflects the "means reasonably likely" to be used today.

For example, when the ONS performed an intruder testing exercise in 2011 they gave the attackers just 3.5 hours to try to re-identify an individual[9]. We believe the standard should be raised to reflect that:

- There is more auxiliary information available today, making re-identification easier

- Compute power has fallen in price, making complex attacks cheaper

- More sophisticated attacks have become more accessible to more users due to new software packages and increased data skills amongst potential attackers. This means more attackers will have "specialist knowledge"

- The value of re-identification may have increased due to online data marketplaces and other societal changes, meaning more parties are motivated to invest more resources in such attacks, including criminals and unethical companies

In addition to potentially being weaker than some modern attackers, the "motivated intruder" is also too simplistic in just imagining one persona without any specialist capabilities. We instead recommend considering a range of attackers in making the assessment, to ensure that the assessment is closer to the actual likely threat. Over time the personas should continue to be updated to take account of new attacker profiles and changing capabilities of existing attackers.

Using a range of personas is in line with other ICO guidance on personal data, which refers to "a determined person with a particular reason to want to identify individuals" listing "investigative journalists, estranged partners, stalkers and industrial spies" as examples[10]. It is also in line with *The Information Commissioner v Miller*[11] which considered whether the information being published would attract the attention of those with specific investigative skills: "The information in the spreadsheets is not such as is likely to attract those with investigative skills, such as a journalist, to attempt to identify individuals."

The table below provides some examples of attacker profiles. It is an illustrative list, not an exhaustive one.

---

[9] ONS, Brief notes on earlier intruder testing exercises for open (public use) datasets, accessed Nov 2019
[10] ICO, Key Definitions: What is personal data, accessed Nov 2019
[11] The Information Commissioner v Miller [2018] UKUT 229 (AC), accessed Dec 2019

PRIVITAR
Research & Policy

Table 1: Illustrative "motivated intruder" profiles

| Profile | Motivation | Capability | Criminality and breach of contract |
|---|---|---|---|
| 1. The Insider | Financial gain, to cause damage to their employer, or to learn about specific individual (such as partner or celebrity) | Varies across and within organisations. Generally limited time and compute. Some may have data science skills but some may not, likewise may have access to raw data or an understanding of internal controls on data use. May have access to significant auxiliary information and domain expertise, including potentially sensitive details. | May steal data, break NDAs, and not abide by contractual requirements for the data they have access to as part of their role, but unlikely to break other kinds of laws, such as hacking into other systems. |
| 2. The Unethical Company | Financial motivation e.g. for targeted marketing, manipulation, competitive advantage, or selling to other companies | Can leverage significant compute. Amount of personnel time will vary by company. Access to data science skills with the option to hire specialists if sufficiently motivated. Able to merge with their existing datasets or scrape data from public sources (e.g. social media). | Unlikely to directly break laws, but may bend them as far as is possible and operate in grey areas. |
| 3. The Academic | Highlight vulnerabilities and demonstrate novel attacks. Likely to publish results in academic venues | May have significant compute but more limited on personnel time. Very strong skills, including all known attacks. Likely to be aware of open data sources, but unlikely to have other information. | Unlikely to break the law or behave unethically. |
| 4. The Criminal | Financial gain, e.g. through fraud or blackmail. | Access to significant compute power and time. Potentially strong technical skills. Unlikely to have access to background information on individuals from offline sources (such as an employee knowing about their clients, or a health researcher inferring a characteristic from medical data) | Likely to ignore data protection law, contracts, and may break the law to access digital information, such as breached data purchased online or through hacking. |
| 5. The State Actor | Surveillance, national security interests, political | Significant time and compute power. Access to expert data scientists, specialised hardware | May break laws of other countries, e.g. |

PRIVITAR
Research & Policy

| | influence, destabilising institutions | and cutting edge research, including attacks which may be unknown to the public and researchers. Access to significant public and private auxiliary data, able to draw on data held by other institutions and gather data from offline sources. | through hacking or spying. |
|---|---|---|---|

We recommend that the ICO explore a range of attacker profiles and then create a minimum set which should always be considered by organisations carrying out an attack based assessment, including profiles one, two, and three above.
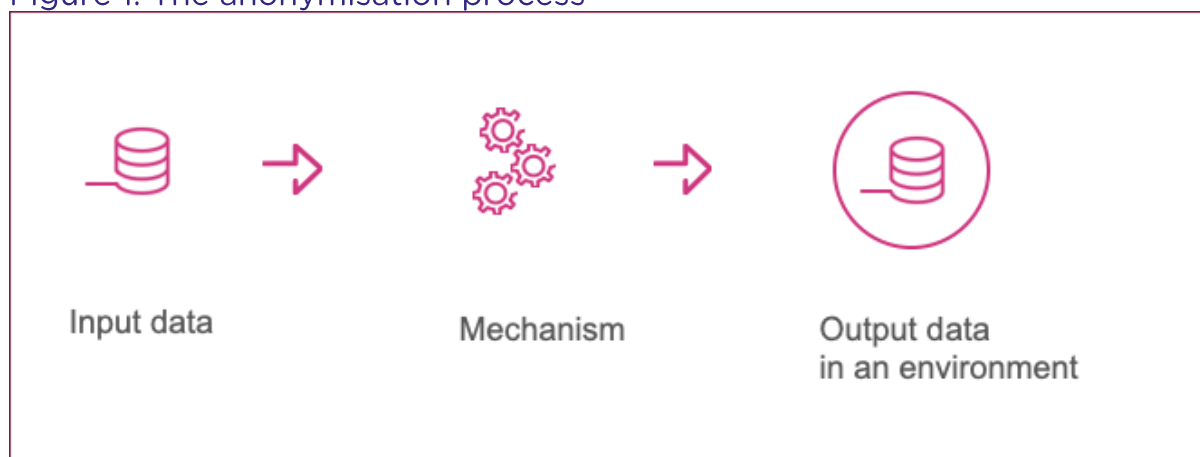
We suggest that additional attackers, such as attacker profiles four and five, be a conditional requirement, whereby they must be considered unless the evaluator can justify that an attacker of this kind would not be reasonably likely to be interested in the dataset in question. In this way the "means reasonably likely" becomes a two part assessment; which attacker profiles is it reasonable to consider, and then what could those attackers reasonably learn?

## Recommendation 1(c) - Provide guidance on "objective factors"

The ICO should provide guidance on what "objective factors" should be taken into account when assessing re-identification risk. Specifically the ICO should make clear that one should consider the environment the data is in, as well as the data.

### Rationale

### Figure 1: The anonymisation process



Input data　　　　Mechanism　　　　Output data in an environment

This diagram shows some input data, being processed through a mechanism, resulting in an output which exists in an environment. Some models consider anonymisation a property of the mechanism, others a property of the output, and others a property of an output in an environment.

**PRIVITAR**
Research & Policy

We see the concept of the environment as a key area of confusion in anonymisation policy. There is international inconsistency both on the factors which constitute the environment and on how those factors relate to de-identification risk and therefore to an assessment of whether data is anonymous (see Recommendation 2(e) and Annex A for more detail).

We believe that an assessment of "objective factors" should include factors relating both to the data and to the environment.

Factors relating to the data could include any transformation which reduces the likelihood that the data on its own could be used to re-identify an individual, for example tokenisation, generalisation, and small count suppression.

Factors relating to the environment could include any other constraints or controls which limit an attacker's ability to use the data for re-identification. For example access controls to prevent unauthorised access, contractual limitations such as nondisclosure agreements, background checks to ensure only trustworthy individuals are given access to the data or internal controls such as requiring two individuals to sign off on data provisioning.

One approach to considering data and environmental factors is the five safes model[12] which looks at:

- Safe people (e.g. vetting, accreditation, training)
- Safe projects (e.g. ethical, legal, or contractual controls)
- Safe settings (e.g. physical and IT security)
- Safe data (e.g. de-identification)
- Safe outputs (e.g. statistical disclosure controls)

The 2012 Code shows that the ICO considers the environment as part of the assessment, noting that "the risk involved will vary according to the local data environment and particularly who has access to information". Likewise, the Anonymisation Decision Making Framework (ADF)[13] notes that anonymisation is "heavily context dependent" and recommends considering the data and its environment as a total system.

We believe the ICO should clarify that anonymisation is a property of data in an environment, and explain what this means, and how different environments relate to the assessment of what is reasonably likely. This should include a description of the factors which make up the environment and the controls one can use to manage the environment, and how these controls and different environments affect the attack based assessment. To be clear, we believe that the concept of the data environment is broad enough to encompass data in the public domain. The public domain is just one type of environment.

Considering environmental controls allows assumptions to be made about an attacker's capabilities. For instance, environmental controls might allow one to assume that the 'criminal' attacker is unable to access the information, due to screening and security protections, meaning they could not re-identify anyone in

---

[12] UK Data service, Regulating Access to Data, accessed Nov 2019
[13] UK Anonymisation Network, Anonymisation Decision Making Framework, published 2016.

PRIVITAR
Research & Policy

the dataset in that environment because they could not access it in that environment.

In some instances it is not possible to apply environmental controls, such as when publishing information online or sharing widely, in these instances one must assume the attacker is able to access the information and can attack it to the best of their ability, meaning a much higher level of de-identification would be required to achieve anonymisation in that context. The level of de-identification required to achieve anonymisation will vary depending on the environment the data is being accessed in. The table below provides a non-exhaustive set of example environments:

Table 2: Illustrative data environments

| Publishing online | Data sharing between two parties | Access in secure trusted third-party location |
|---|---|---|
| Open access, data can be downloaded and linked with other sources without restrictions. | Data shared between organisations under strong contractual and technical controls which may deter and/or limit certain data actions, including linking. | Data held by a trusted third party in a secure location where only screened individuals may access the data, and only in a tightly controlled and monitored environment which does not allow for unauthorised data linking. |
| Must assume that any attacker might access the data and so would require strong de-identification of the data.

For many datasets and use cases there may not be a level of de-identification that results in both anonymisation and sufficient data utility. | If the data is to be deleted after a stated period of time and is watermarked such that if exfiltrated it could be traced and would result in the responsible employee/s facing disciplinary and legal action, then one can perhaps make certain assumptions about the attackers.

However, unlike the trusted third party model, the data is being held by a party who may have other interests and so allow broader data use and weaker environmental controls if they interfere with usability. So the risk will be contingent on the trustworthiness of the recipient and recipient organisation. | Here one could assume many types of attackers will not be able to access the data, and so the de-identification may only need to protect against those who have access, such as a curious insider.

In many of the examples where this is being done the data users are considered highly trustworthy, such as medical researchers where the researchers and their research may have been approved by an ethics panel. |

# Recommendation 1(d) - Provide guidance and support on applying the attack based assessments in practice

The ICO should provide guidance on carrying out attack based assessments, taking account of both the revised personas described above and the guidance on objective factors. The ICO may also wish to consider developing a certification scheme for attack based assessment, which could mirror the National Cyber Security Centre's (NCSC) CHECK standard for approved penetration testing providers.

## Rationale

Guidance on attack based assessments should provide a framework to structure a data controllers' approach. It could resemble guidance e.g. from the NCSC on

**PRIVITAR**
Research & Policy

penetration testing[14]. The guidance should cover different types of assessment, for example, in some situations a paper based assessment may suffice, whereas in others the data should actually be attacked.

The ONS has published an overview of their internal approach to intruder testing[15]. This involved recruiting volunteer intruders, equipping them with internet connected computers and allowing them around half a day with the data. We consider this to be a low bar and recommend that the ICO provide updated guidance on what is reasonable.

The ICO, in line with its position on certification, may wish to consider whether a specific certification scheme for attack based testing might be a useful tool to assist data controllers. This certification scheme could be developed under Art 42(1) of the GDPR. Certification would support organisations who may prefer to rely on external expertise.

# Cluster 2 - Practical Steps

## Recommendation 2(a) - Provide guidance on when controllers should seek to anonymise data and when they can and should take an alternative approach

The ICO should provide further guidance on the types of use cases possible within the data protection regime. This could be based on case studies showing how organisations can process data without needing to meet the very high threshold for the data to be considered anonymous. The ICO should then provide guidance to help organisations identify when anonymisation may be appropriate, and when an alternative route should be taken.

### Rationale
The GDPR allows data controllers to lawfully process data on different legal bases. However, some organisations appear to believe incorrectly that they need to anonymise data in order to achieve their use case. This is fuelled both by a lack of understanding of the trade-offs involved in anonymisation (making it unsuitable for some use cases) and a lack of clarity around what is possible within the remit of the data protection regime e.g. under the legitimate or public interest legal bases.

To use these bases organisations need to demonstrate they have applied appropriate controls to protect the data subjects. However, organisations do not have a clear framework to structure their decisions on whether anonymisation is appropriate in their proposed use case or whether they can achieve an acceptable level of privacy risk via other means (for instance, de-identification and environmental controls demonstrating the risk is sufficiently low to allow

---

14   National Cyber Security Centre, Penetration Testing, published Aug 2017
15   ONS, Guidance on intruder testing, accessed Nov 2019

PRIVITAR
Research & Policy

processing on the basis of a legitimate interest). This may be inhibiting data use and not leading to the best outcomes for data controllers or data subjects.

## Recommendation 2(b) - Incentivise compliance through robust enforcement

The ICO should consider targeting regulatory activity at data controllers purporting to work with 'anonymous' or 'anonymised' data, to clarify issues around anonymisation and incentivise best practices and compliance.

### Rationale

The ICO's Regulatory Action policy sets out the objectives of regulatory action and provides a non-exhaustive list of criteria which the ICO will consider when deciding whether to take regulatory action. That document makes clear that the ICO takes action selectively including where that action is in the public interest, i.e. where it provides a deterrent or clarifies an issue in dispute. We believe that this is such an area as without the threat of enforcement data controllers may not adequately anonymise data, posing significant risk to data subjects, and seeding distrust in the efficacy of anonymisation overall. Enforcement will also provide case studies to help organisations understand the ICO's approach.

## Recommendation 2(c) - Set out the vulnerabilities associated with hashing in some use cases

The ICO should set out the privacy vulnerabilities associated with hashing and recommend that hashing not be used as a means of masking identifiers. If the ICO were minded to recommend an alternative, we suggest tokenisation.

### Rationale

Hashing is commonly used as a method for masking direct identifiers. However, in certain use cases, it will not be an acceptable privacy solution. The New York taxi data release illustrates the problem:

> In response to a freedom of information request, the New York Taxi & Limousine Commission (NYTLC) released a complete set of historical trip and fare logs from its taxis, covering over 173 million trips[16]. The NYTLC attempted to anonymise the data by hashing the license numbers.
>
> The license numbers conform to a specific pattern (either one number, one letter two numbers e.g. 5X55, or two letters, three numbers e.g. XX555, or three letters, three numbers e.g. XXX555). Meaning that the total possible number of license numbers is around 22 million.
>
> The total number of possible values is called the input domain. As there are a known and limited number of secure hashing functions, a privacy researcher was able to calculate the hashes of all possible license numbers in about 2 minutes. He used this lookup table, with rented cloud compute power, to re-identify the entire dataset within about an hour[17].

---

[16]  Chris Wong, FOILing NYC's taxi trip data, accessed Nov 2019

[17] Vijay Pandurangan, On Taxis and Rainbows: Lessons from NYC's improperly anonymized taxi logs, accessed Nov 2019

**PRIVITAR**
Research & Policy

Other common direct identifiers may also have small input domains and well-known patterns and will also be vulnerable to this kind of attack. Credit card numbers are always 16 digits long, with the first 4 digits corresponding to the issuing bank. NHS or National Insurance numbers also follow specific formats with limited number of possible values, as such hashing is not an effective way of protecting privacy for these types of identifiers.

Some organisations also use salting (i.e. appending a random string to the value before hashing it) to defend against the type of attack described above. However, this introduces two new elements of complexity. First, the "salt" value would have to be protected. Second, the "salt" would need to be shared if the data controller needed to create or maintain consistent hashes, e.g. in a use case requiring linkage on a hashed identifier. Both introduce vulnerability and potential usability issues which can be avoided by using other de-identification methods.

In contrast to hashing, tokenisation is nondeterministic, meaning that the output does not depend on the input. This feature means that tokenisation is not vulnerable to dictionary attacks, such as the one used in the NYC taxi example above.

# Recommendation 2(d) - Establish a mechanism to ensure regular consultation with stakeholders

The ICO should establish a forum for regular dialogue with stakeholders. This will ensure that guidance remains current and adapted to the needs of data controllers. The dialogue should include those working on PETs and anonymisation, such as Privitar.

## Rationale
The range and complexity of the preceding recommendations show that continued engagement between the ICO and stakeholders in the data protection community is key. We therefore recommend that the ICO establish a formal mechanism to promote consultation, collaboration and sharing ideas across sectors. This could mirror the format of the AI Council, the panels advising the Financial Conduct Authority or the Sustainable Development Advisory Group helping to guide OFGEM.

# Recommendation 2(e) - Address inconsistencies across Europe

In issuing updated guidance on anonymisation the ICO has an opportunity to drive consistency across Europe and beyond which we believe it should take. We have summarised the major inconsistencies in Annex A.

## Rationale
As we note in the introduction, clear guidance will enable organisations to unlock value in the data they collect and hold. The ICO has an opportunity to position itself as a leading regulator in Europe, in line with its ambition to be an upstream regulator and to drive privacy-compliant innovation.

**PRIVITAR**
Research & Policy

A lack of consistency across jurisdictions creates complexity for organisations and may have a chilling effect on data use. Once the ICO has adopted guidance on anonymisation, we would support outreach to other regulators in Europe and beyond to persuade them to adopt a consistent approach.

# Annex A - Examples of inconsistent interpretation

European guidance on anonymisation outside the UK is scarce. Examples include:

- The Spanish Agencia Española de Protección de Datos's (AEPD) 2016 guidance on anonymisation, which is no longer available on their website. The AEPD has more recently issued specific guidance on *k*-anonymity and hashing.

- CNIL relied on the A29WP definition (singling out, linkability and inference) in rejecting an application by JCDecaux to count pedestrians walking past a billboard using their Wi-Fi MAC addresses[18].

- The Irish Data Protection Commission (DPC) 2019 guidance. Which attempts to reconcile the A29WP opinion (below) and the ICO's 2012 Code. This is problematic because the divergences between the two are not fully addressed by the Irish DPC's document.

- The A29WP opinion, which is an interpretation of the 1995 Directive not the GDPR, and issued before the European Court of Justice (CJEU) decision in *Breyer* and has not been adopted by the European Data Protection Board (EDPB).

The main points of divergence:

1. The basis of the risk assessment and whether it considers only the data or the data in an environment. The A29WP opinion does not consider environmental controls and refers only to techniques applied to the data such as randomisation or generalisation. But the CJEU in *Breyer* appears to consider the legal context, which is one environmental factor and so indicates it cannot be a consideration of solely the data. But what other environmental factors can and should be considered? The basis of the risk assessment links to a consideration of "all objective factors", which we address in Recommendation 1(c).

2. The requirement to delete the original data (the input). We believe that the requirement to consider "all parties", and looking at anonymisation only as a property of the data and not the data in an environment, results in the view that the input must be deleted for the output to be anonymous. A European Commission impact assessment[19] of the GDPR found diverging views between Member States in relation to whether the ability of the controller to identify an individual should be taken into account. We believe that the controller retaining access to the original data does not pose a significant additional privacy risk to the data subject. Requiring controllers to delete the input is highly restrictive on an organisation's ability to anonymise data, without providing a significant benefit.

---

[18] CNIL, Délibération n° 2015-255 du 16 juillet 2015, accessed Dec 2019
[19] European Commission, Impact Assessment SEC(2012) 72 final, published Jan 2012

**PRIVITAR**
Research & Policy

3. Definition of "identifiable" and threshold for re-identification. We noted in Recommendation 1(a) that some European regulators, e.g. CNIL, have in the past appeared to require that re-identification be impossible, whilst others take a risk based approach. As stated in Recommendation 1(a) an absolute definition of re-identification being impossible may not be plausible. The CJEU in *Breyer* referred to re-identification being "practically impossible". We consider the threshold further under Recommendation 1(b).

4. There is the potential for inconsistency with regards to whether 'means reasonably likely' includes breaking the law. In Breyer the court indicated that means are not reasonably likely " ...if the identification of the data subject was prohibited by law...". This could have a perverse effect if combined with Section 171 of the DPA 2018, which makes re-identification of de-identified data an offence, but defines de-identification in such a way as to include pseudonymous data. This could be taken to mean that all a controller need do in the UK to anonymise data is pseudonymise it and then state that no one has permission to re-identify it, making doing so an offence and the data anonymous. This, which would massively weaken the anonymisation standard, we assume was not the intention of the court or the drafters of the DPA 2018.

   More broadly the idea that attackers will not resort to criminality seems too general. Some attackers will resort to some kinds of criminality for some datasets. Criminals wanting to use data to commit fraud, a reasonably likely potential attacker for some datasets, may purchase data illegally online to use in a linkage attack. Certainly the re-identification ban is unlikely to deter them if their end goal is fraud, a more severely punished offence. However, that doesn't mean that the attacker will resort to burglary, which would require physical proximity and much higher risk and effort on their behalf. A better approach might be to consider what types of criminality each attacker may be reasonably likely to resort to (see Recommendation 1(b)).

5. Some jurisdictions beyond the EU have taken a more prescriptive approach to anonymisation. For instance, in Korea, the determination is based on an assessment of the *k*-anonymity of the data set. In the US, data can be removed from the scope of HIPAA if 18 types of identifiers listed in the legislation are removed. For the reasons set out in the introduction, we do not believe that a prescriptive approach is suitable.

PRIVITAR
Research & Policy