



**NIST Special Publication
NIST SP 800-188 3pd**

De-Identifying Government Data Sets

Third Public Draft

Simson Garfinkel
Phyllis Singer
Joseph Near
Aref N. Dajani
Barbara Guttman

This publication is available free of charge from:
<https://doi.org/10.6028/NIST.SP.800-188.3pd>

NIST Special Publication
NIST SP 800-188 3pd

De-Identifying Government Data Sets

Third Public Draft

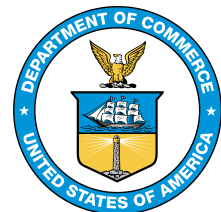
Simson Garfinkel
Barbara Guttman
Software Quality Group
Software and Systems Division

Joseph Near
Department of Computer Science
University of Vermont

Aref N. Dajani
Phyllis Singer
Center for Enterprise Dissemination
US Census Bureau

This publication is available free of charge from:
<https://doi.org/10.6028/NIST.SP.800-188.3pd>

November 2022



US Department of Commerce
Gina M. Raimondo, Secretary

National Institute of Standards and Technology
Laurie E. Locascio, NIST Director and Under Secretary of Commerce for Standards and Technology

Certain commercial entities, equipment, or materials may be identified in this document in order to describe an experimental procedure or concept adequately. Such identification is not intended to imply recommendation or endorsement by the National Institute of Standards and Technology, nor is it intended to imply that the entities, materials, or equipment are necessarily the best available for the purpose.

There may be references in this publication to other publications currently under development by NIST in accordance with its assigned statutory responsibilities. The information in this publication, including concepts and methodologies, may be used by federal agencies even before the completion of such companion publications. Thus, until each publication is completed, current requirements, guidelines, and procedures, where they exist, remain operative. For planning and transition purposes, federal agencies may wish to closely follow the development of these new publications by NIST.

Organizations are encouraged to review all draft publications during public comment periods and provide feedback to NIST. Many NIST cybersecurity publications, other than the ones noted above, are available at <https://csrc.nist.gov/publications>.

This document is presented with the hope that its content may be of interest to the general privacy community. The views in this document are those of the authors, and do not represent those of the US Census Bureau.

Authority

This publication has been developed by NIST in accordance with its statutory responsibilities under the Federal Information Security Modernization Act (FISMA) of 2014, 44 U.S.C. § 3551 et seq., Public Law (P.L.) 113-283. NIST is responsible for developing information security standards and guidelines, including minimum requirements for federal information systems, but such standards and guidelines shall not apply to national security systems without the express approval of appropriate federal officials exercising policy authority over such systems. This guideline is consistent with the requirements of the Office of Management and Budget (OMB) Circular A-130.

Nothing in this publication should be taken to contradict the standards and guidelines made mandatory and binding on federal agencies by the Secretary of Commerce under statutory authority. Nor should these guidelines be interpreted as altering or superseding the existing authorities of the Secretary of Commerce, Director of the OMB, or any other federal official. This publication may be used by nongovernmental organizations on a voluntary basis and is not subject to copyright in the United States. Attribution would, however, be appreciated by NIST.

NIST Technical Series Policies

[Copyright, Fair Use, and Licensing Statements](#)

[NIST Technical Series Publication Identifier Syntax](#)

Publication History

Approved by the NIST Editorial Review Board on YYYY-MM-DD [will be added upon final publication]

How to cite this NIST Technical Series Publication:

Garfinkel S, Guttman B, Near J, Dajani AN, Singer P (2022) De-Identifying Government Data Sets. (National Institute of Standards and Technology, Gaithersburg, MD), NIST Special Publication (SP) NIST SP 800-188 3pd. <https://doi.org/10.6028/NIST.SP.800-188.3pd>

Author ORCID iDs

Simson Garfinkel: 0000-0003-1294-2831

Joseph Near: 0000-0002-3203-3742

Aref N. Dajani: 0000-0003-0361-5409

Phyllis Singer: 0000-0002-8885-7273

Public Comment Period

83 November 15, 2022 – January 15, 2023

84 **Submit Comments**

85 sp800-188-draft@nist.gov

86 National Institute of Standards and Technology

87 Attn: Software and Systems Division, Information Technology Laboratory

88 100 Bureau Drive (Mail Stop 8970) Gaithersburg, MD 20899-8970

89 **All comments are subject to release under the Freedom of Information Act (FOIA).**

Abstract

De-identification is a process that is applied to a dataset with the goal of preventing or limiting informational risks to individuals, protected groups, and establishments while still allowing for meaningful statistical analysis. Government agencies can use de-identification to reduce the privacy risk associated with collecting, processing, archiving, distributing, or publishing government data. Previously, NISTIR 8053, *De-Identification of Personal Information* [51], provided a survey of de-identification and re-identification techniques. This document provides specific guidance to government agencies that wish to use de-identification. Before using de-identification, agencies should evaluate their goals for using de-identification and the potential risks that de-identification might create. Agencies should decide upon a de-identification release model, such as publishing de-identified data, publishing synthetic data based on identified data, or providing a query interface that incorporates de-identification. Agencies can create a Disclosure Review Board to oversee the process of de-identification. They can also adopt a de-identification standard with measurable performance levels and perform re-identification studies to gauge the risk associated with de-identification. Several specific techniques for de-identification are available, including de-identification by removing identifiers and transforming quasi-identifiers and the use of formal privacy models. People performing de-identification generally use special-purpose software tools to perform the data manipulation and calculate the likely risk of re-identification. However, not all tools that merely mask personal information provide sufficient functionality for performing de-identification. This document also includes an extensive list of references, a glossary, and a list of specific de-identification tools, which is only included to convey the range of tools currently available and is not intended to imply a recommendation or endorsement by NIST.

Keywords

data life cycle; de-identification; differential privacy; direct identifiers; Disclosure Review Board; the five safes; k -anonymity; privacy; pseudonymization; quasi-identifiers; re-identification; synthetic data.

Reports on Computer Systems Technology

The Information Technology Laboratory (ITL) at the National Institute of Standards and Technology (NIST) promotes the U.S. economy and public welfare by providing technical leadership for the Nation's measurement and standards infrastructure. ITL develops tests, test methods, reference data, proof of concept implementations, and technical analyses to advance the development and productive use of information technology. ITL's responsibilities include the development of management, administrative, technical, and physical standards and guidelines for the cost-effective security and privacy of other than national security-related information in federal information systems. The Special Publication 800-

127 series reports on ITL's research, guidelines, and outreach efforts in information system
128 security, and its collaborative activities with industry, government, and academic organiza-
129 tions.

130 **Call for Patent Claims**

131 This public review includes a call for information on essential patent claims (claims whose
132 use would be required for compliance with the guidance or requirements in this Information
133 Technology Laboratory (ITL) draft publication). Such guidance and/or requirements may
134 be directly stated in this ITL Publication or by reference to another publication. This call
135 also includes disclosure, where known, of the existence of pending U.S. or foreign patent
136 applications relating to this ITL draft publication and of any relevant unexpired U.S. or
137 foreign patents.

138 ITL may require from the patent holder, or a party authorized to make assurances on its
139 behalf, in written or electronic form, either:

- 140 1. assurance in the form of a general disclaimer to the effect that such party does not
141 hold and does not currently intend holding any essential patent claim(s); or
- 142 2. assurance that a license to such essential patent claim(s) will be made available to ap-
143 plicants desiring to utilize the license for the purpose of complying with the guidance
144 or requirements in this ITL draft publication either:
 - 145 (a) under reasonable terms and conditions that are demonstrably free of any unfair
146 discrimination; or
 - 147 (b) without compensation and under reasonable terms and conditions that are demon-
148 strably free of any unfair discrimination.

149 Such assurance shall indicate that the patent holder (or third party authorized to make assur-
150 ances on its behalf) will include in any documents transferring ownership of patents subject
151 to the assurance, provisions sufficient to ensure that the commitments in the assurance are
152 binding on the transferee, and that the transferee will similarly include appropriate provi-
153 sions in the event of future transfers with the goal of binding each successor-in-interest.

154 The assurance shall also indicate that it is intended to be binding on successors-in-interest
155 regardless of whether such provisions are included in the relevant transfer documents.

156 Such statements should be addressed to: sp800-188-draft@nist.gov

159 **Table of Contents**

160	Executive Summary	1
161	1. Introduction	3
162	1.1. Document Purpose and Scope	7
163	1.2. Intended Audience	7
164	1.3. Organization	7
165	2. Introducing De-Identification	8
166	2.1. Historical Context	8
167	2.2. Terminology	10
168	3. Governance and Management of Data De-Identification	17
169	3.1. Identifying Goals and Intended Uses of De-Identification	17
170	3.2. Evaluating Risks that Arise from De-Identified Data Releases	18
171	3.2.1. Probability of Re-Identification	19
172	3.2.2. Adverse Impacts of Re-Identification	22
173	3.2.3. Impacts Other Than Re-Identification	23
174	3.2.4. Remediation	24
175	3.3. Data Life Cycle	24
176	3.4. Data-Sharing Models	29
177	3.5. The Five Safes	30
178	3.6. Disclosure Review Boards	31
179	3.7. De-Identification Standards	36
180	3.7.1. Benefits of Standards	36
181	3.7.2. Prescriptive De-Identification Standards	36
182	3.7.3. Performance-Based De-Identification Standards	37
183	3.8. Education, Training, and Research	38
184	3.9. Defense in Depth	38
185	3.9.1. Encryption and Access Control	38
186	3.9.2. Secure Computation	38
187	3.9.3. Trusted Execution Environments	39
188	3.9.4. Physical Enclaves	39
189	4. Technical Steps for Data De-Identification	40

190	4.1. Determine the Privacy, Data Usability, and Access Objectives	40
191	4.2. Conducting a Data Survey	41
192	4.3. De-Identification by Removing Identifiers and Transforming Quasi-Identifiers	43
193	4.3.1. Removing or Transforming of Direct Identifiers	44
194	4.3.2. Special Security Note Regarding the Encryption or Hashing of Di-	
195	rect Identifiers	46
196	4.3.3. De-Identifying Numeric Quasi-Identifiers	46
197	4.3.4. De-Identifying Dates	48
198	4.3.5. De-Identifying Geographical Locations and Geolocation Data . . .	49
199	4.3.6. De-Identifying Genomic Information	49
200	4.3.7. De-Identifying Text Narratives and Qualitative Information	51
201	4.3.8. Challenges Posed by Aggregation Techniques	51
202	4.3.9. Challenges Posed by High-Dimensional Data	52
203	4.3.10. Challenges Posed by Linked Data	52
204	4.3.11. Challenges Posed by Composition	52
205	4.3.12. Potential Failures of De-Identification	53
206	4.3.13. Post-Release Monitoring	54
207	4.4. Synthetic Data	54
208	4.4.1. Partially Synthetic Data	55
209	4.4.2. Test Data	56
210	4.4.3. Fully Synthetic Data	56
211	4.4.4. Synthetic Data with Validation	58
212	4.4.5. Synthetic Data and Open Data Policy	58
213	4.4.6. Creating a Synthetic Dataset with Differential Privacy	58
214	4.5. De-Identifying with an Interactive Query Interface	59
215	4.6. Validating a De-Identified Dataset	60
216	4.6.1. Validating Data Usefulness	60
217	4.6.2. Validating Privacy Protection	60
218	4.6.3. Re-Identification Studies	61
219	5. Software Requirements, Evaluation, and Validation	63
220	5.1. Evaluating Privacy-Preserving Techniques	63
221	5.2. De-Identification Tools	64

222	5.2.1. De-Identification Tool Features	64
223	5.2.2. Data Provenance and File Formats	64
224	5.2.3. Data Masking Tools	64
225	5.3. Evaluating De-Identification Software	65
226	5.4. Evaluating Data Accuracy	65
227	6. Conclusion	66
228	References	80
229	Appendix A. Standards	81
230	A.1. NIST Publications	81
231	A.2. Other U.S. Government Publications	82
232	Selected Publications by Other Governments	83
233	Reports and Books	83
234	How-To Articles	84
235	Appendix B. List of Symbols, Abbreviations, and Acronyms	86
236	Appendix C. Glossary	89

237 List of Tables

238	Table 1. Reading levels at a hypothetical school, as measured by entrance exami-	
239	nations, reported at the start of the school year on October 1.	51
240	Table 2. Reading levels at a hypothetical school, as measured by entrance exami-	
241	nations, reported one month into the school year on November 1 after a	
242	new student has transferred to the school.	51
243	Table 3. Adjectives used for describing data in data releases.	55

244 List of Figures

245	Fig. 1. The data life cycle as described by Michener et al. [80]	25
246	Fig. 2. Chisholm's view of the data life cycle is a linear process with a branching	
247	point after data usage [25]	25
248	Fig. 3. Altman's "modern approach to privacy-aware government data releases" [75]	26
249	Fig. 4. Altman's conceptual diagram of the relationship between post-transformation	
250	identifiability, level of expected harm, and suitability of selected privacy	
251	controls for a data release [75]	27
252	Fig. 5. Advice for Practitioners: A Summary	68

Acknowledgments

The authors wish to thank the U.S. Census Bureau for its help in researching and preparing this publication, with specific thanks to John Abowd, Ron Jarmin, Christa Jones, and Laura McKenna. The authors would also like to thank Luk Arbuckle, Andrew Baker, Daniel Barth-Jones, Christi Dant, Khaled El Emam, Robert Gellman, Tom Krenzke, Bradley Malin, Kevin Mangold, John Moehrke, Linda Sanchez, Denise Sturdy, and Chris Traver for providing comments on previous drafts and their valuable insights, all of which were helpful in creating this publication.

The authors also wish to thank several organizations that provided useful comments on previous drafts of this publication: the Defense Contract Management Agency (DCMA) Information Assurance Directorate; the Office of Chief Privacy Officer within the U.S. Department of Education; the Office of Planning, Research, and Evaluation (OPRE) within the Administration for Children and Families at the U.S. Department of Health and Human Services; the Millennium Challenge Corporation (MCC) Department of Policy and Evaluation; Integrating the Healthcare Enterprise (IHE), an ANSI-accredited standards organization focused on healthcare standards; and the Privacy Tools project at Harvard University (including Micah Altman, Stephen Chong, Kobbi Nissim, David O'Brien, Salil Vadhan, and Alexandra Wood).

Author Contributions

Simson Garfinkel: Conceptualization, Supervision, Writing (original draft preparation); **Joseph Near:** Writing (original draft preparation); **Aref N. Dajani:** Writing (original draft preparation of the section on reidentification studies); **Phyllis Singer:** Writing (original draft preparation of the section on reidentification studies).

Executive Summary

Every federal agency creates and maintains internal datasets that are vital for fulfilling its mission. The Foundation for Evidence-based Policymaking Act of 2018 [2] mandates that agencies also collect and publish their government data in open, machine-readable formats, when it is appropriate to do so. Agencies can use de-identification to make government datasets available while protecting the privacy of the individuals whose data are contained within those datasets.

Many Government documents use the phrase *personally identifiable information* (PII) to describe private information that can be linked to an individual [62, 79], although there are a variety of definitions for PII. As a result, it is possible to have information that *singles out* individuals but that does not meet a specific definition of PII. This document therefore presents ways of removing or altering information that can identify individuals that go beyond merely removing PII.

For decades, de-identification based on simply removing of identifying information was thought to be sufficient to prevent the re-identification of individuals in large datasets. Since the mid 1990s, a growing body of research has demonstrated the reverse, resulting in new privacy attacks capable of re-identifying individuals in “de-identified” data releases. For several years the goals of such attacks appeared to be the embarrassment of the publishing agency and achieving academic distinction for the privacy researcher [50]. More recently, as high-resolution de-identified geolocation data has become commercially available, re-identification techniques have been used by journalists and activists [100, 140, 70] with the goal of learning confidential information.

These attacks have become more sophisticated in recent years with the availability of geolocation data, highlighting the deficiencies in traditional

Formal models of privacy, like k -anonymity [122] and differential privacy, [39] use mathematically rigorous approaches that are designed to allow for the controlled use of confidential data while minimizing the privacy loss suffered by the data subjects. Because there is an inherent trade-off between the accuracy of published data and the amount of privacy protection afforded to data subjects, most formal methods have some kind of parameter that can be adjusted to control the “privacy cost” of a particular data release. Informally, a data release with a low privacy cost causes little additional privacy risk to the participants, while a higher privacy cost results in more privacy risk. When they are available, formal privacy methods should be preferred over informal, *ad hoc* methods.

Decisions and practices regarding the de-identification and release of government data can be integral to the mission and proper functioning of a government agency. As such, an agency’s leadership should manage these activities in a way that assures performance and results in a manner that is consistent with the agency’s mission and legal authority. One way that agencies can manage this risk is by creating a formal Disclosure Review Board (DRB) that consists of legal and technical privacy experts, stakeholders within the organization,

and representatives of the organization's leadership. The DRB evaluated applications for data release that describe the confidential data, the techniques that will be used to minimize the risk of disclosure, the resulting protected data, and how the effectiveness of those techniques will be evaluated.

Establishing a DRB may seem like an expensive and complicated administrative undertaking for some agencies. However, a properly constituted DRB and the development of consistent procedures regarding data release should enable agencies to lower the risks associated with each data release, which is likely to save agency resources in the long term.

Agencies can create or adopt standards to guide those performing de-identification, and regarding the accuracy of de-identified data. If accuracy goals exist, then techniques such as differential privacy can be used to make the data sufficiently accurate for the intended purpose but not unnecessarily more accurate, which can limit the amount of privacy loss. However, agencies must carefully choose and implement accuracy requirements. If data accuracy and privacy goals cannot be well-maintained, then releases of data that are not sufficiently accurate can result in incorrect scientific conclusions and policy decisions.

Agencies should consider performing de-identification with trained individuals using software specifically designed for the purpose. While it is possible to perform de-identification with off-the-shelf software like a commercial spreadsheet or financial planning program, such programs typically lack the key functions required for proper de-identification. As a result, they may encourage the use of simplistic de-identification methods, such as deleting sensitive columns and manually searching and removing data that appears sensitive. This may result in a dataset that appears de-identified but that still contain significant disclosure risks.

Finally, different countries have different standards and policies regarding the definition and use of de-identified data. Information that is regarded as de-identified in one jurisdiction may be regarded as being identifiable in another.

1. Introduction

The U.S. Government collects, maintains, and uses many kinds of datasets. Every federal agency creates and maintains internal datasets that are vital for fulfilling its mission, such as delivering services to taxpayers or ensuring regulatory compliance. There are also 13 principal federal statistical agencies, three recognized statistical units, and over 100 other federal statistical programs that collect, compile, process, analyze, and distribute information for statistical purposes [126, 92].

Government programs collect information from individuals and organizations for taxation, public benefits, public health, licensing, employment, censuses, and the production of official statistics. While privacy is integral, many individuals and organizations that provide information to the Government do not typically have the right to opt-out of such requests. For example, people and establishments in the United States are required by law to respond to mandatory U.S. Census Bureau surveys.

Agencies make many of their datasets available to the public. The U.S. Government publishes data to promote commerce, scientific research, and public transparency. Many datasets contain some data elements that should not be made public, and it is necessary to remove such information before making the rest of the dataset available. Some datasets are so sensitive that they cannot be made publicly available at all but can be available on a limited basis to qualified, vetted researchers in protected enclaves. In some cases, agencies may also elect to release summary statistics of sensitive data or create synthetic datasets that resemble the original data but that have a lower disclosure risk [8].

There is frequent tension between the goals of privacy protection and the release of useful data to the public. One way that the Government attempts to resolve this tension is with an official promise of confidentiality to individuals and organizations regarding the information that they provide [102]. A bedrock principle of official statistical programs is that data provided to the Government should generally remain confidential and not be used in a way that could harm the individual or the organization providing the data. One justification for this principle is that it helps to ensure high data accuracy. If data providers did not feel that the information they provide would remain confidential, they might not be willing to provide information that is accurate.

Other information is created by the Government as a consequence of providing government services. This information – sometimes called *administrative data* – is also increasingly being used and made available for statistical purposes and must be protected.

In 2018, the U.S. Congress passed three laws that significantly increased the need for expertise regarding privacy-preserving data analysis and data publishing techniques, such as de-identification:

1. **The Foundations for Evidence-Based Policymaking Act of 2018** [2], commonly called the Evidence Act, requires federal agencies to track all of their data in data

inventories, report public datasets to <https://data.gov>, perform systematic evidence-making and evaluation activities, and engage in capacity-building so that the federal workforce can meet the requirements of data-centric, evidence-based operations. The Evidence Act is based on the findings of the U.S. Commission on Evidence-Based Policymaking [27] and is implemented in part by OMB Memorandum M-19-23 [139].

The Evidence Act contains specific guidance requiring that agencies publishing data take into account “(A) risks and restrictions related to the disclosure of personally identifiable information, including the risk that an individual data asset in isolation does not pose a privacy or confidentiality risk but when combined with other available information may pose such a risk;” and “(B) security considerations, including the risk that information in an individual data asset in isolation does not pose a security risk but when combined with other available information may pose such a risk” [2].

2. **The Open Government Data Act**, which was passed as part of the Evidence Act, requires that the U.S. Government publish data in machine-readable, open, non-proprietary formats when possible. This act largely codified presidential Executive Order 13642 of May 9, 2013, “Making Open and Machine Readable the New Default for Government Information” [88] and its implementation in OMB Memorandum M-13-13 [18].
3. **The Geospatial Data Act of 2018**, which requires that government agencies make inventories of their geospatial data and that public geospatial data be registered on the U.S. Government’s public geospatial platform, <https://www.geoplatform.gov/>.

Other laws, regulations, and policies that govern the release of statistics and data to the public enshrine this principle of confidentiality. For example:

- **The Confidential Information Protection and Statistical Efficiency Act of 2002** states, “data or information acquired by an agency under a pledge of confidentiality for exclusively statistical purposes shall not be disclosed by an agency in identifiable form for any use other than an exclusively statistical purpose, except with the informed consent of the respondent.” [126, §512 (b)(1)] Commonly called CIPSEA, the act further requires that federal statistical agencies “establish appropriate administrative, technical, and physical safeguards to ensure the security and confidentiality of records and to protect against any anticipated threats or hazards to their security or integrity which could result in substantial harm, embarrassment, inconvenience, or unfairness to any individual on whom information is maintained.”
- **US Code Title 13, Section 9** governs the confidentiality of information provided to the Census Bureau and prohibits “any publication whereby the data furnished by any particular establishment or individual under this title can be identified” [130].
- **US Code Title 26, Section 6103** governs the confidentiality of information provided to the U.S. Government on tax returns and other return information. These rules are

now spelled out in IRS Publication 1075, “Tax Information Security Guidelines for Federal, State and Local Agencies,” published by the IRS Office of Safeguards [93].

- **The Privacy Act of 1974** covers the release of personal information of U.S. citizens and Lawful Permanent Residents by the Government. The Act recognizes that the disclosure of records for statistical purposes is acceptable if the data are not “individually identifiable” [103, at a(b)(5)].

Minimizing privacy risk is not an absolute goal of federal laws and regulations. Guidance from the U.S. Department of Health and Human Services (HHS) on the Health Insurance Portability and Accountability Act (HIPAA) de-identification standards notes that “[b]oth methods [the safe harbor and expert determination methods for de-identification], even when properly applied, yield de-identified data that retains some risk of identification. Although the risk is very small, it is not zero, and there is a possibility that de-identified data could be linked back to the identity of the patient to which it corresponds” [136].

U.S. law also balances privacy risk with other factors, such as transparency, accountability, and the opportunity for public good. An example of this balance is the handling of personally identifiable information collected by the Census Bureau as part of the decennial census: this information remains confidential for 72 years and is then transferred to the National Archives and Records Administration where it is released to the public [131, 5].

De-identification is a process that is applied to a dataset with the goal of preventing or limiting privacy risks to individuals, protected groups, and establishments while still allowing for the production of aggregate statistics.¹ De-identification is not a single technique, but a collection of approaches, algorithms, and tools that can be applied to different kinds of data with differing levels of effectiveness. In general, the potential risk to privacy posed by a dataset’s release decreases as more aggressive de-identification techniques are employed, but data accuracy and – in some cases – the ultimate utility of the de-identified dataset decreases as well.

Accuracy is traditionally defined as the “closeness of computations or estimates to the exact or true values that the statistics were intended to measure” [9]. The *data accuracy* of de-identified data, therefore, refers to the degree to which inferences drawn on the de-identified data will be consistent with inferences drawn on the original data. Data accuracy can be measured by the ratio of a value computed with de-identified data to the same value computed using the underlying true confidential value.

In economics, *Utility* is traditionally defined as “the satisfaction derived from consumption of a good or service” [138]. *Data utility* therefore refers to the value that data users can derive from data in general. When speaking of de-identified data, utility comes from two pub-

¹In Europe, the term *data anonymization* is frequently used as a synonym for de-identification, but the terms may have subtly different definitions in some contexts. For a more complete discussion of de-identification and data anonymization, see NISTIR 8053, *De-Identification of Personal Data* [51].

lic goods: the uses of the data and the privacy protection afforded by the de-identification process.

This document uses the phrase *data accuracy* to refer to the abstract characteristic of the data as determined by a specific, measurable statistic, whereas *data utility* refers to the benefit derived from the application of the data to a specific use. Although there has previously been a tendency within official statistical organizations to conflate these two terms, it is important to keep them distinct because they are not necessarily correlated. Data may have low accuracy because they contain errors or substantial noise, yet users may nevertheless derive high value from the data, giving the data high utility. Likewise, data that are very close to the reality of the thing being measured may have high accuracy but may be fundamentally worthless and, thus, have low utility.

In general, data accuracy decreases as more aggressive de-identification techniques are employed. Therefore, any effort that involves the release of data that contain personal information typically involves making a trade-off between identifiability and data accuracy. However, increased privacy protections do not necessarily result in decreased data utility.

Some users of de-identified data may be able to use the data to make inferences about private facts regarding the data subjects. They may even be able to re-identify the data subjects. Both of these uses undo the privacy goals of de-identification. Agencies that release data should understand what data they are releasing, what other data may already be publicly or privately available, and the risk of re-identification. Agencies should aim to make an informed decision about the fidelity of the data that they release by systematically evaluating the risks and benefits and choosing de-identification techniques and data sharing models that are tailored to their requirements. In addition, when telling individuals that their de-identified information will be released, agencies should disclose that privacy risks may remain despite de-identification.

Planning is essential for successful de-identification and data release. In a research environment, this planning should include the research design, data collection, protection of identifiers, disclosure analysis, and data-sharing strategy. In an operational environment, this planning includes a comprehensive analysis of the purpose of the data release and the expected use of the released data, the privacy-related risks, and the privacy protecting controls. Both cases should review the appropriateness of various privacy controls given the risks, intended uses, and the ways that those controls could fail.

De-identification can have significant costs, including time, labor, and data processing costs. However, when properly executed, this effort can result in data that have high value for a research community and the general public while still adequately protecting individual privacy.

1.1. Document Purpose and Scope

This document provides guidance on the selection, use, and evaluation of de-identification techniques for U.S. Government datasets. It also provides a framework that can be adapted by federal agencies to shape the governance of de-identification processes. The ultimate goal of this document is to reduce disclosure risks that might result from an intentional data release.

1.2. Intended Audience

This document is intended for use by government engineers, data scientists, privacy officers, disclosure review boards, and other officials. It is also designed to be generally informative to researchers and academics involved in the technical aspects of the de-identification of government data. While this document assumes a high-level understanding of information system security technologies, it is intended to be accessible to a wide audience.

1.3. Organization

The remainder of this publication is organized as follows:

- **Section 2, “Introducing De-Identification,”** presents a background on the science and terminology of de-identification.
- **Section 3, “Governance and Management of Data De-Identification,”** provides guidance to agencies on the establishment of or improvement to a program that makes privacy-sensitive data available to researchers and the public.
- **Section 4, “Technical Steps for Data De-Identification,”** provides specific technical guidance for performing de-identification using a variety of mathematical approaches.
- **Section 5, “Software Requirements, Evaluation, and Validation,”** provides a recommended set of features that should be in de-identification tools, which may be useful for potential purchasers or developers of such software. This section also provides information for evaluating both de-identification tools and de-identified datasets.
- **Section 6, “Conclusion,”** Section 6 is the conclusion.

Following the conclusion, this document provides a list of all publications referenced in this document, as well as an Appendix that includes standards, related NIST publications, other selected publications by the US and other governments, reports and books, and a few articles of interest. A second appendix provides a list of symbols, abbreviations and acronyms. The third appendix contains a glossary.

2. Introducing De-Identification

This document presents recommendations for de-identifying government datasets.

If the information derived from personal data remains in a de-identified dataset, the dataset might inadvertently reveal attributes related to specific individuals, specific de-identified records could be linked back to specific individuals. When this happens, the privacy protection provided by de-identification is compromised. Even if a specific individual cannot be matched to a specific data record, de-identified data can be used to improve the accuracy of inferences regarding individuals whose de-identified data are in the dataset. This so-called *inference risk* cannot be eliminated if there is any information in the de-identified data, but it can be minimized. Thus, the decision of how or whether to de-identify data should be made in conjunction with decisions over how the de-identified data will be used, shared, or released.

De-identification is especially important for government agencies, businesses, and other organizations that seek to make data available to outsiders. For example, significant medical research resulting in societal benefit is made possible by the sharing of de-identified patient information under the framework established by the HIPAA Privacy Rule, the primary U.S. regulation that provides for the privacy of medical records; billing records; enrollment, payment, and claims records; and “other records that are used, in whole or in part, by or for the covered entity to make decisions about individuals” [90]. The HIPAA Privacy Rule de-identification framework applies to both government organizations charged with protecting government datasets as well as to private sector organizations, such as health plans and health care providers.

Agencies may also be required to de-identify records when responding to a Freedom of Information Act (FOIA) [134, 133] request in a manner that is consistent with Exemption 6, which protects information about individuals in “personnel and medical files and similar files” when the disclosure of such information “would constitute a clearly unwarranted invasion of personal privacy,” and Exemption 7(C), which is limited to information compiled for law enforcement purposes and protects personal information when disclosure “could reasonably be expected to constitute an unwarranted invasion of personal privacy.” The meaning of these exemptions has been clarified by multiple cases before the US Supreme Court [105, 117, 118].

2.1. Historical Context

The modern practice of de-identification comes from three overlapping intellectual traditions.

1. For four decades, official statistical agencies have researched and investigated methods broadly termed *Statistical Disclosure Limitation* (SDL) or *Statistical Disclosure*

555 *Control* [29, 36].² Statistical agencies created these methods so that they could re-
556 lease statistical tables and *public use files* (PUF) to allow users to learn information
557 and perform original research while protecting the privacy of the individuals in the
558 dataset. SDL is widely used in contemporary statistical reporting.

559 2. In the 1990s, there was a significant increase in the release of *microdata* files for
560 public use in the form of both individual responses from surveys and administrative
561 records. Initially, these releases merely stripped obviously identifying information,
562 such as names and social security numbers (what are now called *direct identifiers*).
563 Following some releases, researchers discovered that it was possible to re-identify
564 individuals' data by triangulating with some of the remaining data (now called *quasi-*
565 *identifiers* or *indirect identifiers* [28]). The research resulted in the invention of the
566 *k*-anonymity model for protecting privacy [124, 108, 109, 123], which is reflected
567 in the Office of Civil Rights guidance on how to apply de-identification in a manner
568 consistent with the HIPAA Privacy Rule [89]. Today, variants of *k*-anonymity are
569 commonly used to allow for the sharing of medical microdata. This intellectual tra-
570 dition is typically called *de-identification*, although this document uses that term to
571 describe all three intellectual traditions.

572 3. In the 2000s, research in theoretical computer science and cryptography developed
573 the theory of *differential privacy* [40], which is based on a mathematical definition
574 of the privacy loss to an individual that results from queries on a database containing
575 that individual's personal information. Differential privacy is termed a *formal model*
576 *for privacy protection* because its definitions for privacy and privacy loss are based on
577 mathematical proofs.³ This does not mean that algorithms that implement differen-
578 tial privacy cannot result in increased privacy risk. Rather, it means that the amount
579 of privacy risk that results from the use of these algorithms can be mathematically
580 bounded. These mathematical limits on privacy risk have created considerable inter-
581 est in differential privacy in academia, commerce, and business. To date, however,
582 only a few systems that utilize differential privacy have been operationally deployed.

583 During the first decade of the 21st century, there was a growing awareness within the U.S.
584 Government about the risks that could result from the improper handling and inadvertent
585 release of personal identifying and financial information. This realization, combined with
586 a growing number of inadvertent data disclosures within the U.S. Government, resulted
587 in President George Bush signing Executive Order 13402, which established an Identity
588 Theft Task Force on May 10, 2006 [19]. One year later, the Office of Management and
589 Budget issued Memorandum M-07-16 [62], which required federal agencies to develop
590 and implement breach notification policies. As part of this effort, NIST issued Special

²A summary of the history of Statistical Disclosure Limitation can be found in *Private Lives and Public Policies: Confidentiality and Accessibility of Government Statistics* [102].

³Other formal methods for privacy include cryptographic algorithms and techniques with provably secure properties, privacy-preserving data mining, Shamir's secret sharing, and advanced database techniques. A summary of such techniques appears in [128].

Publication (SP) 800-122, *Guide to Protecting the Confidentiality of Personally Identifiable Information (PII)* [79]. These policies and documents had the specific goal of limiting the accessibility of information that could be directly used for identity theft but did not create a framework for processing government datasets so that they could be released without impacting the privacy of the data subjects.

In 2015, NIST published NISTIR 8053, *De-Identification of Personal Information* [51], which provided an overview of de-identification issues and terminology. It also summarized significant publications involving de-identification and re-identification. However, NISTIR 8053 did not make recommendations regarding the appropriateness of de-identification or specific de-identification algorithms. The following year, NIST convened a Government Data De-Identification Stakeholder’s Meeting [52].

De-identification is one of several models for allowing the controlled sharing of personal data and other kinds of sensitive data.⁴ Other models include the use of data processing enclaves, where computations are performed with confidential data using computers that are physically isolated from the outside world. That isolation might be performed with locked doors and guards, or it might be performed using silicon and encryption, as is the case with enclaves implemented on some modern microprocessors. Another approach is to use mathematical techniques – such as secure multiparty computation – so that computations can be carried out on confidential data held by multiple parties without ever bringing all of the confidential data together in a single location.

Techniques for privacy-preserving data-sharing and analysis can be layered to provide stronger protection than any single technique would provide in isolation. Such complementary models are discussed in Section 3.4. For a more complete description of data-sharing models, privacy-preserving data publishing, and privacy-preserving data mining, see NISTIR 8053.

Many of the techniques discussed in this publication (e.g., fully synthetic data and differential privacy) have limited use within the Federal Government due to cost, time constraints, and the sophistication required of practitioners. However, these techniques are likely to see increased use as agencies seek to make datasets that include identifying information available.

2.2. Terminology

While each of the de-identification traditions has developed its own terminology and mathematical models, they share many underlying goals and concepts. Where terminology differs, this document relies on the terminology developed in previous documents by the U.S. Government and standards organizations.

⁴For information on characterizing the sensitivity of information, see NIST SP 800 Volume I, Revision 1 [119].

De-identification is a process that is applied to a dataset with the goal of preventing or limiting informational risks to individuals, protected groups, and establishments while still allowing for the production of aggregate statistics.⁵ *De-identification* takes an *original dataset* and produces *de-identified data*.

Re-identification is the general term for any process that restores the association between a set of de-identified data and the data subject. *Re-identification* is not the only way that de-identification techniques can fail to protect privacy. Improperly de-identified information can also be used to infer private facts about individuals that were thought to have been protected.

Re-identification risk is the likelihood that a third party can re-identify data subjects in a de-identified dataset. *Re-identification risk* is typically a function of the adverse impacts that would arise if the re-identification were to occur and the likelihood of occurrence. *Re-identification risk* is a specific form of privacy risk.

Redaction is the removal of information from a document or dataset for legal or security purposes. Also known as *suppression*, redaction is a kind of de-identifying technique that relies on the removal of information. In general, redaction alone is not sufficient to provide formal privacy guarantees, such as differential privacy. Redaction may also reduce the data accuracy of the dataset since the use of selective redaction may result in the introduction of non-ignorable bias.

Anonymization is a “process that removes the association between the identifying dataset and the data subject” [66]. This term is reserved for de-identification processes that cannot be reversed.

Some authors use the terms *de-identification* and *anonymization* interchangeably. In some contexts, the term *anonymization* is used to describe the destruction of a table that maps pseudonyms to real identifiers.⁶ Both of these uses are potentially misleading, as many de-identification procedures can be readily reversed if a dataset is discovered that maps a unique attribute or combination of attributes to identities. For example, a medical dataset may contain a list of names, medical identifiers, the rooms where a patient was seen, the time that the patient was seen, and the results of a medical test. Such a dataset could be de-identified by removing the name and medical identification numbers. However, the dataset of medical test results should not be considered anonymized because the tests can be re-identified if the dataset is joined with a second dataset of room numbers, times, and

⁵ISO/TS 25237:2008 defines de-identification as the “general term for any process of removing the association between a set of identifying data and the data subject.” [66]. This document intentionally adopts a broader definition for de-identification that allows for noise-introducing techniques, such as differential privacy and the creation of synthetic datasets that are based on privacy-preserving models.

⁶For example, “Anonymization is a step subsequent to de-identification that involves destroying all links between the de-identified datasets and the original datasets. The key code that was used to generate the new identification code number from the original is irreversibly destroyed (i.e., destroying the link between the two code numbers)” [127].

names. Since it is not possible to know whether such an auxiliary dataset exists, this publication recommends avoiding the word *anonymization* and using the word *de-identification* instead.

Because of the inconsistencies in the use and definitions of the word “anonymization,” this document avoids the term except in this section and in the titles of some references. Instead, it uses the term “de-identification” with the understanding that sometimes de-identified information can be re-identified, and sometimes it cannot.⁷

Pseudonymization is a “particular type of [de-identification]⁸ that both removes the association with a data subject and adds an association between a particular set of characteristics relating to the data subject and one or more pseudonyms” [66]. The term *coded* is frequently used in healthcare settings to describe data that has been pseudonymized. Pseudonymization is commonly used so that multiple observations of an individual over time can be matched and so that an individual can be re-identified if there is a policy reason to do so. Although pseudonymous data are typically re-identified by consulting a key that may be highly protected, the existence of the pseudonym identifiers frequently increases the risk of re-identification through other means.

Many U.S. Government documents use the phrase *personally identifiable information* (PII) to describe private information that can be linked to an individual [62, 79], although there are a variety of definitions for PII in various laws, regulations, and agency guidance documents. Because of these differing definitions, it is possible to have information that *singles out* individuals but that does not meet a specific definition of PII. An added complication is that some documents use the term PII to denote any information that is attributable to individuals or information that is uniquely attributable to a specific individual, while others use the term strictly for data that are directly identifying.

This document avoids the term *personally identifiable information*. Instead, it uses the phrases *personal data* or *personal information* to denote information related to individuals and *identifying information* for “information that can be used to distinguish or trace an individual’s identity, such as their name, social security number, biometric records, etc., alone, or when combined with other personal or identifying information which is linked or linkable to a specific individual, such as date and place of birth, mother’s maiden name, etc.” [62]. Under this definition, identifying information is personal information, but personal information is not necessarily identifying information.

Non-public personal information is used to describe personal information that is in a dataset that is not publicly available. Non-public personal information is not necessarily identifying.

⁷Thus, where other references (e.g. [104]) might use the term *anonymized file* or *anonymized dataset* to describe a dataset that has been de-identified, this publication will use the terms *de-identified file* and *de-identified dataset* since the term *de-identified* is descriptive while the term *anonymized* is aspirational.

⁸Here, the word *anonymization* in the ISO 25237 definition is replaced with the more accurate and descriptive term *de-identification*.

The definition of identifying information above suggests that it is easy – or at least possible – to distinguish personal information from identifying information. Indeed, many techniques for de-identification require an expert to make this distinction and protect only the identifying information. However, as understanding of privacy risk develops, it is increasingly apparent that *all* information is potentially identifying information.

This document envisions a *de-identification process* in which an *original dataset* that contains personal information is algorithmically processed to produce *de-identified data*. The result may be a *de-identified dataset*, *aggregate statistics* such as summary tables, or a *synthetic dataset*, in which the data are created by a model. This kind of de-identification is envisioned as a batch process. Alternatively, the de-identification process may be a system that accepts queries and returns responses that do not leak more identifying information than is allowable by policy. De-identified results may be corrected or updated and re-released on a periodic basis. The accumulated leakage of information from multiple releases may be significant, even if the leakage from a single release is small. Issues that arise from multiple releases are discussed in Section 3.4, “Data-Sharing Models.”

Disclosure is generally the exposure of data beyond the original collection use case. However, when the goal of de-identification is to protect privacy, disclosure

...relates to inappropriate attribution of information to a data subject, whether an individual or an organization. Disclosure occurs when a specific individual can be associated with a corresponding record(s) in the released dataset with high probability (*identity disclosure*), when an attribute described in a dataset is held by a specific individual, even if the record(s) associated with that individual is (are) not identified (*attribute disclosure*), or when it is possible to make an inference about an individual, even if the individual was not in the dataset prior to de-identification (*inferential disclosure*). [47, emphasis in original]

More information about disclosure can be found in Section 3.2.1, “Probability of Re-Identification.”

Disclosure limitation is a general term for the practice of allowing summary information or queries on data within a dataset to be released without revealing information about specific individuals whose personal information is contained within the dataset. Thus, de-identification is a kind of disclosure limitation technique. Every disclosure limitation process introduces inaccuracy into the results [14, 11].

A primary goal of disclosure limitation is to protect the privacy of individuals while avoiding the introduction of *non-ignorable biases* [7] (e.g., bias that might lead a social scientist to come to the wrong conclusion) into the de-identified dataset. One way to measure the amount of bias that has been introduced by the de-identification process is to compare statistics or models generated by analyzing the original dataset with those that are generated by analyzing the de-identified datasets. Such biases introduced by the de-identification

process are typically unrelated to any statistical biases that may also exist in the original data.

Formal models of privacy can quantify the amount of privacy protection offered by a de-identification process. With methods based on differential privacy, this measurement takes the form of a number called *privacy loss*, which quantifies the additional risk that an adversary might learn something new about an individual as a result of a de-identified data release. When a de-identification process is associated with low privacy loss, releasing the data it produces results in little additional risk for individuals in the input dataset. Some formal models, such as differential privacy, allow *composing* the privacy losses of multiple data releases to quantify the *total risk* to individuals of the combined releases, while others – such as *k*-anonymity – do not have this capability.

An upper bound on the total acceptable privacy loss of many data releases is often called a *privacy loss budget* or simply a *privacy budget*. This number quantifies the *total* privacy risk to an individual who participates in all of the releases.

Differential privacy [40] is a model based on a mathematical definition of privacy that considers the risk to an individual from the release of a query on a dataset containing their personal information. Statisticians, mathematicians, and other kinds of privacy engineers then develop mathematical algorithms, called mechanisms, that process data in a way that is consistent with the definition. Differential privacy limits both identity and attribute disclosure by adding non-deterministic noise (random values) to the results of mathematical operations before the results are reported. Unlike *k*-anonymity and other de-identification frameworks, differential privacy is based on information theory and makes no distinction between what is private data and what is not. Differential privacy does not require that values be classified as direct identifiers, quasi-identifiers, and non-identifying values. Instead, differential privacy assumes that *all values* in a record might be identifying and therefore all must be de-identified.

Differential privacy’s mathematical definition requires that the result of an analysis of a dataset should be roughly the same with or without the data of any single individual. The definition is usually satisfied by adding random noise to the result of a query, ensuring that the added noise masks the contribution of any individual. The degree of sameness is defined by the parameter ϵ (epsilon). The smaller the parameter ϵ , the more noise is added, and the more difficult it is to distinguish the contribution of a single individual. The result is increased privacy for all individuals – both those in the sample and those in the population from which the sample is drawn who are not present in the dataset. The research literature describes differential privacy being used to solve a variety of tasks, including statistical analysis, machine learning, and data sanitization [38]. Differential privacy can be implemented in an online query system or in a batch mode in which an entire dataset is de-identified at one time. In common usage, the phrase “differential privacy” is used to describe both the formal mathematical framework for evaluating privacy loss and for algorithms that provably provide those privacy guarantees.

The use of differential privacy algorithms does not guarantee that privacy will be preserved. Instead, the algorithms guarantee that the amount of privacy risk introduced by data processing or data release will reside within specific mathematical bounds. It is also important to remember that the impact on privacy risk is limited to reducing the risk of identity and attribute disclosures (see §3.2.1, “Probability of Re-Identification”) and not inferential disclosure.

K-anonymity [108, 123] is a framework for quantifying the amount of manipulation required of the quasi-identifiers to achieve a desired level of privacy. The technique is based on the concept of an *equivalence class* – the set of records that have the same values on the quasi-identifiers⁹. A dataset is said to be *k-anonymous* if there are no fewer than *k* matching records for every specific combination of quasi-identifiers. For example, if a dataset that has the quasi-identifiers (birth year) and (state) has *k*=4 anonymity, then there must be at least four records for every combination of (birth year, state). Subsequent work has refined *k-anonymity* by adding requirements for diversity of the sensitive attributes within each equivalence class (known as *l-diversity* [76]) and requiring that the resulting data be statistically close to the original data (known as *t-closeness* [73]).

K-anonymity and its subsequent refinements define formal privacy models but come with two important drawbacks. First, they require an expert to determine the set of quasi-identifiers by distinguishing between identifying and non-identifying information. As described earlier, this task can be difficult or impossible in some contexts. If identifying information is not marked as a quasi-identifier, then the resulting *k-anonymous* dataset will not prevent the re-identification of data subjects. Second, *k-anonymity* and related techniques are not compositional – they do not quantify the cumulative privacy loss of multiple data releases, and multiple releases can result in a catastrophic loss of privacy.

When data releases containing information about the same individual accumulate, then privacy loss accumulates. This accumulation of privacy loss is not reflected in *k-anonymity*, nor is it reflected in HIPAA privacy rule guidance [136]. Nevertheless, the accumulation is real. In 2003, Dinur and Nissim discovered that it was possible to reconstruct private microdata from a query interface even if the results of each query were systematically infused with small amount of noise [33]. The researchers showed that the amount of noise added to prevent an accurate reconstruction increases as the amount of queries on the dataset increase. If a query interface allows for an unlimited number of queries, no amount of noise is sufficient. Organizations should keep this in mind and try to assess the overall accumulated risk. The discovery in this paper led directly to the invention of differential privacy.

Some agencies (notably those that publish data for accountability and enforcement purposes) view perturbative Statistical Disclosure Limitation methods (e.g., those that add noise, such as differential privacy) as being inherently unacceptable, since the noise intro-

⁹A quasi-identifier is a variable that can be used to identify an individual through association with other information.

duced by the methods can void their ability to be used for accountability. For example, if a school would lose funding if the promotion rate for any class fell below a certain threshold, then a method that protects the privacy of students within each class by introducing noise could mask whether the school did or did not make that target. Thus, despite their weaknesses and flaws, program agencies often prefer to use suppression as the preferred protection method for these purposes because the data are either reported as is or suppressed, eliminating the uncertainty. Agencies should realize that suppression alone is not sufficient to protect privacy, and if a large enough number of queries is released based on the same confidential dataset, it is frequently possible to reconstruct even data that have been suppressed.

Traditional disclosure limitation and *k*-anonymity start with specific disclosure limitation mechanisms that were designed to hide information while allowing for useful data analysis and attempting to reach the goal of privacy protection. In contrast, differential privacy starts with an information-theoretic definition of privacy and has attempted to evolve mechanisms that produce useful (but privacy-preserving) results. These techniques are currently the subject of academic research, so it is reasonable to expect new techniques to be developed in the coming years that simultaneously increase privacy protection while providing for the high accuracy of resulting de-identified data. Indeed, some authors have shown that the models can be viewed synergistically [114] under some circumstances.

Finally, privacy harms are not the only kinds of harms that can result from the release of de-identified data. Analysts working with de-identified data often have no way of knowing how inaccurate their statistical results are due to statistical distortions introduced by the de-identification process. Thus, de-identification operations intended to shield individuals from harm could result in inaccurate research findings. Such research might also cause harm if it is used to support harmful policies.

3. Governance and Management of Data De-Identification

The decisions and practices regarding the de-identification and release of government data can be integral to the mission and proper functioning of a government agency. As such, these activities should be managed by an agency's leadership in a way that assures that performance and results that are consistent with the agency's mission and legal authority. As discussed above, the need for attention arises because of the conflicting goals of data transparency and privacy protection. Although many agencies once assumed that it was relatively straightforward to remove privacy-sensitive data from a dataset so that the remainder could be released without restriction, history shows that this is not the case [51, §2.4, §3.6].

Given this history, there may be a tendency for government agencies to either over-protect data or to simply avoid its release. Limiting the release of data clearly limits the privacy risk that might result from a data release. However, limiting the release of data also creates costs and risks for other government agencies (which will then not have access to the identified data), external organizations, and society. For example, absent the data release, external organizations will suffer the cost of recollecting the data (if it is possible to do so) or the risk of incorrect decisions that might result from having insufficient information.

This section begins with a discussion of why agencies might wish to de-identify data and how agencies should balance the benefits of data release with risks to the data subjects. It then discusses where de-identification fits within the data life cycle. Finally, it discusses options that agencies have for adopting de-identification standards.

3.1. Identifying Goals and Intended Uses of De-Identification

Before engaging in de-identification, agencies should clearly articulate their goals regarding transparency and disclosure limitation in making a data release. They should then develop a written plan that explains how de-identification will be used to accomplish those goals.

For example:

- **Federal Statistical Agencies** collect, process, and publish data for use by researchers, business planners, and other well-established purposes. These agencies are likely to have established standards and methodologies for de-identification. As these agencies evaluate new approaches for de-identification, they should document their rationale for adopting legacy versus new approaches, evaluate how successful their approaches have been over time, and address inconsistencies between data releases.
- **Federal Awarding Agencies** are allowed under OMB Circular A-110 to require that institutions of higher education, hospitals, and other non-profit organizations that receive federal grants provide the U.S. Government with "the right to (1) obtain, reproduce, publish or otherwise use the data first produced under an award; and

871 (2) authorize others to receive, reproduce, publish, or otherwise use such data for
872 Federal Purposes” [91, see §36 (c) (1) and (2)]. To realize this policy, awarding
873 agencies can require that awardees establish data management plans for making re-
874 search data publicly available. Such data are used for a variety of purposes, including
875 transparency and reproducibility. In general, research data that contain personal in-
876 formation should be de-identified by the awardee prior to public release. Awarding
877 agencies may establish de-identification standards to ensure the protection of per-
878 sonal information and may consider audits to assure that awardees have performed
879 de-identification in an appropriate manner.

- 880 • **Federal Research Agencies** may wish to make de-identified data available to the
881 public to further the objectives of research transparency and allow others to reproduce
882 and build upon their results. These agencies are generally prohibited from publishing
883 research data that contain personal information, requiring the use of de-identification.
- 884 • **All Federal Agencies** that wish to make administrative or operational data available
885 for transparency, accountability, or program oversight or to enable academic research
886 may wish to employ de-identification to avoid sharing sensitive personally identifi-
887 able information of employees, customers, or others. These agencies may wish to
888 evaluate the effectiveness of simple field suppression, de-identification that involves
889 aggregation, and the creation and release of synthetic data as alternatives for realizing
890 their commitment to open data.

891 3.2. Evaluating Risks that Arise from De-Identified Data Releases

892 Once the purpose of the data release is understood, agencies should identify the risks that
893 might result from the data release. As part of this risk analysis, agencies should specifically
894 evaluate the anticipated negative actions that might result from re-identification, as well as
895 strategies for remediation. NIST provides detailed information on how to conduct risk
896 assessments in NIST SP 800-30 [23].

897 Risk assessments should be based on objective scientific factors and consider the best inter-
898 ests of the individuals in the dataset, the responsibilities of the agency holding the data, and
899 the anticipated benefits to society. The goal of a risk evaluation is not to eliminate risk but
900 to identify which risks can be reduced while still meeting the objectives of the data release
901 and then deciding whether the residual risk is justified by the goals of the data release. An
902 agency decision-making process may choose to accept or reject the risk that might result
903 from a release of de-identified data, but participants in the risk assessment should not be
904 empowered to prevent risk from being documented and discussed. Centralized processes
905 also allow for standardization of the risk assessment and the amount of “acceptable risk”
906 across different programs’ releases.

907 It is difficult to measure re-identification risk in ways that are both general and meaningful.
908 For example, it is possible to measure the similarity between individuals in the dataset

under a variety of different parameters and to model how that similarity is impacted when the larger population is considered. However, such calculations may result in different levels of risk for different groups. There may be some individuals in a dataset who would be significantly adversely impacted by re-identification and for whom the likelihood of re-identification might be quite high, but these individuals might represent a tiny fraction of the entire dataset. This represents an important area for research in the field of risk communication.

3.2.1. Probability of Re-Identification

As discussed in Section 2.2, “Terminology,” the potential impacts on individuals from the release and use of de-identified data include [143]:

Identity disclosures: Associating a specific individual with the corresponding record(s) in the dataset with high probability. Identity disclosure can result from insufficient de-identification, re-identification by linking, or pseudonym reversal.

Attribute disclosure: Determining that an attribute described in the dataset is held by a specific individual with high probability, even if the records associated with that individual are not identified. Attribute disclosure can occur without identity disclosure if the de-identified dataset contains data from a significant number of relatively homogeneous individuals [51, p.13]. In these cases, traditional de-identification does not protect against attribute disclosure, although differential privacy can. Membership inference is an example of attribute disclosure.

Inferential disclosure: Being able to make an inference about an individual (typically a member of a group) with high probability, even if the individual was not in the dataset prior to de-identification. “Inferential disclosure is of less concern in most cases as inferences are designed to predict aggregate behavior, not individual attributes, and thus are often poor predictors of individual data values” [60]. Traditional de-identification does not protect against inferential disclosure. Such disclosures can never be eliminated; they can only be controlled.

*Re-identification probability*¹⁰ is the estimated probability that an outside party will be able to use information contained in a de-identified dataset to make identity-related inferences about individuals. This outside party was originally termed a *data intruder*, although the terms *adversary* and *attacker* are also used, borrowing from the colorful language of information security. Different kinds of re-identification probabilities for this data intruder can be calculated.

¹⁰Previous publications described identification probability as “re-identification risk” and used scenarios such as a journalist seeking to discredit a national statistics agency or a prosecutor seeking to find information about a suspect as the bases for probability calculations. That terminology is not presented in this document because of the possible unwanted connotations of those terms and in the interest of bringing the terminology of de-identification into agreement with the terminology used in contemporary risk analysis processes [42].

Here are several kinds of probabilities, as well as proposals for new, declarative, self-describing names:

Known inclusion re-identification probability (KIRP) is the probability of finding the record that matches a specific individual known to be in the sample corresponding to a specific record. KIRP can be expressed as the probability for a specific individual or the probability averaged over the entire dataset (AKIRP).¹¹

Unknown inclusion re-identification probability (UIRP) is the probability of finding the record that matches a specific individual without first knowing whether the individual is in the dataset. UIRP can be expressed as a probability for an individual record in the dataset averaged over the entire population (AUIRP).¹²

Record matching probability (RMP) is the probability of finding the record that matches a specific individual chosen from the population. RMP can be expressed as the probability for a specific record (RMP), the probability averaged over the entire dataset (ARMP), or the maximum probability over the entire dataset.

Inclusion probability (IP) is the probability that a specific individual's presence in the dataset can be inferred.

Whether it is necessary to quantitatively estimate these probabilities depends on the specifics of each intended data release. For example, many cities publicly disclose whether taxes have been paid on a property. Given that this information is already a matter of public record, it may not be necessary to consider inclusion probability when a dataset of property taxpayers for a specific dataset is released. Likewise, there may be some attributes in a dataset that are already public and may not need to be protected with disclosure limitation techniques. However, the existence of such attributes may pose a re-identification risk for other information in the dataset or in other de-identified datasets. The fact that information is public may not negate the responsibility of an agency to provide protection for that information, as the aggregation and distribution of information may cause privacy risk that was not otherwise present. Agencies may also be legally prohibited from releasing copies of information that is similar to information that is already in the public domain.

Although disclosures are commonly thought to be discrete events involving the release of specific data, such as an individual's name matched to a record, disclosures can result from the release of data that merely changes a data intruder's probabilistic belief. For example, a disclosure might change an intruder's estimate that a specific individual is present in a dataset from a 50% probability to 90%. The intruder still does not *know* if the individual is in the dataset or not (and the individual might not, in fact, be in the dataset), but a

¹¹Some texts refer to KIRP as "prosecutor risk." The scenario is that a prosecutor is looking for records that belong to a specific, named individual.

¹²Some texts refer to UIRP as "journalist risk." The scenario is that a journalist has obtained a de-identified file and is trying to identify one of the data subjects, but the journalist fundamentally does not care *who* is identified.

976 probabilistic disclosure has still occurred because the intruder's estimate of the individual
977 has been changed by the data release.

978 It may be difficult to estimate specific re-identification probabilities, as the ability to re-
979 identify depends on the original dataset, the de-identification technique, the technical skill
980 of the data intruder, the intruder's available resources, and the availability of additional
981 data (publicly available or privately held) that can be linked with the de-identified data.
982 It is likely that the true probability of re-identification increases over time as techniques
983 improve and more contextual information becomes available to potential data intruders.
984 Indeed, some researchers have claimed that computing these probabilities "is a fundamen-
985 tally meaningless exercise" because the calculations are based on assumptions that cannot
986 be validated (e.g., the lack of a database that could link specific quasi-identifiers or sensi-
987 tive, non-identifying values to identities) [83].

988 De-identification practitioners have traditionally quantified re-identification probability, in
989 part, based on the skills and abilities of a potential data intruder. Datasets that were thought
990 to have little possibility for exploitation were deemed to have a lower re-identification
991 probability than datasets containing sensitive or otherwise valuable information. Such ap-
992 proaches are not appropriate when attempting to evaluate the re-identification probability
993 of government datasets that will be publicly released.

- 994 • Although a specific de-identified dataset may not be recognized as sensitive, re-
995 identifying that dataset may be an important step in re-identifying another dataset
996 that is sensitive. Alternatively, the data intruder may merely wish to embarrass the
997 government agency. Thus, adversaries may have a strong incentive to re-identify
998 datasets that are seemingly innocuous.
- 999 • Although the public may not generally be skilled in re-identification, many resources
1000 on the internet make it easy to acquire specialized datasets, tools, and experts for
1001 specific re-identification challenges. Family members, friends, colleagues, and oth-
1002 ers may also possess substantial personal knowledge about individuals in the data
1003 that can be used for re-identification.

1004 Instead, de-identification practitioners should assume that de-identified government datasets
1005 could be subjected to sustained, worldwide re-identification attempts, and they should
1006 gauge their de-identification requirements accordingly. Of course, it is unrealistic to as-
1007 sume that all of the world's resources will be used to attempt to re-identify every publicly
1008 released file. Therefore, de-identification requirements should be gauged using a risk as-
1009 sessment [75]. More information on conducting risk assessments can be found in NIST SP
1010 800-30, *Guide for Conducting Risk Assessments* [23].

1011 Members of vulnerable populations (e.g., prisoners, children, people with disabilities) may
1012 be more susceptible to having their identities disclosed by de-identified data than non-
1013 vulnerable populations because the thing that makes these individuals vulnerable may also
1014 make them stand out in the dataset. Likewise, residents of areas with small populations

may be more susceptible to having their identities disclosed than residents of urban areas. Individuals with multiple traits will generally be more identifiable if the individual's location is geographically restricted. For example, data belonging to a person who is labeled as a pregnant, unemployed female veteran will be more identifiable if restricted to Baltimore County, Maryland, than to all of North America.

If agencies determine that the potential for harm is large in a contemplated data release, one way to manage the risk is by increasing the level of de-identification and accepting a lower data accuracy level. Other options include data controls, such as restricting the availability of data to qualified researchers in a data enclave.

3.2.2. Adverse Impacts of Re-Identification

As part of a risk analysis, agencies should attempt to enumerate specific kinds of adverse impacts that can result from the re-identification of de-identified information. These can include potential impacts on individuals, the agency, and society.

Potential adverse impacts on individuals include:

- Increased availability of personal information that leads to an increased risk of fraud, identity theft, discrimination, or abuse
- Increased availability of an individual's location that puts that person at risk for burglary, property crime, assault, or other kinds of violence
- Increased availability of an individual's non-public personal information that causes psychological harm by exposing potentially embarrassing information or information that the individual may not otherwise choose to reveal to the public or to family members and that potential affects opportunities in the economic marketplace (e.g., employment, housing, college admission)

Potential adverse impacts on agencies include:

- Mandatory reporting under breach reporting laws, regulations, or policies
- Embarrassment or reputational damage
- Harm to agency operations if some aspect of those operations required that the de-identified data remain confidential (e.g., an agency that is forced to discontinue a scientific experiment because the data release may have biased the study participants)
- Financial impacts that result from the harm to the individuals (e.g., lawsuits)
- Civil or criminal sanctions against employees or contractors that result from a data release contrary to U.S. law

Potential adverse impacts on society include:

- Undermining the reputation of researchers in general and the willingness of the public to support/tolerate research and provide accurate information to government agencies and researchers
- Engendering a lack of trust in government – individuals may stop consenting to the use of their data, may stop providing data, or may provide false data
- Damaging the practice of using de-identified information – de-identification is an important tool for promoting research and accountability, and poorly executed de-identification efforts may negatively impact the public’s view of this technique and limit its use

One way to calculate an upper bound on impact to an individual or the agency is to estimate the impact that would result from the inadvertent release of the original dataset. This approach will not calculate the upper bound on the societal impact, however, since that impact includes reputational damage to the practice of de-identification itself. That is, every time data are compromised because of a poorly executed de-identification effort, it becomes harder to justify the use of de-identification in future data releases.

As part of a risk analysis process, organizations should enumerate specific measures that they will take to minimize the risk of successful re-identification. Organizations may wish to consider both the actual risk and the perceived risk to those in the dataset and in the broader community.

As part of the risk assessment, an organization may determine that there is no way to achieve the de-identification goal in terms of data accuracy and identifiability. In these cases, the organization will need to decide whether it should adopt additional measures to protect privacy (e.g., administrative controls or data use agreements), accept a higher level of risk, or choose not to proceed with the project.

3.2.3. Impacts Other Than Re-Identification

The use of de-identified data can lead to adverse impacts other than those that might result from re-identification. Risk assessments that evaluate the risks of re-identification can address these other risks as well. Such risks might include:

- The risk of excessive inferential disclosures
- The risk that the de-identification process might introduce bias or inaccuracies into the dataset that result in incorrect decisions¹³

¹³For example, a personalized warfarin dosing model created with data that had been modified in a manner consistent with the differential privacy de-identification model produced higher mortality rates in simulation than a model created from unaltered data [49]. Educational data de-identified with the *k*-anonymity model can also result in the introduction of bias that leads to spurious results [14, 125].

- The risk that releasing a de-identified dataset might reveal non-public information about an agency’s policies or practices

It is preferable to use de-identification processes that include assessments of accuracy (e.g., confidence intervals) with respect to the bias and precision of statistical properties of the data. Where it does not provide information that may aid data intruders, it is also useful to reveal the de-identification process itself so that analysts can understand any potential inaccuracies that might be introduced by the de-identification. This is consistent with *Kerckhoffs’ principle* [67], a widely accepted system design principle that holds that the security of a system should not rely on the secrecy of the methods that it employs.

3.2.4. Remediation

As part of a risk analysis process, agencies should attempt to enumerate techniques that could be used to mitigate or remediate harm that would result from a successful re-identification of de-identified information. Remediation could include victim education, the procurement of monitoring or security services, the issuance of new identifiers, or other measures.

3.3. Data Life Cycle

The *NIST Big Data Interoperability Framework* defines the data life cycle as “the set of processes in an application that transform raw data into actionable knowledge” [85]. The data life cycle can be used in the de-identification process to help analyze the expected benefits, intended uses, privacy threats, and vulnerabilities of de-identified data. As such, the data life cycle concept can be used to select appropriate privacy controls based on a reasoned analysis of the threats. For example, privacy-by-design concepts [22] can be employed to decrease the number of identifiers collected, minimizing requirements for de-identification prior to data release. The data life cycle can also be used to design a tiered access mechanism based on this analysis [12].

Several data life cycles have been proposed, but none are widely accepted as a standard.

Michener et al. [80] (Figure 1) describe the data life cycle as a true cycle:

→ Assure → Describe → Deposit → Preserve → Discover → Integrate → Analyze →
Collect

Stobierski [120] also describes the data life cycle as a cycle with different steps:

Generation → Collection → Processing → Storage → Management → Analysis →
Visualization → Interpretation → Generation

De-identification does not fit into a circular data life cycle model, as the data owner typically retains access to the identified data. However, if the organization employs de-identification, it could be performed during Collect or between Collect and Assure if identified data were collected but the identifying information was not actually needed. Alter-

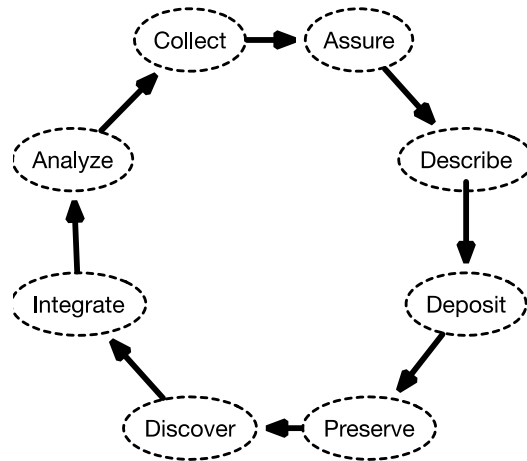


Fig. 1. The data life cycle as described by Michener et al. [80]

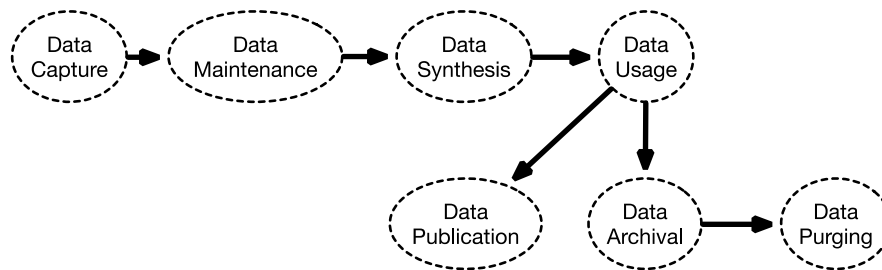


Fig. 2. Chisholm's view of the data life cycle is a linear process with a branching point after data usage [25]

1114 natively, de-identification could be applied after Describe and prior to Deposit to avoid
1115 archiving identifying information.

1116 Chisholm and others [25] (Figure 2) describe the data life cycle as a linear process with a
1117 fork for data publication:

1118 Data Capture → Data Maintenance → Data Synthesis → Data Usage →
1119 {Data Publication & Data Archival → Data Purging}

1120 Using this formulation, de-identification can take place either during Data Capture or fol-
1121 lowing Data Usage. However, agencies should consider data release requirements from
1122 the very beginning of the planning process for each new data collection. By knowing in
1123 advance how they intend to publish and for what purposes and by having a plan for how
1124 disclosure limitation will be applied, agencies can tailor information collection accordingly.

1125 For example, if specific identifiers are not needed for maintenance, synthesis, and usage,
1126 then those identifiers should not be collected. If fully identified data are needed within the

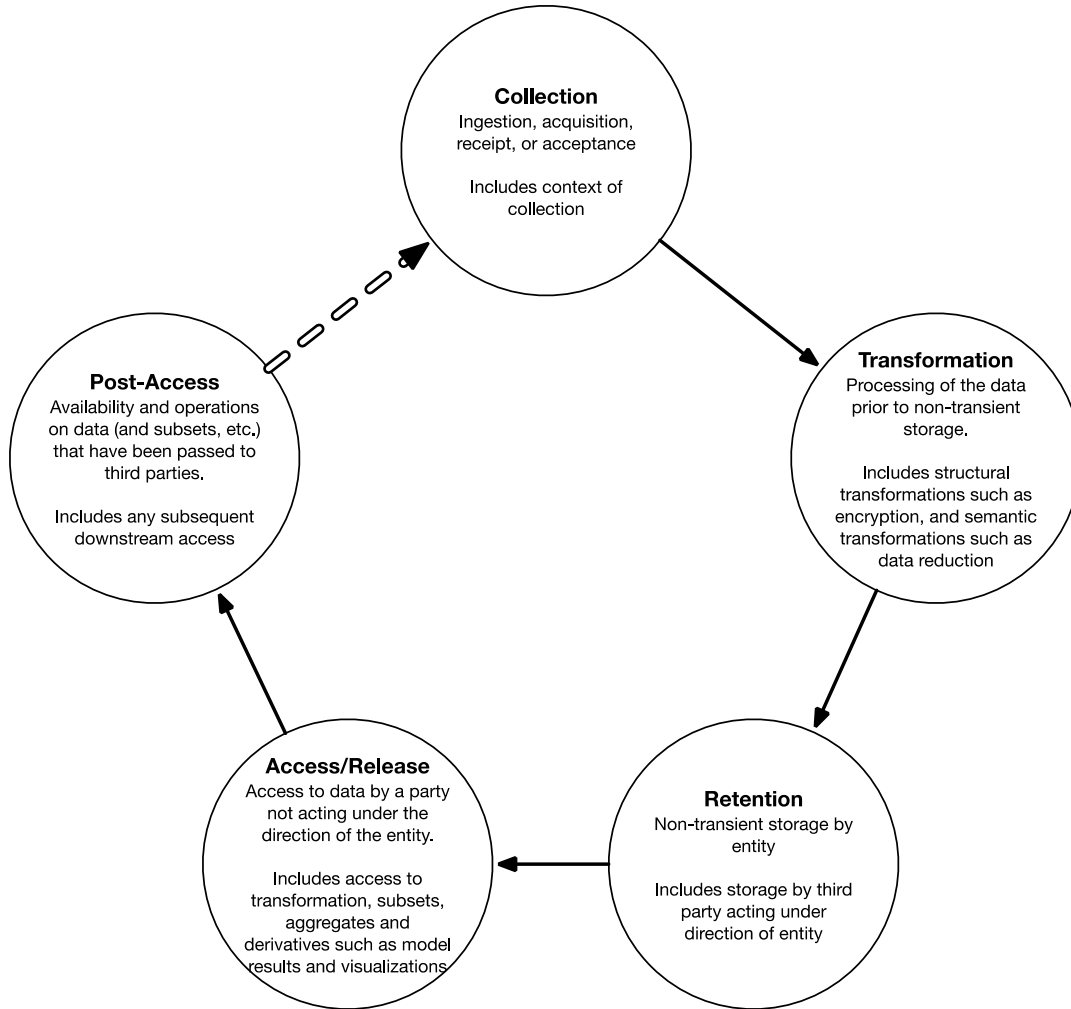


Fig. 3. Altman’s “modern approach to privacy-aware government data releases” [75]

1127 organization, the identifying information can be removed prior to the data being published,
1128 shared, or archived. Applying de-identification throughout the data life cycle minimizes
1129 privacy risk and significantly eases the process of public release. However, agencies should
1130 be cognizant of the potential loss of future utility if identifiers are permanently removed.
1131 For this reason, agencies may wish to retain an identified dataset or data linking informa-
1132 tion, as it may be difficult to predict future needs.

1133 Altman et al. [75] (Figures 3 and 4) propose a “modern approach to privacy-aware gov-
1134 ernment data releases” that incorporates progressive levels of de-identification as well as
1135 different kinds of access and administrative controls in line with the sensitivity of the data.

1136 Agencies that perform de-identification should document that:

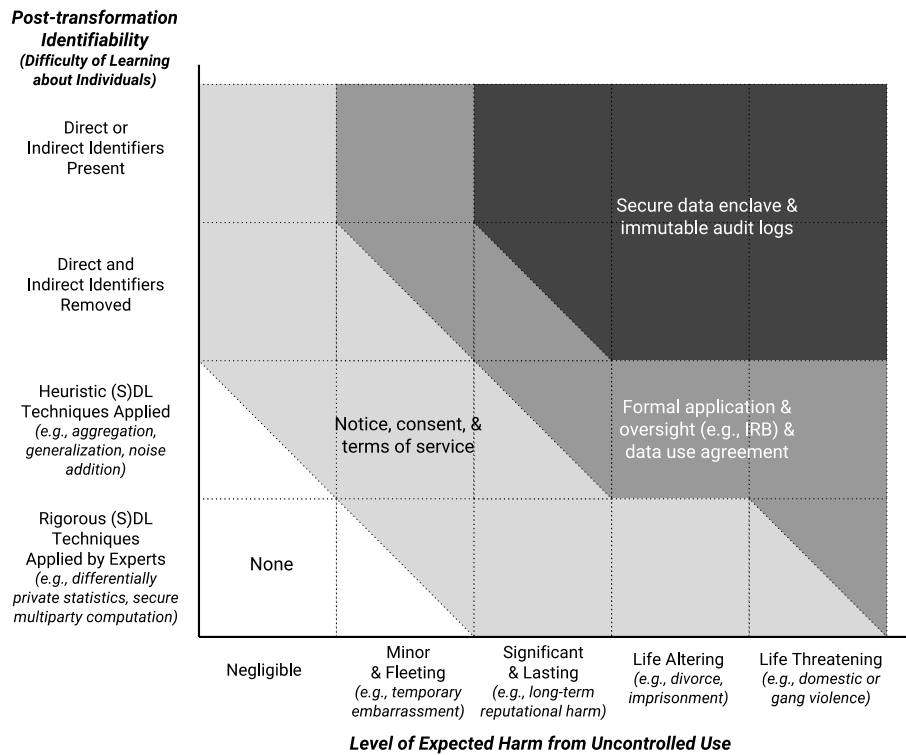


Fig. 4. Altman's conceptual diagram of the relationship between post-transformation identifiability, level of expected harm, and suitability of selected privacy controls for a data release [75]

- 1137 • The techniques used to perform the de-identification are theoretically sound and gen-
1138 erally accepted.¹⁴
 - 1139 • The software used to perform the de-identification is reliable for the intended task.
 - 1140 • The individuals who performed the de-identification were suitably qualified.
 - 1141 • The tests that were used to evaluate the effectiveness of the de-identification were
1142 validated for that purpose.
 - 1143 • Ongoing monitoring is in place to ensure the continued effectiveness of the de-
1144 identification strategy.
- 1145 No matter where de-identification is applied in the data life cycle, agencies should docu-
1146 ment the answers to the following questions for each de-identified dataset:
- 1147 • Are direct identifiers collected with the dataset?
 - 1148 • Even if direct identifiers are not collected, is it still possible to identify the data
1149 subjects through the presence of quasi-identifiers?
 - 1150 • Where in the data life cycle is de-identification performed? Is it performed in only
1151 one place or in multiple places?
 - 1152 • Is the original dataset retained after de-identification?
 - 1153 • Is there a key or map retained so that specific data elements can be re-identified later?
 - 1154 • How are decisions made regarding de-identification and re-identification?
 - 1155 • Are there specific datasets that can be used to re-identify the de-identified data? If so,
1156 what controls are in place to prevent intentional or unintentional re-identification?
 - 1157 • Is it a problem if some records in a dataset are re-identified?

¹⁴To determine that a technique is theoretically sound and generally accepted, agencies that wish to adopt guidance that mirrors the language that the HHS November 26, 2012 *Guidance Regarding Methods for De-identification of Protected Health Information in Accordance with the Health Insurance Portability and Accountability Act (HIPAA) Privacy Rule* [136]. uses in its discussion of the Privacy Rule’s “expert determination method,” which states on page 7:

“A covered entity may determine that health information is not individually identifiable health information only if:

- (1) A person with appropriate knowledge of and experience with generally accepted statistical and scientific principles and methods for rendering information not individually identifiable:
 - (i) Applying such principles and methods, determines that the risk is very small that the information could be used, alone or in combination with other reasonably available information, by an anticipated recipient to identify an individual who is a subject of the information; and
 - (ii) Documents the methods and results of the analysis that justify such determination;”

- Is there a mechanism that will inform the de-identifying agency if there is an attempt to re-identify the de-identified dataset? Is there a mechanism that will inform the agency if the attempt is successful?

3.4. Data-Sharing Models

Agencies should decide on the data-sharing model that will be used to make the data available outside of the agency after the data have been de-identified [51, p.14]. Specific models combine *security* and *privacy* techniques to reduce privacy risks to individuals. *Security* refers to techniques that limit *who* can view the data. Encryption is an example of a security technique – it allows only the party holding the encryption key to view the data. *Privacy* refers to techniques that limit *what information* the data contains. The two concepts can be considered orthogonally. In practice, however, *who* has access to the data makes a significant difference in the expected risk of disclosure and therefore influences the extent to which privacy techniques must be used to limit the presence of sensitive personally identifiable information in the data.

A number of possible models exist at different points in the spectrum of security and privacy protections. Figure 4 summarizes this spectrum: its x-axis describes various privacy techniques that can limit the informational content of the data; its y-axis describes how much harm would occur if the underlying information were disclosed; and the regions of the graph are labeled with suggested security techniques. Some common combinations of security and privacy techniques include:

The Release and Forget Model [94]. The de-identified data may be released to the public, typically by being published on the internet. It can be difficult or impossible for an organization to recall the data once released in this fashion and may limit information for future releases.

The Data Use Agreement (DUA) Model. The de-identified data may be made available under a legally binding data use agreement that details what can and cannot be done with the data. Typically, data use agreements may prohibit attempted re-identification, linking to other data, and redistribution of the data without a similarly binding DUA. A DUA will typically be negotiated between the data holder and qualified researchers (the “qualified investigator model” [44]) or members of the general public (e.g., citizen scientists or the media), although they may be simply posted on the internet with a click-through license agreement that must be agreed to before the data can be downloaded (the “click-through model” [44]).

The Synthetic Data with Verification Model. Statistical disclosure limitation techniques are applied to the original dataset and used to create a synthetic dataset that contains many of the aspects of the original dataset but does not contain disclosing information. The synthetic dataset is released, either publicly or to vetted researchers. The synthetic dataset can then be used as a proxy for the original dataset, and if

constructed well, the results of statistical analyses should be similar. If used in conjunction with an enclave model as below, researchers may use the synthetic dataset to develop queries and/or analytic software. These queries and/or software can then be taken to the enclave or provided to the agency and be applied on the original data.

The Enclave Model [44, 87, 113]. The de-identified data may be kept in a segregated enclave that restricts the export of the original data and instead accepts queries from qualified researchers, runs the queries on the de-identified data, and responds with results. Enclaves can be physical or virtual and can operate under a variety of different models. For example, vetted researchers may travel to the enclave to perform their research, as is done with the Federal Statistical Research Data Centers operated by the U.S. Census Bureau. Enclaves may be used to implement the verification step of the Synthetic Data with Verification Model. Queries made in the enclave model may be vetted automatically or manually (e.g., by the DRB). Vetting can try to screen for queries that might violate privacy or are inconsistent with the stated purpose of the research.

Sharing models should consider the possibility of multiple or periodic releases. Just as repeated queries to the same dataset may leak personal data from the dataset, repeated de-identified releases (whether from the same dataset or from different datasets containing some of the same individuals) by an agency may result in compromising the privacy of individuals unless each subsequent release is viewed in light of the previous release. Even if a contemplated release of a de-identified dataset does not directly reveal identifying information, federal agencies should ensure that the release – combined with previous releases – will also not reveal identifying information [137].

Instead of sharing an entire dataset, the data owner may choose to release a sample. If only a sample is released, the probability of re-identification decreases because a data intruder will not know if a specific individual from the data universe is present in the de-identified dataset [43]. However, releasing only a sample may decrease the statistical power of tests on the data, may cause users to draw incorrect inferences if proper statistical sampling methods are not used, and may not align with agency goals regarding transparency and accountability.

3.5. The Five Safes

Agencies that make data available to outsiders should use a repeatable methodology for evaluating the terms under which that data will be made available. The Five Safes [31] is such a framework.

The Five Safes was created in the United Kingdom to assist a national statistical agency in evaluating proposed collaborative projects with the larger research community. The framework is designed to assist in “designing, describing and evaluating” data access systems. Here, the term “data access system” is viewed broadly as any mechanism that allows out-

siders to gain access to the agency’s confidential data. That is, a data access system might include setting up an enclave for academic researchers who undergo extensive background checks, but it also includes publishing data on the internet.

The Five Safes framework gets its name from the use of five categories (called “risk” or “access” dimensions) that are used in the evaluation. They are:

1. **Safe projects** Is this use of the data appropriate?
2. **Safe people** Can the researchers be trusted to use it in an appropriate manner?
3. **Safe data** Is there a disclosure risk in the data itself?
4. **Safe settings** Does the access facility limit unauthorized use?
5. **Safe outputs** Are the statistical results non-disclosive?

Each of these dimensions is independent. That is, the legal, moral, and ethical review of each dimension is independent of the others. In practice, this might mean that the project is safe (the proposed use of the data is appropriate), the people are safe (the researchers are noted academics with respected histories of collaborative work), the data are safe (there is no disclosure risk in the data), and the output is safe (it will not disclose personal information). However, because the setting is not safe (perhaps the facility has poor internal security), the project should not go forward. In this example, the Five Safes framework would provide a decision-maker with the tools to separate each of these dimensions and resolve the problems so that the project could proceed.

One of the positive aspects of the Five Safes framework is that it forces data controllers to consider many different aspects of data release when evaluating data access proposals. Frequently, the authors write, it is common for data owners to “focus on one, and only one, particular issue (such as the legal framework surrounding access to their data or IT solutions).” With the Five Safes, people who may be specialists in one area are forced to consider (or to explicitly not consider) aspects of privacy protection with which they may not be familiar and might otherwise overlook.

The Five Safes framework can be used as a tool for designing access systems, for evaluating existing systems, for communication, and for training. Agencies should consider using a framework such as The Five Safes for organizing risk analyses of data release efforts.

3.6. Disclosure Review Boards

Disclosure Review Boards (DRBs), also known as Data Release Boards, are administrative bodies created within an organization that are charged with ensuring that intended disclosures meet the policy and procedural requirements of that organization. DRBs should be governed by a written *mission statement* and *charter* (or equivalent document) that are ideally approved by the same mechanisms that the organization uses to approve other organization-wide policies.

The DRB should have a mission statement that guides its activities. For example, the U.S. Department of Education’s DRB has the mission statement:

The Mission of the Department of Education Disclosure Review Board (ED-DRB) is to review proposed data releases by the Department’s principal offices (POs) through a collaborate technical assistance, aiding the Department to release as much useful data as possible, while protecting the privacy of individuals and the confidentiality of their data, as required by law. [41]

The DRB charter specifies the mechanics of how the mission is implemented. A formal, written charter promotes transparency in the decision-making process and ensures consistency in the applications of its policies.

Most DRBs will be established to weigh the interests of data release against those of individual privacy protection. However, a DRB may also be chartered to consider *group harms* [51, p.13] that can result from the release of a dataset. Such harms go beyond the harm to the privacy interests of a specific individual.

The DRB charter should frame the DRB’s responsibilities in reference to existing organizational policies, regulations, and laws. Some agencies may balance these concerns by employing data use models other than de-identification (e.g., by establishing data enclaves where a limited number of vetted researchers can access sensitive datasets in a way that provides data value while minimizing the possibility for harm or by authorizing the use of secure multi-party computation, homomorphic encryption, or other Privacy Preserving Data Analytics to compute various statistics). In those agencies, a DRB would be empowered to approve the use of such mechanisms.

Certain agencies may engage in data disclosure on a routine basis (such as research and evaluation agencies), in which case it may be beneficial for the DRB to establish policies and procedures for de-identification rather than being responsible for every review. In these cases, the DRB charter should clearly specify how the group will provide oversight and ensure organizational accountability to the agreed-upon policies.

The DRB charter should specify the DRB’s composition. To be effective, the DRB should include representatives from multiple groups and experts in both technology and privacy policy. Specifically, DRBs may wish to have as members:

- Individuals who represent the interests of potential users (such individuals need not come from outside of the organization)
- Representation from among the public, specifically from groups represented in the datasets if they have a limited scope
- Representation from the organization’s leadership team, such as a representation of the Senior Agency Official for Privacy [4, Appendix II, section 4] (such representation helps to establish the DRB’s credibility with the rest of the organization)

- 1307 • A representative of the organization’s senior privacy official
- 1308 • Subject matter experts
- 1309 • Outside experts

1310 The charter should establish rules for ensuring a quorum and specify whether members can
1311 designate alternates on a standing or meeting-by-meeting basis. The DRB should specify
1312 the mechanism by which members are nominated and approved, their tenure, conditions
1313 for removal, and removal procedures.¹⁵

1314 The charter should set policy expectations for record keeping and reporting, including
1315 whether records and reports are considered public or restricted. For example, the char-
1316 ter could specify that a DRB issue an annual report with a list of every dataset that was
1317 approved for release. The charter should indicate whether it is possible to exclude sensitive
1318 decisions from these reporting requirements and the mechanism for doing so. Ideally, the
1319 charter should be a public document to promote transparency.

1320 To meet its requirement of evaluating data releases, the DRB should require that writ-
1321 ten applications be submitted to the DRB that specify the nature of the dataset, the de-
1322 identification methodology, and the result. An application may require that the proposer
1323 present the re-identification risk, the risk to individuals if the dataset is re-identified, and
1324 a proposed plan for detecting and mitigating successful re-identification. In addition, the
1325 DRB should require that when individuals are informed that their information will be de-
1326 identified, they also be informed that privacy risks may remain despite de-identification.

1327 The DRB should keep accurate records of its request memos, their associated documen-
1328 tation, the DRB decision, and the actual files released. These records should be appropri-
1329 ately archived and curated so that they can be recovered. In the case of large data releases,
1330 the definitive version of the released data should be curated using an externally validated
1331 procedure, such as a recorded cryptographic hash value or signature, and a digital object
1332 identifier (DOI) [64].

1333 DRBs may wish to institute a two-step process in which the applicant first proposes and
1334 receives approval for a specific de-identification process that will be applied to a specific
1335 dataset and then submits and receives approval for the release of the dataset that has been
1336 de-identified according to the proposal. However, because it is theoretically impossible
1337 to predict the results of applying an arbitrary process to an arbitrary dataset [26, 129],
1338 the DRB should be empowered to reject a proposed release of a dataset even if it has
1339 been de-identified in accordance with an approved procedure because performing the de-
1340 identification may demonstrate that the procedure was insufficient to protect privacy. The

¹⁵For example, in 2022, the Census Bureau’s DRB had 12 voting members: two technical co-chairs, a repre-
sentative from the Policy Coordination Office, a representative from the Associate Director for Communica-
tions, two representatives from the Center for Enterprise Dissemination-Disclosure Avoidance (CED-DA),
two representatives from the Economic Programs Directorate, two representatives from the Demographic
Programs Directorate, and two representatives from the Decennial Programs Directorate [24].

DRB should be able to delegate the responsibility of reviewing the de-identified dataset, but such responsibility should not be delegated to the individual or group that performed the de-identification.

The DRB charter should specify whether the DRB needs to approve each data release by the organization or if it may grant blanket approval for all data of a specific type that is de-identified according to a specific methodology. The charter should specify the duration of the approval. Given advances in the science and technology of de-identification, it is inadvisable that a Board be empowered to grant release authority for an indefinite or unlimited amount of time.

In most cases, a single privacy protection methodology will be insufficient to protect the varied datasets that an agency may wish to release. That is, different techniques might best optimize the trade-off between re-identification risk and data usability, depending on the specifics of each kind of dataset. Nevertheless, the DRB may wish to develop guidance, recommendations, and training materials regarding specific de-identification techniques that are to be used. Agencies that standardize on a small number of de-identification techniques will gain familiarity with these techniques and are likely to have results with a higher level of consistency and success than those that have no such guidance or standardization.

Although it is envisioned that DRBs will work in a cooperative, collaborative, and congenial manner with those inside an agency seeking to release de-identified data, there will at times be a disagreement of opinion. For this reason, the DRB's charter should state whether the DRB has the final say over disclosure matters or if the DRB's decisions can be overruled, by whom, and by what procedure. For example, an agency might give the DRB final say over disclosure matters but allow the agency's leadership to replace members of the DRB as necessary. Alternatively, the DRB's rulings might merely be advisory, with all data releases being individually approved by agency leadership or its delegates.¹⁶

Finally, agencies should decide whether the DRB charter will include any kind of performance timetables or be bound by a service-level agreement (SLA) that defines a level of service to which the DRB commits.

The key elements of a Disclosure Review Board include:

- A written mission statement and charter
- Members represent different groups within the organization, including leadership
- The Board receives written applications to release de-identified data
- The Board reviews *both* the proposed methodology *and* the results of applying the methodology

¹⁶At the Census Bureau, "staff members [who] are not satisfied with the DRB's decision . . . may appeal to a steering committee consisting of several Census Bureau Associate Directors. Thus far, there have been few appeals, and the Steering Committee has never reversed a decision made by the Board" [130, p.35].

- 1375 • Applications should identify the risks associated with data release, including re-
1376 identification probability, potentially adverse events that would result if individuals
1377 are re-identified, and a mitigation strategy if re-identification takes place
- 1378 • Approvals may be valid for multiple releases but should not be valid indefinitely
- 1379 • Reliable records management for applications, approvals, and released data
- 1380 • Mechanisms for dispute resolution
- 1381 • Timetable or service-level agreement (SLA)
- 1382 • Legal and technical understanding of privacy

1383 Example outputs of a DRB include specifying access methods for different kinds of data
1384 releases, establishing acceptable levels of re-identification risk, and maintaining detailed
1385 records of previous data releases that ideally include the dataset that was released and the
1386 privacy-preserving methodology that was employed.

1387 There is some similarity between DRBs as envisioned here and the Institutional Review
1388 Board (IRBs) system created by the Common Rule¹⁷ for regulating human subject research
1389 in the United States. However, there are also important differences:

- 1390 • While the purpose of IRBs is to protect human subjects involved in human subject
1391 research, DRBs are charged with protecting data subjects, institutions, and – poten-
1392 tially – society as a whole.
- 1393 • Whereas IRBs are required to have “at least one member whose primary concerns
1394 are in nonscientific areas” and “at least one member who is not otherwise affiliated
1395 with the institution and who is not part of the immediate family of a person who is
1396 affiliated with the institution,” there does not appear to be a requirement for such
1397 members on a DRB.
- 1398 • Whereas IRBs give approval for research and then typically receive reports only dur-
1399 ing an annual review or when a research project terminates, DRBs may be involved
1400 at multiple points during the process.
- 1401 • Whereas approval of an IRB is required before research with human subjects can
1402 commence, DRBs are typically involved after research has taken place and prior to
1403 data or other research findings being released.
- 1404 • Whereas service on an IRB requires knowledge of the Common Rule and an under-
1405 standing of ethics, service on a DRB requires knowledge of statistics, computation,
1406 public policy, and some familiarity with the data being considered for release.

¹⁷The Federal Policy for the Protection of Human Subjects or the “Common Rule” was published in 1991 and codified in separate regulations by 15 federal departments and agencies. The Revised Common Rule was published in the Federal Register (FR) on January 19, 2017, and was amended to delay the effective and compliance dates on January 22, 2018, and June 19, 2018 [135].

3.7. De-Identification Standards

Agencies can rely on de-identification standards to provide standardized terminology, procedures, and performance criteria for de-identification efforts. Agencies can adopt existing de-identification standards or create their own. De-identification standards can be prescriptive or performance-based.

3.7.1. Benefits of Standards

De-identification standards assist agencies with the process of de-identifying data prior to public release. Without standards, data owners may be unwilling to share data, as they may be unable to assess whether a procedure for de-identifying data is sufficient to minimize privacy risk.

Standards can increase the availability of individuals with appropriate training by identifying a specific body of knowledge and practice that training should address. Absent standards, agencies may forego opportunities to share data. De-identification standards can help practitioners develop a community, as well as certification and accreditation processes.

Standards decrease uncertainty and provide data owners and custodians with best practices to follow. Courts can consider standards as acceptable practices that should generally be followed. In the event of litigation, an agency can point to the standard and say that it followed good data practice.

3.7.2. Prescriptive De-Identification Standards

A prescriptive de-identification standard specifies an algorithmic procedure that – if followed – results in data that are de-identified to an established benchmark.

The “Safe Harbor” method of the HIPAA Privacy Rule [3] is an example of a prescriptive de-identification standard. The intent of the Safe Harbor method is to “provide covered entities with a simple method to determine if the information is adequately de-identified” [89]. It does this by specifying that health information is considered to be de-identified through the removal of 18 kinds of identifiers and the assurance that the entity does not have *actual knowledge* that the remaining information can be used to identify an individual who is the subject of the information. Once de-identified, the dataset is no longer subject to HIPAA privacy, security, and breach notification regulations. Nevertheless, “a covered entity may require the recipient of de-identified information to enter into a data use agreement to access files with known disclosure risk” [89].

The Privacy Rule states that a covered entity that employs the Safe Harbor method must have no “actual knowledge” that the information – once de-identified – could still be used to re-identify individuals. However, covered entities are not obligated to employ experts or mount re-identification attacks against datasets to verify that the use of the Safe Harbor method has in fact resulted in data that cannot be re-identified.

Prescriptive standards have the advantage of being relatively easy for users to follow, but developing, testing, and validating such standards can be burdensome. Because prescriptive de-identification standards do not depend on the particulars of a specific case, there is a tendency for them to be more conservative than is necessary, resulting in an unnecessary decrease in data for corresponding levels of risk. Even so, there is no assurance that following a prescriptive standard actually produces the intended outcome.

Agencies that create prescriptive de-identification standards should ensure that data de-identified according to the standards have a sufficiently small risk of being re-identified consistent with the intended level of privacy protection. Such assurances frequently cannot be made unless formal privacy techniques, such as *differential privacy*, are employed. However, agencies may determine that public policy goals furthered by having an easy-to-use prescriptive standard outweighs the risk of a standard that does not have provable privacy guarantees.

Prescriptive de-identification standards carry the risk that the standard may not sufficiently de-identify to avoid the risk of re-identification, especially as methodology advances and more data sources become available.

A second risk when adopting prescriptive standards is that different agencies (or governments) may adopt inconsistent rules. In such a case, information that is legally de-identified for one purpose or in one jurisdiction may not be legally de-identified in another.

3.7.3. Performance-Based De-Identification Standards

Performance-based de-identification standards specify the properties that de-identified data must have. For example, under the “Expert Determination” method of the HIPAA Privacy Rule, a technique for de-identifying data is sufficient if an appropriate expert applying generally accepted statistical and scientific principles and methods “determines that the risk is very small that the information could be used, alone or in combination with other reasonably available information, by an anticipated recipient to identify an individual who is a subject of the information” [89]. The rule requires that experts document their methods and the results of their analyses.

Performance-based standards have the advantage of allowing users many different ways to solve a problem by leaving room for innovation. Another advantage is that they can require the desired outcome rather than specifying an aspirational mechanism.

Performance-based standards should be sufficiently detailed to perform in a manner that is reliable and repeatable. For example, standards that call for the use of experts can specify how an expert’s expertise should be determined. Standards that call for the reduction of risk to an acceptable level should provide a procedure for determining that level.

3.8. Education, Training, and Research

De-identifying data in a manner that preserves privacy can be a complex mathematical, statistical, administrative, and data-driven process. Frequently, the opportunities for identity disclosure will vary from dataset to dataset. Privacy-protecting mechanisms developed for one dataset may not be appropriate for others. For these reasons, agencies that engage in de-identification should ensure that their workers have adequate education and training in the subject domain. Agencies may wish to establish education or certification requirements for those who work directly with the datasets or to adopt industry standards such as the HITrust De-Identification Framework [77]. Because de-identification techniques are modality-dependent, agencies using de-identification may need to institute research efforts to develop and test appropriate data release methodologies.

3.9. Defense in Depth

In addition to de-identification, there are other technologies and methodologies that can secure sensitive data. Many of these approaches can complement de-identification and further reduce privacy risk to data subjects. Combining techniques is an example of *defense in depth* and should be considered whenever possible.

3.9.1. Encryption and Access Control

Encrypting sensitive data *at rest* can prevent attackers from obtaining the data directly (e.g., by compromising the server that stores it). Encryption can also serve as a form of *access control* (i.e., it can control *who* can access the data) because examining the data requires access to the encryption keys. If the original data (with identities) are retained, they should be stored encrypted, and access should be limited. Even after de-identification, more sensitive data not intended for public release can be provided to select individuals by limiting access via encryption.

3.9.2. Secure Computation

Two technologies enable computing on encrypted data *without decrypting it*:

1. **Fully-homomorphic encryption (FHE)** [55] allows a server to compute a function $f(x)$ on an encrypted value x without decrypting it. The result is a new encrypted value that can only be decrypted by someone who holds the original encryption key.
2. **Secure multi-party computation (MPC)** [74] allows multiple servers to *jointly* compute a function $f(x_1, \dots, x_k)$, where each server provides one of the inputs x_i , and no server learns any of the others' inputs.

Both of these approaches are general-purpose in that they can be used to compute any function, and both are considerably slower than performing the equivalent computation

with unencrypted data on a single computer. Nevertheless, both approaches are now sufficiently performant that they can be used for many practical kinds of privacy-preserving data analysis.¹⁸

3.9.3. Trusted Execution Environments

Trusted Execution Environments (TEEs) (also called *trusted hardware enclaves* or *secure hardware enclaves*) are another approach for computing on encrypted data. TEEs are implemented in computer hardware, typically within the silicon of a modern CPU, and protect programs that run on that CPU from the surrounding environment. For example, a TEE can cause data from a computer's CPU to be automatically encrypted when written to main memory and decrypted when read back to the CPU. In this way, data in memory are protected from other devices that can access memory, such as a network interface card. In addition to encryption, TEEs typically support attestation so that a program running on a TEE can attest to a remote system that the program is a true, legitimate, and faithful execution of the program.

Traditional cloud services require trusting the cloud provider, who may have a compromised environment (e.g., an operating system that records encryption keys). A TEE decreases the need for trust because it allows a user to validate that they are communicating with the remote program and offers assurance that no other program running in the cloud provider can access the program's data. Secure enclaves can thus be used to allow untrusted infrastructure to operate on sensitive data in much the same way as technologies like FHE and MPC.

Intel's Software Guard Extensions (SGX) [112], ARM's TrustZone [101], and AMD's Secure Encrypted Virtualization (SEV) [13] are all examples of secure hardware enclaves. All of these products are designed to provide similar security to cryptographic techniques while also providing performance similar to a single CPU operating on unencrypted data. These secure hardware products are necessarily complex, and various implementation errors have been discovered that can allow attackers to defeat their security protections. Secure hardware enclaves certainly offer increased security for data compared to plaintext computation, but agencies should carefully consider the trade-off between performance and security when choosing between secure hardware and cryptographic techniques.

3.9.4. Physical Enclaves

For extremely sensitive data, a *physical enclave* (see Section 3.4) may provide additional security. In this model, data are stored on a computer not connected to any network and are accessible only via physical access to a particular room. Access to the data is then controlled by limiting access to the room. This approach can be quite cumbersome.

¹⁸More information about these and other kinds of secure computation can be found on the NIST Privacy-Enhancing Cryptography (PEC) project website at <https://csrc.nist.gov/projects/pec>.

4. Technical Steps for Data De-Identification

The goal of de-identification is to transform data in a way that protects privacy while preserving the validity of inferences drawn on that data within the context of a target use-case. This section discusses technical options for performing de-identification and verifying the result of a de-identification procedure.

Agencies should adopt a detailed, written process for de-identifying data prior to commencing work on a de-identification project. The details of the process will depend on the particular de-identification approach that is pursued. In developing technical steps for data de-identification, agencies may wish to consider existing de-identification standards, such as the HIPAA Privacy Rule, the IHE De-Identification Handbook [61], or the HITRUST De-Identification Framework [77].

4.1. Determine the Privacy, Data Usability, and Access Objectives

Agencies intent on de-identifying data for release should understand the nature of the data that they intend to de-identify and determine the policies and standards that will be used to determine acceptable levels of data accuracy, de-identification, and the risk of re-identification. For example:

- Where did the data come from?
- What promises were made when the data were collected?
- What are the legal and regulatory requirements regarding data privacy and release?
- What is the purpose of the data release?
- What is the intended use of the data?
- What data-sharing model (Section 3.4) will be used?
- Which standards for privacy protection or de-identification will be used?
- What is the level of risk that the project is willing to accept?
- What are the goals for limiting re-identification? For example:
 - No one can be re-identified.
 - Only a few people can be re-identified.
 - Only a few people can be re-identified in theory, but no one will actually be re-identified in practice.
 - Only outliers can be re-identified.
 - Only people who are not outliers can be re-identified.

- 1578 – There is a small percentage chance of re-identification that is shared by every-
1579 one in the dataset.
- 1580 – There is a small percentage chance of re-identification, but some people in the
1581 dataset are significantly more likely to be re-identified, and the re-identification
1582 probability is somehow bounded.
- 1583 • What harm might result from re-identification, and what techniques will be used to
1584 mitigate those harms?
- 1585 • How should compliance with that level of risk be determined?

1586 Some goals and objectives are synergistic, while others are in opposition.

1587 4.2. Conducting a Data Survey

1588 Different kinds of data require different kinds of de-identification techniques. As a result,
1589 an important early step in the de-identification of government data is to identify the data
1590 modalities that are present in the dataset and formulate a plan for de-identification that takes
1591 into account goals for data release, data accuracy, privacy protection, and the best available
1592 science.

1593 For example:

- 1594 • **Tabular numeric and categorical data** is the subject of the majority of de-identification
1595 research and practice. These datasets are most frequently de-identified by using tech-
1596 niques based on the designation and removal of direct identifiers and the manipula-
1597 tion of quasi-identifiers. The chief criticism of de-identification based on direct and
1598 quasi-identifiers is that administrative determinations of quasi-identifiers may miss
1599 variables that can be uniquely identifying when combined and linked with external
1600 data, including data that are not available at the time the de-identification is per-
1601 formed but become available in the future.

1602 *K*-anonymity [122] is a common framework for performing and evaluating the de-
1603 identification of tabular numeric and categorical data. However, *risk determinations*
1604 *based on this kind of de-identification will be incorrect if direct and quasi-identifiers*
1605 *are not properly classified.* For example, if there exist quasi-identifiers that are not
1606 identified as such and not subjected to *k-anonymity*, then it may be easy to re-identify
1607 records in the de-identified dataset.

1608 Tabular data may also be used to create a synthetic dataset that preserves some infer-
1609 ence validity but does not have a one-to-one correspondence to the original dataset.

- 1610 • **Dates and times** require special attention when de-identifying because temporal in-
1611 formation is inherently linked to an external dataset: the natural progression of time.
1612 Some dates and times (e.g., February 22, 1732) are highly identifying, while others
1613 are not. Dates that refer to matters of public record (e.g., date of birth, death, or home

purchase) should be routinely taken as having high re-identification potential. Dates may also form the basis of linkages between dataset records or even within a record. For example, a record may contain the date of admission, the date of discharge, and the number of days in residence. Thus, care should be taken when de-identifying dates to locate and properly handle potential linkages and relationships. Applying different techniques to different fields may result in information being left in a dataset that can be used for re-identification. Specific issues regarding date de-identification are discussed in Section 4.3.4, “De-Identifying Dates.”

- **Geographic and map data** also require special attention when de-identifying, as some locations can be highly identifying, other locations are not identifying at all, and some locations are only identifying at specific times. As with dates and times, de-identifying geographic locations is challenging because locations inherently link to an external reality, and some locations during specific time periods are highly correlated with specific individuals (e.g., 38.8977° N, 77.0365° W). Identifying locations can be de-identified through the use of perturbation or generalization. The effectiveness of such de-identification techniques for protecting privacy in the presence of external information has not been well-characterized [51, p.37][115]. Specific issues regarding geographical de-identification are discussed in Section 4.3.5, “De-Identifying Geographical Locations.”
- **Unstructured text** may contain direct identifiers, such as a person’s name, or may contain additional information that can serve as a quasi-identifier. Finding such identifiers invariably requires domain-specific knowledge [51, p. 30]. Note that unstructured text may be present in tabular datasets and require special attention.¹⁹
- **Photos and video** may contain identifying information, such as printed names (e.g., name tags), as well as metadata in the file format. A range of biometric techniques also exists for matching photos of individuals against a dataset of photos and identifiers [51, p. 32].
- **Medical imagery** poses additional problems over photographs and video due to the presence of technical, medically specific information. For example, identifying information may be present in the image itself (e.g., a photo may show an identifying scar or tattoo), an identifier may be “burned in” to the image area (e.g., an identification plate containing a patient name that is included in an X-Ray), or an identifier may be present in the file metadata. The body part in the image itself may also be recognized using a biometric algorithm and dataset [51, p.35].
- **Genetic sequences** and other kinds of sequence information can be identified by using existing databanks that match sequences and identities. There is also evi-

¹⁹For an example of how unstructured text fields can damage the policy objectives and privacy assurances of a larger structured dataset, see Andrew Peterson’s article, “Why the names of six people who complained of sexual assault were published online by Dallas police” [98].

1650 dence that genetic sequences from individuals who are not in datasets can be matched
1651 through genealogical triangulation – a process that uses genetic information and other
1652 information as quasi-identifiers to single out a specific identity [51, p.36]. At present,
1653 there is no known method to reliably de-identify genetic sequences. Specific issues
1654 regarding the de-identification of genetic information is discussed in Section 4.3.6,
1655 “De-Identifying Genomic Information.”

1656 In many cases, data are complex and contain multiple modalities. Such mixtures may
1657 complicate risk determinations.

1658 **4.3. De-Identification by Removing Identifiers and Transforming Quasi-Identifiers**

1659 De-identification based on the removal of identifiers and the transformation of quasi-identifiers
1660 is one of the most common approaches currently in use. It has the advantage of being con-
1661 ceptually straightforward, and there is a long institutional history of using this approach
1662 within both federal statistical agencies and the healthcare industry. This approach has the
1663 disadvantage of not being based on formal methods for assuring privacy protection. The
1664 lack of formal methods does not mean that this approach cannot protect privacy, but it does
1665 mean that privacy protection is not assured.

1666 Below is a sample process for de-identifying data by removing identifiers and transforming
1667 quasi-identifiers.²⁰

1668 1. Determine the re-identification risk threshold. The organization determines accept-
1669 able risk for working with the dataset and possibly mitigating controls based on
1670 strong precedents and standards.²¹

1671 2. Determine the information in the dataset that could be used to identify the data sub-
1672 jects. Identifying information can include:

1673 **Direct identifiers** including names, phone numbers, and other information that un-
1674 ambiguously identifies an individual.

1675 **Quasi-identifiers** that could be used in a linkage attack. Typically, quasi-identifiers
1676 identify multiple individuals and can be used to triangulate a specific individual.

1677 **High-dimensional data** [10] that can be used to single out data records and thus
1678 constitute a unique pattern that could be identifying if the values exist in a
1679 secondary source to link against.²²

²⁰This protocol is based on a protocol developed by Professors Khaled El Emam and Bradley Malin [44].

²¹See the Federal Committee on Statistical Methodology’s Data Protection Toolkit at <https://nces.ed.gov/fcsm/dpt>.

²²For example, Narayanan and Shmatikov demonstrated that the set of movies that a person had watched could be used as an identifier given the existence of a second dataset of movies that had been publicly rated [84].

- 1680 3. Determine the direct identifiers in the dataset. An expert determines the elements in
1681 the dataset that only serve to identify the data subjects.
- 1682 4. Mask (transform) direct identifiers. The direct identifiers are either removed or re-
1683 placed with pseudonyms. Options for performing this operation are discussed in
1684 Section 4.3.1.
- 1685 5. Perform threat modeling. The organization determines the additional information
1686 they might be able to use for re-identification, including both quasi-identifiers and
1687 non-identifying values that a data intruder might use for re-identification.
- 1688 6. Determine minimal acceptable data accuracy. The organization determines what uses
1689 can or will be made with the de-identified data.
- 1690 7. Determine the transformation process that will be used to manipulate the quasi-
1691 identifiers. Pay special attention to the data fields that contain dates and geographical
1692 information, removing or recoding as necessary.
- 1693 8. Import (sample) data from the source dataset. Because the effort to acquire data from
1694 the source (identified) dataset may be substantial, some researchers recommend a test
1695 data import run to assist in planning [44].
- 1696 9. Review the results of the trial de-identification. Correct any coding or algorithmic
1697 errors that are detected.
- 1698 10. Transform the quasi-identifiers for the entire dataset.
- 1699 11. Evaluate the actual re-identification risk, which is calculated. As part of this evalua-
1700 tion, every aspect of the released dataset should be considered in light of the question,
1701 “Can *this* information be used to identify someone?”
- 1702 12. Compare the actual re-identification risk with the threshold specified by the policy-
1703 makers.
- 1704 13. If the data do not pass the actual risk threshold, adjust the procedure and repeat Steps
1705 11 and 12. For example, additional transformations may be required. Alternatively,
1706 it may be necessary to remove outliers. Removing data will of course impact data
1707 quality, but it will also protect the privacy of the individuals whose data has been
1708 removed.

1709 4.3.1. Removing or Transforming of Direct Identifiers

1710 There are many possible processes for removing direct identifiers from a dataset, including:

- 1711 • **Removal and replacement.** Replace identifiers with the value used by the database
1712 to indicate a missing value, such as NULL or NA.
- 1713 • **Masking.** Replace identifiers with a repeating character, such as XXXXXX or 999999.

- 1714 • **Encryption.** Encrypt the identifiers with a strong encryption algorithm. After en-
1715 cryptation, the key can be discarded the cryptographic key to prevent decryption. How-
1716 ever, if there is a desire to employ the same transformation at a later point in time,
1717 the key should not be discarded but rather stored in a secure location separate from
1718 the de-identified dataset. Encryption used for this purpose carries special risks that
1719 need to be addressed with specific controls (see Section 4.3.2 below for further in-
1720 formation).
 - 1721 • **Hashing with a keyed hash.** A keyed hash is a special kind of hash function that
1722 produces different hash values for different keys. The hash key should have sufficient
1723 randomness to defeat a brute force attack aimed at recovering the hash key (e.g.,
1724 SHA-256 HMAC [20] with a 256-bit randomly generated key). As with encryption,
1725 the key should be discarded unless there is a desire for repeatability. Hashing used
1726 for this purpose carries special risks that need to be addressed with specific controls
1727 (see Section 4.3.2 below for further information).
 - 1728 • **Replacement with keywords.** This approach transforms identifiers such as George
1729 Washington to PATIENT. Note that some keywords may be equally identifying, such
1730 as transforming George Washington to PRESIDENT.
 - 1731 • **Replacement with realistic surrogate values.** This approach transforms identifiers
1732 such as George Washington to surrogates that blend in, such as Abraham Polk.²³
- 1733 Encryption, hashing with a keyed hash, and replacement with realistic surrogate values are
1734 pseudonymization techniques. The technique used to remove direct identifiers should be
1735 clearly documented for users of the dataset – especially if the technique of replacement by
1736 realistic surrogate names is used – so that future data users have documentation that the
1737 dataset has been de-identified.
- 1738 If the agency plans to make data available for longitudinal research and contemplates mul-
1739 tiple data releases, then the transformation process should be repeatable, and the resulting
1740 transformed identities should be *pseudonyms*. The mapping between the direct identifier
1741 and the pseudonym is performed using a lookup table or a repeatable transformation. In
1742 either case, the release of the lookup table or the information used for the repeatable trans-
1743 formation will result in compromised identities. Thus, the lookup table or the information
1744 for the transformation must be highly protected. When using a lookup table, the pseudonym
1745 must be randomly assigned.
- 1746 A significant risk of using a repeatable transformation is that a data intruder may be able
1747 to determine the transformation and – in so doing – gain the capability to re-identify all of
1748 the records in the dataset.

²³A study by Carrell et. al found that using realistic surrogate names in de-identified text like John Walker and 3900 Pennsylvania Ave instead of generic labels like PATIENT and ADDRESS could decrease or mitigate the risk of re-identifying the few names that remained in the text because “the reviewers were unable to distinguish the residual (leaked) identifiers from the...surrogates” [21].

When multiple organizations use the same pseudonymization scheme, they can trade data and perform matching on the pseudonyms. However, this practice also allows the organizations to re-identify each other's shared datasets. As an alternative, organizations can participate in a *private set intersection protocol*, of which there are many in the cryptographic literature [78, 34, 69].

4.3.2. Special Security Note Regarding the Encryption or Hashing of Direct Identifiers

The transformation of direct identifiers through encryption or hashing carries special risks, as errors in procedure or the release of the encryption key can compromise identities for the entire dataset.

When information is protected with encryption, the security of the encrypted data depends entirely on the security of the encryption key. If a key is improperly chosen, it may be possible for a data intruder to discover the key using a brute force search. Because there is no visual difference between data that are encrypted with a strong encryption key and data that are encrypted with a weak key, organizations must utilize administrative controls to ensure that keys are both unpredictable and suitably protected. The use of encryption or hashing to protect direct identifiers is, therefore, *not recommended*.

4.3.3. De-Identifying Numeric Quasi-Identifiers

Once a determination is made regarding quasi-identifiers, they should be transformed. A variety of techniques are available to transform quasi-identifiers:

- **Top and bottom coding.** Outlier values that are above or below certain values are coded appropriately. For example, the HIPAA Privacy Rules calls for ages over 89 to be “aggregated into a single category of age 90 or older” [132, § 164.514 (b)].
- **Micro aggregation.** Individual microdata are combined into small groups that preserve some data analysis capability while providing for some disclosure protection [110].
- **Generalize categories with small counts.** When preparing contingency tables, several categories with small values may be combined. For example, rather than reporting that there is one person with blue eyes, two people with green eyes, and one person with hazel eyes, it may be reported that there are four people with blue, green, or hazel eyes.
- **Data suppression.** Cells in contingency tables with counts lower than a predefined threshold can be suppressed to prevent the identification of attribute combinations with small numbers [141].
- **Blanking and imputing.** Specific values that are highly identifying can be removed and replaced with imputed values.

- **Attribute or record swapping.** Attributes or data values are swapped within a set of similar records. For example, data that represent families in two similar towns within a county might be swapped with each other. “Swapping has the additional quality of removing any 100-percent assurance that a given record belongs to a given household” [130, p.31] while preserving the accuracy of regional statistics, such as sums and averages. In this case, the average number of children per family in the county would be unaffected by data swapping. However, swapping may damage or destroy important relationships within the data and introduce systematic biases, depending on how the swapping candidates are selected.
- **Noise infusion.** Also called “partially synthetic data,” this approach adds small random values to attributes. For example, instead of reporting that a person is 84 years old, the person may be reported as being 79 years old. Noise infusion increases variance in reported statistics and leads to attenuation bias in estimated regression coefficients and correlations among attributes [36, 7]. When combined with a requirement for non-negative reporting of attributes, such as age or population, noise infusion also introduces systematic bias since more values are increased in value than decreased.

These techniques (and others) are described in detail in several publications, including:

- **Statistical Policy Working Paper #22.** (Second version, 2005) by the Federal Committee on Statistical Methodology [47]. This 137-page paper includes worked examples of disclosure limitation, specific recommended practices for federal agencies, profiles of federal statistical agencies conducting disclosure limitation, and an extensive bibliography. This document has been superseded by the Data Protection Toolkit.
- **The Data Protection Toolkit (BETA).** A website maintained by the Federal Committee on Statistical Methodology for the purpose of promoting data access while protecting confidentiality throughout the federal statistical system [48]. <https://nces.ed.gov/fcsm/dpt>
- **The Anonymisation Decision-Making Framework.** By Mark Elliot, Elaine MacKey, Kieron O’Hara and Caroline Tudor, UKAN, University of Manchester, Manchester, UK, 2016. This 156-page book provides tutorials and worked examples for de-identifying data and calculating risk.
- **IHE IT Infrastructure Handbook: De-Identification.** (Integrating the Healthcare Enterprise, June 6, 2014) IHE offers a variety of guides, including one on de-identification at http://www.ihe.net/User_Handbooks/.

Swapping and noise infusion both introduce noise into the dataset, such that records literally contain incorrect data. Certain kinds of noise infusion have been mathematically proven to provide formal privacy guarantees. Swapping has no such guarantees.

All of these techniques impact data accuracy, but whether they impact data *utility* depends on the downstream uses of the data. For example, top-coding household incomes will not impact a measurement of the 90-10 quantile ratio, but it will impact a measurement of the top 1% of household incomes [99].

Prior to the adoption of differential privacy by the U.S. Census Bureau, federal statistical agencies largely did not document the specific statistical disclosure techniques they used when performing statistical disclosure limitation. Likewise, statistical agencies did not document the parameters used in the transformations nor the amount of data that have been transformed, as documenting these techniques can allow a data intruder to reverse-engineer the specific values, eliminating privacy protection [7]. This lack of transparency sometimes resulted in erroneous conclusions on the part of data users. This is another example of why it is important for documentation of the de-identification process to accompany the release of de-identified data and is one of the motivations for the U.S. Census Bureau to adopt data privacy techniques that do not rely on secrecy for their effectiveness [56, 121, 58, 6].

4.3.4. De-Identifying Dates

Dates can exist in many ways in a dataset. Dates may be in particular kinds of typed columns, such as a date of birth or the date of an encounter. Dates may be present as a number, such as the number of days since an epoch like January 1, 1900. Dates may be present in the free text narratives or in photographs (e.g., a photograph that shows a calendar or a picture of a computer screen with date information).

Several strategies have been developed for de-identifying dates:

- Under the HIPAA Privacy Rule, dates must be generalized to no greater specificity than the year (e.g., July 4, 1776, becomes 1776).
- Dates within a single person's record can be systematically adjusted by a random amount. For example, dates of a hospital admission and discharge might be systematically moved the same number of days – a date of admission and discharge of July 4, 1776, and July 9, 1776, become Sept. 10, 1777, and Sept. 15, 1777 [89]. However, this does not eliminate the risk that a data intruder will make inferences based on the interval between dates.
- In addition to a systematic shift, the intervals between dates can be perturbed to protect against re-identification attacks that involve identifiable intervals while still maintaining the order of events.
- Some dates cannot be arbitrarily changed without compromising data accuracy. For example, it may be necessary to preserve the day of the week, whether a day is a workday or a holiday, or a relationship to a holiday or event.
- Some ages can be randomly adjusted without impacting data accuracy while others cannot. For example, in many cases, the age of an individual can be randomly ad-

1860 justed ± 2 years if the person is over the age of 25 but not if their age is between
1861 one and three. However, individuals become eligible for specific benefits at specific
1862 ages, such as Social Security retirement at age 62, so changes to ages around these
1863 milestones may also result in data accuracy problems.

1864 **4.3.5. De-Identifying Geographical Locations and Geolocation Data**

1865 Geographical data can exist in many ways in a dataset. Geographical locations may be
1866 indicated by map coordinates (e.g., 39.1351966, -77.2164013), a street address (e.g., 100
1867 Bureau Drive), or a postal code (e.g., 20899). Geographical locations can also be embedded
1868 in textual narratives.

1869 Some geographical locations are not identifying (e.g., a crowded train station), while others
1870 may be highly identifying (e.g., a house in which a single person lives). Other locations
1871 may be identifying at some times of day and not others or during some months or some
1872 years. The amount of noise required to de-identify geographical locations significantly
1873 depends on the availability of external data, including geographical surveys. Identity may
1874 be shielded in an urban environment by adding ± 100 m, whereas a rural environment may
1875 only require ± 5 km or more to introduce sufficient ambiguity.

1876 A prescriptive de-identification rule – even one that accounts for varying population densi-
1877 ties – may still be insufficient for de-identification if the rule fails to consider the interaction
1878 between geographic locations and other quasi-identifiers in the dataset. Noise should be
1879 added with caution to avoid the creation of inconsistencies in underlying data (e.g., mov-
1880 ing the location of a residence along a coast into a body of water or across geopolitical
1881 boundaries).

1882 Single locations may become identifying if they represent locations linked to a single in-
1883 dividual that are recorded over time (e.g., a work/home commuting pair). Such behavioral
1884 time-location patterns can be quite distinct and allow for re-identification even with a small
1885 number of recorded locations per individual [82, 81]. Research in 2021 concluded that
1886 “[t]he risk of re-identification remains high even in country-scale location datasets” [46].

1887 Data that are of higher resolution are typically more identifying. For example, in July
1888 2021, the Catholic publication *The Pillar* published a report in which it had purchased
1889 the de-identified geolocation information for users of a homosexual dating platform. With
1890 this data, the journalists identified a prominent Catholic official as a user of the platform
1891 by simply matching the geolocation data to the official’s official residence. The official
1892 promptly resigned [100].

1893 **4.3.6. De-Identifying Genomic Information**

1894 Deoxyribonucleic acid (DNA) is the molecule inside human cells that carries genetic in-
1895 structions used for the proper functioning of living organisms. DNA present in the cell

1896 nucleus is inherited from both parents, while DNA present in the mitochondria is only
1897 inherited from an organism's mother.

1898 DNA is a repeating polymer that is made from four chemical bases: adenine (A), guanine
1899 (G), cytosine (C), and thymine (T). Human DNA consists of roughly 3 billion bases, of
1900 which 99% are the same in all people [54]. Modern technology allows for the complete
1901 specific sequence of an individual's DNA to be chemically determined, although this is
1902 rarely done in practice. With current technology, it is far more common to use a DNA
1903 microarray to probe for the presence or absence of specific DNA sequences at predeter-
1904 mined points in the genome. This approach is typically used to determine the presence
1905 or absence of specific single nucleotide polymorphisms (SNPs) [53]. DNA sequences and
1906 SNPs are the same for monozygotic (identical) twins, individuals resulting from divided
1907 embryos, and clones. With these exceptions, it is believed that no two humans have the
1908 same complete DNA sequence.

1909 Individual SNPs may be shared by many individuals, but a sufficiently large number of
1910 SNPs that show sufficient variability is generally believed to produce a combination that
1911 is unique to an individual. Thus, there are some sections of the DNA sequence and some
1912 combinations of SNPs that have high variability within the human population and oth-
1913 ers that have significant conservation between individuals within a specific population or
1914 group. When there is high variability, DNA sequences and SNPs can be used to match an
1915 individual with a historical sample that has been analyzed and entered into a dataset. The
1916 inheritability of genetic information has also allowed researchers to determine the surnames
1917 and even the complete identities of some individuals [57].

1918 As the number of individuals who have their DNA and SNPs measured increases, scientists
1919 are realizing that the characteristics of DNA and SNPs in individuals may be more com-
1920 plicated than the preceding paragraphs imply. DNA changes as individuals age because of
1921 senescence, transcription errors, and mutation. DNA methylation, which can impact the
1922 functioning of DNA, also changes over time [17]. Individuals who are made up of DNA
1923 from multiple individuals – typically the result of the fusion of twins in early pregnancy
1924 – are known as *chimera* or *mosaic*. In 2015, a man in the United States failed a paternity
1925 test because the genes in his saliva were different from those in his sperm [68]. A hu-
1926 man chimera was identified in 1953 because the person's blood contained a mixture of two
1927 blood types: A and O [37]. The incidence of human chimeras is unknown.

1928 Because of the high variability inherent in DNA, complete DNA sequences may be iden-
1929 tifiable by linking with an external dataset. Likewise, biological samples for which DNA
1930 can be extracted may be identifiable. Subsections of an individual's DNA sequence and
1931 collections of highly variable SNPs may be identifiable unless it is known that there are
1932 many individuals who share the region of DNA or those SNPs. Furthermore, genetic infor-
1933 mation may not only identify an individual but could also identify an individual's ancestors,
1934 siblings, and descendants.

Reading Level at Start of School Year	# of Students
Below grade level	30-39
At grade level	50-59
Above grade level	20-29

Table 1. Reading levels at a hypothetical school, as measured by entrance examinations, reported at the start of the school year on October 1.

Reading Level at Start of School Year	# of Students
Below grade level	30-39
At grade level	50-59
Above grade level	30-39

Table 2. Reading levels at a hypothetical school, as measured by entrance examinations, reported one month into the school year on November 1 after a new student has transferred to the school.

4.3.7. De-Identifying Text Narratives and Qualitative Information

Researchers must devote specific attention when they de-identify text narratives and other kinds of qualitative information. Many approaches developed in the 1980s and 1990s that provided reasonable privacy assurances at the time may no longer provide adequate protection in an era with high-quality internet search and social media [96, 97]. This is an area of active research.

4.3.8. Challenges Posed by Aggregation Techniques

Aggregation does not necessarily provide privacy protection, especially when data are presented in multiple data releases. Consider a hypothetical example of a school that reports on its website the number of students performing below, at, and above grade level at the start of the school year (table 1). Then consider that a new student enrolls at the school on October 15, and the school updates the table on its website (table 2).

By comparing the two tables, it is possible to infer that the student who joined the school is likely performing above grade level. This reveals protected information. Moreover, if a person who views both tables knows the specific student who enrolled in October, they have learned a private fact about that student.

Aggregation does not inherently protect privacy, and thus aggregation alone is not sufficient to provide formal privacy guarantees. However, the differential privacy literature does provide methods for performing aggregation that are both formally private and highly accurate when applied to large datasets. These methods work through the addition of carefully calibrated noise.

4.3.9. Challenges Posed by High-Dimensional Data

Even after removing all of the unique identifiers and manipulating the quasi-identifiers, data can still be identifying if it is of sufficiently high dimensionality and if there exists a way to link the supposedly non-identifying values to an identity.²⁴

4.3.10. Challenges Posed by Linked Data

Data can be linked in many ways. Pseudonyms allow data records from the same individual to be linked together over time. Family identifiers allow data from parents to be linked with their children. Device identifiers allow data to be linked to physical devices and potentially link together all data coming from the same device. Data can also be linked to geographical locations.

Data linkage increases the risk of re-identification by providing more attributes that can be used to distinguish the true identity of a data record from others in the population. For example, survey responses that are linked together by household are more readily re-identified than survey responses that are not linked. Heart rate measurements may not be considered identifying, but given a long sequence of tests, each individual in a dataset would have a unique constellation of heart rate measurements, and the dataset could be susceptible to being linked with another dataset that contains the same values.²⁵ Geographical location data can – when linked over time – create individual behavioral time-location patterns that can be used to classify and identify unlabeled data, even with a small number of recorded locations per individual [82, 81].

Dependencies between records may result in record linkages even when there is no explicit linkage identifier. For example, it may be that an organization has new employees take a proficiency test within seven days of being hired. This information would allow links to be drawn between an employee dataset that accurately reported an employee's start date and a training dataset that accurately reported the date that the test was administered, even if the sponsoring organization did not intend for the two datasets to be linkable.

4.3.11. Challenges Posed by Composition

In computer science, the term *composition* refers to combining multiple functions to create more complicated ones. One of the defining characteristics of complex systems is that they have unpredictable behavior, even when they are composed of very simple components. A challenge of composition is to develop approaches for limiting or eliminating such unpredictable behavior. Typically, this is done by proactively limiting the primitives that can be

²⁴For example, consider a dataset of an anonymous survey that links together responses from parents and their children. In such a dataset, a child might be able to find their parents' confidential responses by searching for their own responses and then following the link [84].

²⁵This is a different approach than characterizing an individual's heartbeat pattern so that it can be used as a biometric. In this case, it is a specific sequence of heartbeats that is recognized.

composed. De-identification is such a primitive that statisticians and data scientists must carefully control to ensure that the results of de-identification efforts can be composed. Without such controls, the results of composition can become unpredictable.

Specifically, it is important to understand whether the techniques used for de-identifying will retain their privacy guarantees when they are subject to composition. For example, if the same dataset is made available through two different de-identification regimes, what will happen to the privacy guarantees if the two downstream datasets are recombined? One of the primary advantages of differential privacy is that its operators are composable. This is not true of most other de-identification techniques.

Composition concerns can arise when:

- The same dataset is provided to multiple downstream users.
- Snapshots of a dataset are published on a periodic basis.
- Changes in computer technology result in new aspects of a dataset being made available.
- Legal proceedings require that aspects of the dataset (attributes or a subset of records) are made available without de-identification.

Privacy risk can result from unanticipated composition, which is one of the reasons that released datasets should be subjected to periodic review and reconsideration.

4.3.12. Potential Failures of De-Identification

The de-identification process outlined in this section can fail to prevent a disclosure for a number of different reasons. In addition, failures of data *utility* can also occur, in which the de-identification process removes *too much* information, and the de-identified dataset is not useful for its intended purpose.

- If an **inappropriate risk threshold** is selected, then the risk of re-identification may be higher than intended. Agencies should select risk thresholds conservatively to address this issue.
- If **direct or quasi-identifiers are missed**, then identifying information may remain in the de-identified dataset, leading to increased re-identification risk. Agencies should be mindful of the ways in which personal information can be used to identify individuals and – in ambiguous situations – assume that such information is identifying.
- If **threats are missed** during threat modeling, then the re-identification risk could be higher than intended. In particular, if other datasets that could be linked with the de-identified dataset are not considered, then the risk could be much higher than anticipated. Agencies should carefully consider existing and future data releases during threat modeling.

- If the **selected transformations fail to remove identifying information**, then the risk of de-identification could be higher than intended. Agencies should select transformations with well-understood properties and a history of successful use.
- If the de-identified dataset **does not produce accurate results for its intended use**, then it may not satisfy the goals of the data release. Future data custodians may be forced to oversee additional data releases, and those future releases might be combined with the already released datasets in ways that are unforeseen. Agencies should understand how the de-identified data will be used and make sure to carefully evaluate its utility for those purposes before releasing it.

4.3.13. Post-Release Monitoring

Following the release of a de-identified dataset, the releasing agency should monitor it to ensure that the assumptions made during the de-identification remain valid. This is because the identifiability of a dataset can only increase over time. For example, the de-identified dataset may contain information that can be linked to an internal dataset that is later the subject of a data breach. In such a situation, the data breach could also result in the re-identification of the de-identified dataset. The de-identified dataset might also be linked to an external dataset released by a completely separate organization. Agencies have no control over the release of such datasets, and even monitoring may be challenging in this situation. In some cases, the de-identified dataset might be linked with privately held data, making monitoring impossible.

Agencies may wish to make releasing units responsible for post-release monitoring or to centralize the post-release monitoring in a single location. However, proper post-release monitoring requires knowledge of the datasets that have been released and the kinds of data that would allow for a re-identification attack. These requirements are likely to increase costs to organizations that wish to delegate post-release monitoring to other organizations or third parties. One way to decrease the requirement for post-release monitoring is to perform the de-identification using a formal privacy model (e.g., differential privacy) that provides for privacy without making assumptions about background information available to the data intruder.

4.4. Synthetic Data

An alternative to de-identifying using the technique presented in the previous section is to use the original dataset to create a synthetic dataset [35, p.8].

Synthetic data can be created by two approaches:

1. Sampling an existing dataset and either adding noise to specific cells likely to have a high risk of disclosure or replacing those cells with imputed values. This is known as a “partially synthetic” dataset (see Table 3).

Data adjective	Description
<i>Datasets without formal guarantees:</i>	
Partially synthetic	Data for which there may be one-to-one mappings between records in the original dataset and the synthetic dataset but for which some attributes may have been altered or swapped between records. This approach is sometimes called <i>blank-and-impute</i> .
<i>Datasets with formal guarantees if the original dataset is not used to create the data:</i>	
Test	Data that resemble the original dataset in terms of structure and the range of values but for which there is no attempt to ensure that inferences drawn on the test data will be like those drawn on the original data. Test data may also include extreme values that are not in the original data but are present for testing software.
Realistic	Data that have a characteristic that is like the original data but that is not developed by modifying original data and which contains no information that is privacy-sensitive.
<i>Datasets with formal guarantees when formal techniques are used:</i>	
Fully synthetic	Data for which there is no one-to-one mapping between any record in the original dataset and the synthetic dataset.

Table 3. Adjectives used for describing data in data releases.

2059 2. Using the existing dataset to create a model and then using that model to create a
2060 synthetic dataset. This is known as a “fully synthetic” dataset (see Table 3).

2061 In both cases, formal privacy techniques can be used to quantify the privacy protection
2062 offered by the synthetic dataset.

2063 4.4.1. Partially Synthetic Data

2064 A partially synthetic dataset is one in which some of the data have been altered from the
2065 original dataset using probabilistic models. For example, data that belong to two families
2066 in adjoining towns may be swapped to protect the identity of the families. Alternatively, the
2067 data for an outlier variable may be removed and replaced with a range value that is incorrect
2068 (e.g., replacing the value “60” with the range “30-35”). It is considered best practice for
2069 the data publisher to indicate that some values have been modified or otherwise imputed
2070 but not to reveal the specific values that have been modified.

4.4.2. Test Data

It is also possible to create *test data* that is syntactically valid but does not convey accurate information when analyzed. Such data can be used for software development. When creating test data, it is useful for the names, addresses, and other information in the data to be conspicuously non-natural so that the test data are not inadvertently confused with true confidential data. For example, use the name “FIRSTNAME1 LASTNAME2” rather than “JOHN SMITH.”

4.4.3. Fully Synthetic Data

A fully synthetic dataset is a dataset for which there is no one-to-one mapping between data in the original dataset and data in the de-identified dataset. One approach to creating a fully synthetic dataset is to use the original dataset to create a high-fidelity model and then to use a simulation to produce individual data elements that are consistent with the model. Special efforts must be taken to maintain marginal and joint probabilities when creating partially or fully synthetic data.

Fully synthetic datasets cannot provide more information to the downstream user than was contained in the original model. Nevertheless, some users may prefer to work with the fully synthetic dataset instead of the model for a variety of reasons:

- Synthetic data provides users with the ability to develop queries and other techniques that can be applied to the real data without exposing real data to users during the development process. The queries and techniques can then be provided to the data owner, who can run the queries or techniques on the real data and provide the results to the users.
- Many hypotheses not represented exactly in the original model may be informed by the synthetic data because they are correlated with hypotheses (effects) that are in the model.
- Some users may place more trust in a synthetic dataset than in a model.
- When researchers form their hypotheses from synthetic data and then verify their findings on actual data, they can be protected from pretest estimation and false-discovery bias [7, p.257].

Because of the possibility of false discovery, analysts should be able to validate their discoveries against the original data to ensure that the things they discover are in the original data and not artifacts of the data generation process.

Both high-fidelity models and synthetic data generated from models may leak personal information that is potentially re-identifiable. The amount of leakage can be controlled using formal privacy models (e.g., differential privacy) that typically involve the introduction

2106 of noise. Section 4.4.6 describes the construction of fully synthetic data with differential
2107 privacy.

2108 There are several advantages for agencies that choose to release de-identified data as a fully
2109 synthetic dataset:

- 2110 • It can be very difficult or even impossible to map records to actual people.
- 2111 • The privacy guarantees can potentially be mathematically established and proven (cf.
2112 the section below on “Creating a synthetic dataset with differential privacy”).
- 2113 • The privacy guarantees can remain in force even if there are future data releases.

2114 Fully synthetic data also have these disadvantages and limitations:

- 2115 • It is not possible to create pseudonyms that map back to actual people because the
2116 records are fully fabricated.
- 2117 • The data release may be less useful for accountability or transparency. For example,
2118 investigators equipped with a synthetic data release would be unable to find the actual
2119 “people” who make up the release because they would not actually exist.
- 2120 • It is difficult to find meaningful correlations or abnormalities in synthetic data that
2121 are not represented in the model. For example, if a model contains only main effects
2122 and first-order interactions, then all second-order interactions can only be estimated
2123 from the synthetic data to the extent that their design is correlated with the main or
2124 first-order interactions.
- 2125 • Users of the data may not realize that the data are synthetic. Simply providing doc-
2126 umentation that the data are fully synthetic may not be sufficient public notification
2127 since the dataset may be separated from the documentation. Instead, it is best to
2128 indicate in the data itself that the values are synthetic. For example, names like
2129 “SYNTHETIC PERSON” or “FIRSTNAME1 LASTNAME1” may be placed in the
2130 data.
- 2131 • Releasing a synthetic dataset may not be regarded by the public as a legitimate act of
2132 transparency, or the public may question the validity of the data based on its perceived
2133 lack of relationship to the original dataset. These concerns can be addressed with
2134 public education and by documenting the accuracy of the synthetic dataset.

2135 In addition, it can be extremely challenging to construct the high-fidelity models that enable
2136 good synthetic datasets. The best known techniques for constructing these models are
2137 designed around ensuring that specific properties of the data (e.g., correlations between
2138 certain data attributes) are preserved when the model is constructed. Models constructed
2139 this way may not necessarily reflect *other* properties that were present in the original data.

2140 It is often possible to construct very high-fidelity models when the desirable properties
2141 of the synthetic data are known in advance (e.g., when it is known what questions future

analysts will want to answer using the synthetic data). Constructing synthetic data that faithfully represents *all* properties of the original data is much more challenging.

4.4.4. Synthetic Data with Validation

Agencies that share or publish synthetic data can optionally provide a validation service that takes queries or algorithms developed with synthetic data and applies them to actual data. The results of these queries or algorithms can then be compared with the results of running the same queries on the synthetic data, and the researchers can be warned if the results are different. Alternatively, results can be provided to the researchers after the application of additional statistical disclosure limitation.

4.4.5. Synthetic Data and Open Data Policy

Releases of synthetic data can be confusing to the lay public.

- It may not be clear to data users that the synthetic data release is actually synthetic. Members of the public may assume instead that the synthetic data are simply an operational dataset that has had identifying columns suppressed.
- Synthetic data may contain synthetic individuals who appear similar to actual individuals in the population.
- Fully synthetic datasets do not have a zero-disclosure risk because they still contain information derived from non-public information about individuals. The disclosure risk may be greater when synthetic data are created with traditional statistical modeling or data imputing techniques rather than those based on formal privacy models, such as differential privacy, as the formal models have provisions for tracking the accumulated privacy loss that results from multiple data operations, as discussed in Section 4.4.6.

4.4.6. Creating a Synthetic Dataset with Differential Privacy

A growing number of mathematical algorithms have been developed for creating synthetic datasets that meet the mathematical definition of privacy provided by differential privacy [40]. Most of these algorithms will transform a dataset containing private data into a new dataset that contains synthetic data that nevertheless provides reasonably accurate results in response to a variety of queries. However, there is no algorithm or implementation currently in existence that can be used by a person who is unskilled in the area of differential privacy.

The idea of differential privacy is that the result of a data analysis function κ applied to a dataset should not change very much if an arbitrary person p 's data is added to or removed from a dataset D . That is, $\kappa(D) \approx \kappa(D - p)$. The degree to which the two values are approximately equal is determined by the privacy loss parameter ϵ .

In the mathematical formulation of differential privacy, ϵ can range from 0 to ∞ . When $\epsilon = 0$, the output of κ does not depend on the input dataset. When $\epsilon = \infty$, the output of κ is entirely dependent upon the input dataset, such that changing a single record results in an unambiguous measurable change in κ 's output. Thus, larger values of ϵ provide for more accuracy but result in increased privacy loss.

When ϵ is set appropriately, differential privacy limits the privacy loss that a data subject experiences from the use of their private data to the maximum privacy loss necessary for a given statistical purpose. Note that this particular notion of privacy does *not* protect all secrets about a person. It only protects the secrets that an observer would not have been able to learn if the person's data was not present in the dataset. Stated another way, differential privacy protects individuals from *additional harm resulting from their participation in the data analysis* but does not protect them from harm that would have occurred *even if their data were not present*. For example, if a study concludes that residents of Vermont overwhelmingly drive 4-wheel-drive vehicles, one might conclude that a *particular* Vermonter drives a 4-wheel-drive vehicle even if that individual did not participate in the study. Differential privacy does not attempt to prevent inferences of this type.

Many academic papers on differential privacy assume a value of 1.0 for ϵ but do not explain the rationale of the choice. Some researchers working in the field of differential privacy have tried mapping existing privacy regulations to the choice of ϵ , but these efforts invariably result in values of $\epsilon \neq 1$. Principled approaches for setting ϵ is a subject of current academic research [72].

There are relatively few scholarly publications regarding the deployment of differential privacy in real-world situations, and there are few papers that provide guidance in choosing appropriate values of ϵ . Thus, agencies that are interested in using differential privacy algorithms to allow for querying of sensitive datasets or the creation of synthetic data should ensure that the techniques are appropriately implemented and that the privacy protections are appropriate to the desired application.

4.5. De-Identifying with an Interactive Query Interface

Another model for granting public access to de-identified agency information is to construct an interactive query interface that allows members of the public or qualified investigators to run queries over the agency's dataset. This option has been developed by several agencies, and there are many ways that it can be implemented. For example:

- If the queries are run on actual data, the results can be altered through the injection of noise to protect privacy, potentially satisfying a formal privacy model such as differential privacy. Alternatively, individual queries can be reviewed by agency staff to verify that privacy thresholds are maintained.
- Queries can be run on synthetic data. In this case, the agency can also run queries on the actual data and warn the external researchers if the queries run on synthetic

2215 data deviate significantly from the queries run on the actual data (ensuring that the
2216 warning itself does not compromise the privacy of some individual).

- 2217 • Query interfaces can be made freely available on the public internet, or they can be
2218 made available in a restricted manner to qualified researchers operating in secure
2219 locations.

2220 A significant privacy risk with interactive queries is that each query results in additional
2221 privacy loss [33].²⁶ For this reason, query interfaces should also log both queries and
2222 query results in order to deter and detect malicious use.

2223 One of the advantages of synthetic data is that the privacy loss budget can be spent on
2224 creating the synthetic dataset rather than on responding to interactive queries.

2225 **4.6. Validating a De-Identified Dataset**

2226 Agencies should validate datasets after they are de-identified to ensure that the resulting
2227 dataset meets the agency's goals in terms of both data usefulness and privacy protection.

2228 **4.6.1. Validating Data Usefulness**

2229 De-identification decreases data accuracy and the usefulness of the resulting dataset. It is
2230 therefore important to ensure that the de-identified dataset is still useful for the intended
2231 application. Otherwise, there is no reason to go through the expense and added risk of
2232 de-identification.

2233 Several approaches exist for validating data usefulness. For example, insiders can perform
2234 statistical calculations on both the original dataset and the de-identified dataset and compare
2235 the results to see if the de-identification resulted in unacceptable changes. Agencies can
2236 engage trusted outsiders to examine the de-identified dataset and determine whether the
2237 data could be used for the intended purpose.

2238 Recognizing that there is an inherent trade-off between data accuracy and privacy protec-
2239 tion, agencies can adopt accuracy goals for the data that they make available to a broad
2240 audience. An accuracy goal specifies how accurate data must be in order to be fit for an
2241 intended use. Limiting data accuracy to this goal is an important technique for protecting
2242 the privacy of data subjects.

2243 **4.6.2. Validating Privacy Protection**

2244 Several approaches exist for validating the privacy protection provided by de-identification,
2245 including:

²⁶If a finite privacy loss budget is allocated, the data controller needs to respond by increasing the amount of noise added to each response, accepting a higher level of privacy risk, or ceasing to answer questions as the budget nears exhaustion. This can result in equity issues if the first users to query the dataset obtain better answers than later users.

- 2246 • Examining the resulting data files to make sure that no identifying information is
2247 unintentionally included in file data or metadata.
- 2248 • Examining the resulting data files to make sure that the data meet stated goals for
2249 ambiguity under a k -anonymity model, if such a standard is desired.
- 2250 • Critically evaluating all default assumptions used by software that performs data
2251 modification or modeling.
- 2252 • Conducting a *motivated intruder test* to see if reasonably competent outside indi-
2253 viduals can perform re-identification using publicly available datasets, commercially
2254 available datasets, or even private datasets that might be available to certain data
2255 intruders. Motivations for an intruder can include prurient interest, causing embar-
2256 rassment or harm, revealing private facts about public figures, or engaging in a rep-
2257 utation attack. Details for how to conduct a motivated intruder test can be found in
2258 *Anonymisation: Managing data protection risk code of practice*, published by the
2259 United Kingdom’s Information Commissioner’s Office [63].
- 2260 • Providing the team conducting the motivated intruder test with some confidential
2261 agency data to understand how a data intruder might be able to take advantage of
2262 data leaked as a result of a breach or a hostile insider.

2263 These approaches do not provide provable guarantees on the protection offered by de-
2264 identification, but they may be useful as part of an overall agency risk assessment.²⁷ Ap-
2265 plications that require provable privacy guarantees should rely on formal privacy methods,
2266 such as differential privacy, when planning their data releases.

2267 Validating the privacy protection of de-identified data is greatly simplified by using vali-
2268 dated de-identification software, as discussed in Section 5, “Evaluation.”

2269 4.6.3. Re-Identification Studies

2270 Re-identification studies are motivated intruder tests. These studies can identify issues that
2271 would allow external actors to successfully re-identify de-identified data. Re-identification
2272 studies look for vulnerabilities in a dataset that could be used for re-identifying data sub-
2273 jects. They do not determine whether someone with intimate knowledge of a specific re-
2274 spondent can find that respondent in the database. The only way to protect a single specific

²⁷ Although other documents that discuss de-identification use the term *risk assessment* to refer to a specific calculation of ambiguity using the k -anonymity de-identification model, this document uses the term *risk assessment* to refer to a much broader process. Specifically, risk assessment is defined as, “The process of identifying, estimating, and prioritizing risks to organizational operations (including mission, functions, image, reputation), organizational assets, individuals, other organizations, and the Nation, resulting from the operation of an information system. Part of risk management incorporates threat and vulnerability analyses and considers mitigations provided by security controls planned or in place. Synonymous with risk analysis” [23].

individual perceived to be at high risk of re-identification is through data perturbation (e.g., noise injection) or information reduction (e.g., removing the observation altogether).

The key statistic calculated in re-identification studies is the *conditional re-identification rate*. This statistic is a proxy for disclosure risk. The rate is the number of confirmed links between the dataset and another dataset divided by the number of putative (suspected) links, unduplicated by “defender” ID, expressed as a percentage. If the conditional re-identification rate falls above an agreed upon threshold for any publication strata, it suggests that the data should not be released outside of a controlled environment.

Re-identification studies are often an iterative process. If a re-identification study uncovers problems with the de-identified data, the data curator can engage with subject matter experts, make changes to the dataset, and perform another re-identification study. Changes to the dataset might involve coarsening linking variables, eliminating highly disclosive linking variables from the microdata to be released, or coarsening strata. This continues until the study concludes that the de-identified data can be disseminated.

There are two very different types of re-identification studies:

1. **Micro (or targeted) re-identification studies**, where one is looking for a specific person. A well-known example is that of former Governor William Weld of Massachusetts, whose medical records in a hospital discharge summary were record linked to voter records [16]. As noted earlier, individual targets are supremely hard to protect as there is often extensive publicly available information about them.
2. **Macro (or wholesale) re-identification studies**, where one seeks to embarrass or discredit the organization releasing the data. This is accomplished by linking easily procurable external intruder data to the protected microdata that are being released. Several metrics can be calculated to uncover putative links, and several methods can be used to confirm putative links. Python has record linkage objects that probabilistically link files using a wide variety of metrics.

Formal privacy parameters often appear opaque and elusive to non-theoreticians. Subject matter experts and decision-makers more clearly understand disclosure risk after reviewing the results of re-identification studies.

External intruders may calculate low or high suspected re-identification rates, given the information they have available to them. They may even purport to have successfully linked their external data to a de-identified dataset. By conducting a re-identification study *a priori*, those seeking to disseminate the de-identified data know how successful the external intruder’s re-identification attempt was if both parties have access to the same external internal data.

The conditional re-identification rate is identical to the metric of *precision* in the record linkage and health science literature. It represents the ratio of true positives to the sum of true positives and false positives. Data owners should not be alarmed if an external

organization reports a relatively high suspected re-identification rate as long as they know that the conditional re-identification rate is low [45, 59, 111].

Confirmed re-identification rates are defined in Section 3.2.1 as *re-identification probabilities*. On its own, a low confirmed re-identification probability does not indicate that an organization should disseminate a de-identified dataset. Even when a confirmed rate is low, a high *conditional* rate should direct an organization to not disseminate the de-identified microdata.

Re-identification studies may identify problems that can direct improvements to any organization's disclosure avoidance methods. Re-identification studies are not designed to replace legacy or modern provable privacy methods but to act as a quality control to validate that the methods – old and new – protect as they were designed.

5. Software Requirements, Evaluation, and Validation

Agencies should clearly define the requirements for de-identification algorithms and the software that implements those algorithms. They should be sure that the algorithms that they intend to use are validated, that the software implements the algorithms as expected, and that the data that result from the operation of the software are correct.

Today, there a growing number of algorithms and tools for performing de-identification, data masking, and performing other privacy-preserving operations. NIST maintains a list of some of these tools at <https://www.nist.gov/itl/applied-cybersecurity/privacy-engineering/collaboration-space/focus-areas/de-id/tools>. Such tools are also increasingly being evaluated in academic literature [116] and by NIST [107, 1], although there are no widely accepted performance standards or certification procedures at present.

5.1. Evaluating Privacy-Preserving Techniques

There have been decades of research in the field of statistical disclosure limitation and de-identification, and understanding in the field has evolved over time. Agencies should not base their technical evaluation of a technique solely on the fact that the technique has been published in peer-reviewed literature or that the agency has a long history of using the technique and has not experienced any problems. Instead, it is necessary to evaluate proposed techniques through the totality of scientific experience and with regard to current threats.

Traditional statistical disclosure limitation and de-identification techniques base their risk assessments – in part – on an expectation of what kinds of data are available to a data intruder to conduct a linkage attack. Where possible, these assumptions should be documented and published along with a description of the privacy-preserving techniques that were used to transform the datasets prior to release so that they can be reviewed by external experts and the scientific community.

Because our understanding of privacy technology and the capabilities of privacy attacks are rapidly evolving, techniques that have been previously established should be periodically reviewed. New vulnerabilities may be discovered in techniques that have been previously accepted. Alternatively, new techniques may be developed that allow agencies to re-evaluate the trade-offs they have made with respect to privacy risk and data usability.

5.2. De-Identification Tools

A de-identification tool is a program that is involved in the creation of de-identified datasets.

5.2.1. De-Identification Tool Features

De-identification tools may perform many functions, including:

- Detecting identifying information
- Calculating re-identification risk
- Performing de-identification
- Mapping identifiers to pseudonyms
- Providing for the selective revelation of pseudonyms

De-identification tools may handle a variety of data modalities. For example, tools may be designed for tabular data or for multimedia. Tools may attempt to de-identify all data types or be developed for specific modalities. A potential risk of using de-identification tools is that a tool could be equipped to handle some but not all of the different modalities in a dataset. For example, a tool could de-identify the categorical information in a table according to a de-identification standard but might not detect or attempt to address the presence of identifying information in a text field. For this reason, de-identification tools should be validated for the specific kinds of data that the agency intends to use.

5.2.2. Data Provenance and File Formats

Output files created by de-identification tools and data masking tools can record provenance information, such as metadata regarding input datasets, the de-identification methods used, and the resulting decrease in data accuracy. Output files can also be explicitly marked to indicate that they have been de-identified. For example, de-identification profiles that are part of the Digital Imaging and Communications in Medicine (DICOM) specification indicate which elements are direct versus quasi-identifiers and which de-identification algorithms have been employed [32, Appendix E, “Attribute Confidentiality Profiles”].

5.2.3. Data Masking Tools

Data masking tools are programs that can remove or replace designated fields in a dataset while maintaining relationships between tables. These tools can be used to remove direct

identifiers but generally cannot identify or modify quasi-identifiers in a manner consistent with a privacy policy or risk analysis.

Data masking tools were developed to allow software developers and testers access to datasets that contain realistic data while providing minimal privacy protection. Absent additional controls or data manipulations, data masking tools should not be used for the de-identification of datasets that are intended for public release nor as the sole mechanism to ensure confidentiality in non-public data sharing.

5.3. Evaluating De-Identification Software

Once techniques are evaluated and approved, agencies should ensure that the techniques are faithfully executed by their chosen software. Privacy software evaluation should consider the trade-off between data usability and privacy protection. Privacy software evaluation should also seek to detect and minimize the chances of tool error and user error.

For example, agencies should verify:

- **Correctness.** The software properly implements the chosen algorithms.
- **Containment.** The software does not leak identifying information in expected or unexpected ways, such as through the inaccuracies of floating-point arithmetic or the differences in execution time (if observable to a data intruder).
- **Usability.** The software can be operated efficiently and with minimal error, and users can detect and correct errors when they happen.

Agencies should also evaluate the performance of the de-identification software, such as:

- **Efficiency.** How long does it take to run on a dataset of a typical size?
- **Scalability.** How much does it slow down when moving from a dataset of N to $100N$?
- **Repeatability.** If the tool is run twice on the same dataset, are the results similar? If two different people run the tool, do they get similar results?

Ideally, software should be able to track the accumulated privacy leakage from multiple data releases.

5.4. Evaluating Data Accuracy

Finally, agencies should evaluate the accuracy of the de-identified data to verify that it is sufficient for the intended use. For example, researchers at MIT and Harvard applied k -anonymity de-identification to educational data collected by a massive open online course operated by MITx and HarvardX on the edX platform and found that de-identification resulted in meaningful biases that changed the meaning of some statistics. For example, in one case, de-identification decreased the number of enrolled female students from 29% to 26% because of the need to suppress attributes for specific microdata [30].

The field of statistical disclosure control has developed approaches for gauging the impact of SDC techniques on microdata [142]. The literature examines the mathematical impact of SDC procedures (e.g., sampling, recoding, suppression, rounding, and noise infusion) and computes the possible impact on various statistical measurements.

Approaches for evaluating data accuracy include [71]:

- Demonstrating that machine learning algorithms trained on the de-identified data can accurately predict the original data and vice versa
- Verifying that statistical distributions do not incur undue bias because of the de-identification procedure
- Publishing sufficient information about the statistical properties of the disclosure limitation methods to permit the correction of inferences using these properties

Agencies can create or adopt standards regarding the accuracy of de-identified data. If data accuracy cannot be well-maintained along with data privacy goals, then the release of data that is inaccurate for statistical analyses could potentially result in incorrect scientific conclusions and incorrect policy decisions.

6. Conclusion

Government agencies can use de-identification technology to make datasets available to researchers and the public without compromising the privacy of the people contained within the data.

There are currently three primary models available for de-identification:

1. agencies can make data available with traditional de-identification techniques that rely on the suppression of identifying information (direct identifiers) and the manipulation of information that partially identifies (quasi-identifiers);
2. agencies can create synthetic datasets; and
3. agencies can make data available through a query interface.

These models can be mixed within a single dataset to provide different kinds of access for different users or intended uses.

Privacy protection can be strengthened when agencies employ formal models for privacy protection, such as differential privacy, because the mathematical models that these systems use are designed to ensure privacy protection irrespective of future data releases or developments in re-identification technology. However, the mathematics underlying these systems is very new, and there is little experience within the Government in using these systems. Thus, agencies should understand the implications of these systems before deploying them in place of traditional de-identification approaches that do not offer formal privacy guarantees.

2451 Agencies that use de-identification should establish appropriate governance structures to
2452 support de-identification, data release, and post-release monitoring. Such structures will
2453 typically include a Disclosure Review Board as well as appropriate education, training,
2454 and research efforts.

2455 A summary of this document's advice for practitioners appears in Figure 5.

2456 In closing, it is important to remember that different jurisdictions may have different stan-
2457 dards and policies regarding the definition and use of de-identified data. Information that
2458 is considered de-identified in one jurisdiction may be regarded as being identifiable in an-
2459 other.

Governance and Management (Section 3) The management of de-identification includes identifying the goals of the de-identification process and considering risks to participants in the data release. To guide this process, this document describes several tools:

- Consider all phases of the **Data Life Cycle** (Section 3.3).
- Consider different **Data Sharing Models** (Section 3.4), including complementary protections like Data Use Agreements, Synthetic Data, and Enclaves.
- Leverage the **Five Safes** (Section 3.5), a methodology for evaluating risk.
- Form a **Disclosure Review Board** (Section 3.6) to oversee the implementation of de-identification policies.
- Follow existing **de-identification standards** when possible (Section 3.7).

Technical Steps (Section 4) The technical process of de-identification should leverage the best practices developed over the past several decades. In particular, NIST recommends that agencies:

- Conduct a **Data Survey** (Section 4.2) to identify de-identification requirements specific to the data.
- Determine **identifiers** and **quasi-identifiers** in the data, and select a method for de-identifying each one (Section 4.3).
- Consider the existing **auxiliary data** (Section 4.3) that could be used to enable a re-identification attack.
- Practice **defense in depth** by combining security measures with de-identification when possible, and consider using **Synthetic Data** (Section 4.4) or an **Interactive Query Interface** (Section 4.5).
- When possible, use **formal privacy techniques** to quantify privacy loss associated with the release of de-identified data (Section 4.4.6).
- Validate the **utility** and **privacy** of the de-identified data (Section 4.6). In particular, establish accuracy goals for de-identification data so that the data is not more accurate than required for the intended purpose.

Software (Section 5) In general, agencies should:

- Utilize automated, repeatable, software-based approaches for performing de-identification.
- Carefully consider the software used to implement de-identification to ensure that the algorithms used have been validated and that the software correctly implements those algorithms.
- Consider the efficiency, scalability, and repeatability properties of software tools, and evaluate the accuracy of the tool's output.

Fig. 5. Advice for Practitioners: A Summary

References

- [1] July 2020. URL: <https://www.nist.gov/blogs/cybersecurity-insights/differential-privacy-privacy-preserving-data-analysis-introduction-our>.
- [2] 115th Congress (2017–2018). *Public Law 115-435: The Foundations for Evidence-Based Policymaking Act of 2018*. 2018. URL: <https://www.congress.gov/bill/115th-congress/house-bill/4174>.
- [3] *45 CFR 164 Health Insurance Portability and Accountability Act of 1996 (HIPAA) Privacy Rule Safe Harbor method Standard: De-identification of protected health information*.
- [4] *81 FR 49689: Revision of OMB Circular No. A-130, "Managing Information as a Strategic Resource"*. July 2016. URL: <https://www.cio.gov/policies-and-priorities/circular-a-130/>.
- [5] 95th Congress. *Public Law 95-416*. Oct. 1978. URL: https://www.census.gov/history/pdf/NARA_Legislation.pdf.
- [6] John Abowd et al. "The 2020 Census Disclosure Avoidance System TopDown Algorithm". In: *Harvard Data Science Review* Special Issue 2 (June 2022). URL: <https://hdsr.mitpress.mit.edu/pub/7evz361i>.
- [7] John M. Abowd and Ian M. Schmutte. *Economic Analysis and Statistical Disclosure Limitation*. Mar. 2015. URL: <https://www.brookings.edu/bpea-articles/economic-analysis-and-statistical-disclosure-limitation/>.
- [8] John M. Abowd and Lars Vilhuber. "How Protective are Synthetic Data?" In: *Lecture Notes in Computer Science: Privacy in Statistical Databases* 5262 (2008), pp. 239–246.
- [9] "Accuracy". In: *Glossary of Statistical Terms* (Sept. 2001). Last accessed June 23, 2022. URL: <https://stats.oecd.org/glossary/detail.asp?ID=21>.
- [10] Charu C. Aggarwal. "On K-Anonymity and the Curse of Dimensionality". In: *Proceedings of the 31st International Conference on Very Large Data Bases*. VLDB '05. Trondheim, Norway: VLDB Endowment, 2005, pp. 901–909. ISBN: 1595931546.
- [11] J. Trent Alexander, Michael Davern, and Betsey Stevenson. "Inaccurate age and sex data in the census PUMS files: Evidence and implications". In: *Public Opinion Quarterly* 74 (3 2010), pp. 551–569. URL: <https://doi.org/10.1093/poq/nfq033>.
- [12] Micah Altman et al. "Towards a Modern Approach to Privacy-Aware Government Data Releases". In: *Berkeley Technology Law Journal* 30 (3), pp. 1967–2072. URL: http://papers.ssrn.com/sol3/papers.cfm?abstract_id=2779266.
- [13] AMD. *AMD Secure Encrypted Virtualization (SEV)*. Last accessed July 13, 2022. URL: <https://developer.amd.com/sev/>.

- 2496 [14] Olivia Angiuli, Joe Blitzstein, and Jim Waldo. “How to De-Identify Your Data”.
2497 In: *Commun. ACM* 58.12 (Nov. 2015), pp. 48–55. ISSN: 0001-0782. DOI: [10.1145/](https://doi.org/10.1145/2814340)
2498 [2814340](https://doi.org/10.1145/2814340). URL: <https://doi.org/10.1145/2814340>.
- 2499 [15] ASTM International. *ASTM E1869-04 (Reapproved 2014) Standard Guide for Con-*
2500 *fidentiality, Privacy, Access, and Data Security Principles for Health Information*
2501 *Including Electronic Health Records*. 2014.
- 2502 [16] Daniel Barth-Jones. *The ‘Re-Identification’ of Governor William Weld’s Medical*
2503 *Information: A Critical Re-Examination of Health Data Identification Risks and*
2504 *Privacy Protections, Then and Now*. July 2012. URL: [https://ssrn.com/abstract=](https://ssrn.com/abstract=2076397%20or%20http://dx.doi.org/10.2139/ssrn.2076397)
2505 [2076397%20or%20http://dx.doi.org/10.2139/ssrn.2076397](https://ssrn.com/abstract=2076397%20or%20http://dx.doi.org/10.2139/ssrn.2076397).
- 2506 [17] Hans T Bjornsson et al. “Intra-individual Change Over Time in DNA Methylation
2507 with Familial Clustering”. In: *JAMA* 299 (24 June 2008), pp. 2877–2833. URL:
2508 <https://pubmed.ncbi.nlm.nih.gov/18577732/>.
- 2509 [18] Sylvia M. Burwell. *Open Data Policy-Managing Information as an Asset*. May
2510 2013. URL: [https://obamawhitehouse . archives . gov / sites / default / files / omb /](https://obamawhitehouse.archives.gov/sites/default/files/omb/memoranda/2013/m-13-13.pdf)
2511 [memoranda/2013/m-13-13.pdf](https://obamawhitehouse.archives.gov/sites/default/files/omb/memoranda/2013/m-13-13.pdf).
- 2512 [19] George Bush. *Executive Order 13402:Strengthening Federal Efforts to Protect Against*
2513 *Identity Theft*. May 2006. URL: [https://www.gpo.gov/fdsys/pkg/FR-2006-05-](https://www.gpo.gov/fdsys/pkg/FR-2006-05-15/pdf/06-4552.pdf)
2514 [15/pdf/06-4552.pdf](https://www.gpo.gov/fdsys/pkg/FR-2006-05-15/pdf/06-4552.pdf).
- 2515 [20] G. Camarillo, C. Holmberg, and Y. Gao. *Re-INVITE and Target-Refresh Request*
2516 *Handling in the Session Initiation Protocol (SIP)*. RFC 6141 (Proposed Standard).
2517 Internet Engineering Task Force, Mar. 2011. URL: www.ietf.org/rfc/rfc6141.txt.
- 2518 [21] David Carrell et al. “Hiding in plain sight: Use of realistic surrogates to reduce
2519 exposure of protected health information in clinical text”. In: *Journal of the Ameri-*
2520 *can Medical Informatics Association : JAMIA* 20 (2 July 2012), pp. 342–348. DOI:
2521 [10.1136/amiajnl-2012-001034](https://doi.org/10.1136/amiajnl-2012-001034).
- 2522 [22] Ann Cavoukian. *Privacy by Design: The 7 Foundational Principles*. Ontario, CA,
2523 Jan. 2011. URL: [https://www.ipc.on.ca/wp-content/uploads/Resources/7foundationalprinciples.](https://www.ipc.on.ca/wp-content/uploads/Resources/7foundationalprinciples.pdf)
2524 [pdf](https://www.ipc.on.ca/wp-content/uploads/Resources/7foundationalprinciples.pdf).
- 2525 [23] Jennifer Cawthra et al. *Securing Telehealth Remote Patient Monitoring Ecosystem*.
2526 2022. DOI: [10.6028/NIST.SP.1800-30](https://doi.org/10.6028/NIST.SP.1800-30). URL: [https://nvlpubs.nist.gov/nistpubs/](https://nvlpubs.nist.gov/nistpubs/SpecialPublications/NIST.SP.1800-30.pdf)
2527 [SpecialPublications/NIST.SP.1800-30.pdf](https://nvlpubs.nist.gov/nistpubs/SpecialPublications/NIST.SP.1800-30.pdf).
- 2528 [24] Census Bureau Data Stewardship Program. *DS025: Organization of the Disclosure*
2529 *Review Board*. Dec. 2019. URL: [https://www2.census.gov/foia/ds_policies/ds025.](https://www2.census.gov/foia/ds_policies/ds025.pdf)
2530 [pdf](https://www2.census.gov/foia/ds_policies/ds025.pdf).
- 2531 [25] Malcolm Chisholm. “7 Phases of a Data Life Cycle”. In: *Information Manage-*
2532 *ment* (July 2015). URL: [http://www.information-management.com/news/data-](http://www.information-management.com/news/data-management/Data-Life-Cycle-Defined-10027232-1.html)
2533 [management/Data-Life-Cycle-Defined-10027232-1.html](http://www.information-management.com/news/data-management/Data-Life-Cycle-Defined-10027232-1.html).

- 2534 [26] A. Church. “A Note on the ‘Entscheidungsproblem’”. In: *Journal of Symbolic*
2535 *Logic* 1 (1936), pp. 40–41.
- 2536 [27] Commission on Evidence-Based Policymaking. *The Promise of Evidence-Based*
2537 *Policymaking*. Sept. 2017. URL: [https://www.acf.hhs.gov/opre/project/commission-](https://www.acf.hhs.gov/opre/project/commission-evidence-based-policymaking-cep)
2538 [evidence-based-policymaking-cep](https://www.acf.hhs.gov/opre/project/commission-evidence-based-policymaking-cep).
- 2539 [28] Tore Dalenius. “Finding a Needle in a Haystack, or Identifying Anonymous Census
- 2540 Records”. In: *Journal of Official Statistics* 2 (3 1986), pp. 329–336.
- 2541 [29] Tore Dalenius. “Towards a methodology for statistical disclosure control”. In: *Statis-*
2542 *tik Tidskrift* 15 (1977), pp. 429–444.
- 2543 [30] Jon P. Daries et al. “Privacy, Anonymity, and Big Data in the Social Sciences”. In:
2544 *Communications of the ACM* 57 (6 Sept. 2014), pp. 56–63.
- 2545 [31] Tanvi Desai, Felix Ritchie, and Richard Welpton. Economics Working Paper Series
2546 1601. 2016. URL: <http://dx.doi.org/10.13140/RG.2.1.3661.1604>.
- 2547 [32] DICOM Standards Committee. *DICOM PS3.15 2016e — Security and System Man-*
2548 *agement Profiles*. 2016. URL: [http://dicom.nema.org/medical/dicom/current/](http://dicom.nema.org/medical/dicom/current/output/html/part15.html#chapter_E)
2549 [output/html/part15.html#chapter_E](http://dicom.nema.org/medical/dicom/current/output/html/part15.html#chapter_E).
- 2550 [33] Irit Dinur and Kobbi Nissim. “Revealing Information While Preserving Privacy”.
2551 In: *Proceedings of the Twenty-second ACM SIGMOD-SIGACT-SIGART Sympo-*
2552 *sium on Principles of Database Systems*. PODS ’03. San Diego, California: ACM,
2553 2003, pp. 202–210. ISBN: 1-58113-670-6. DOI: [10.1145/773153.773173](https://doi.org/10.1145/773153.773173). URL:
2554 [doi.acm.org/10.1145/773153.773173](https://doi.org/10.1145/773153.773173).
- 2555 [34] Changyu Dong, Liqun Chen, and Zikai Wen. “When Private Set Intersection Meets
- 2556 Big Data: An Efficient and Scalable Protocol”. In: *Proceedings of the 2013 ACM*
2557 *SIGSAC Conference on Computer and Communications Security*. CCS ’13. Berlin,
2558 Germany: Association for Computing Machinery, 2013, pp. 789–800. ISBN: 9781450324779.
2559 DOI: [10.1145/2508859.2516701](https://doi.org/10.1145/2508859.2516701). URL: <https://doi.org/10.1145/2508859.2516701>.
- 2560 [35] Jörg Drechsler, Stefan Bender, and Susanne Rässler. *Comparing fully and partially*
2561 *synthetic datasets for statistical disclosure control in the German IAB Establish-*
2562 *ment Panel (Working paper 11)*. New York, 2007. URL: [http://fdz.iab.de/342/](http://fdz.iab.de/342/section.aspx/Publikation/k080530j05)
2563 [section.aspx/Publikation/k080530j05](http://fdz.iab.de/342/section.aspx/Publikation/k080530j05).
- 2564 [36] George T. Duncan, Mark Elliot, and Juan-José Salazar-Gonzalez. *Statistical Confi-*
2565 *dentiality: Principles and Practice*. Springer, 2011, p. 113.
- 2566 [37] I. Dunsford et al. “A human blood-group chimera”. In: *British Medical Journal* 81
2567 (July 1953). URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2028470/>.
- 2568 [38] Cynthia Dwork and Aaron Roth. “The Algorithmic Foundations of Differential Pri-
- 2569 vacy”. In: *Foundations and Trends in Theoretical Computer Science*. Vol. 9. 3–4.
2570 NOW, 2014, pp. 211–407.

- 2571 [39] Cynthia Dwork et al. “Calibrating Noise to Sensitivity in Private Data Analysis”.
2572 In: *Theory of Cryptography*. Ed. by Shai Halevi and Tal Rabin. Berlin, Heidelberg:
2573 Springer Berlin Heidelberg, 2006, pp. 265–284. ISBN: 978-3-540-32732-5.
- 2574 [40] Cynthia Dwork et al. “Calibrating Noise to Sensitivity in Private Data Analysis”.
2575 In: *Proceedings of the Third Conference on Theory of Cryptography*. TCC’06. New
2576 York, NY: Springer-Verlag, 2006, pp. 265–284. ISBN: 3540327312. DOI: [10.1007/](https://doi.org/10.1007/11681878_14)
2577 [11681878_14](https://doi.org/10.1007/11681878_14). URL: doi.org/10.1007/11681878_14.
- 2578 [41] Department of Education (ED) Disclosure Review Board (DRB). *The Data Disclo-*
2579 *sure Decision*. Version 1.0. 2015. URL: [https://s3.amazonaws.com/sitesusa/wp-](https://s3.amazonaws.com/sitesusa/wp-content/uploads/sites/1151/2016/10/The-Data-Disclosure-Decision-Department-of-Education-Case-Study_Mar-2015.pdf)
2580 [content/uploads/sites/1151/2016/10/The-Data-Disclosure-Decision-Department-](https://s3.amazonaws.com/sitesusa/wp-content/uploads/sites/1151/2016/10/The-Data-Disclosure-Decision-Department-of-Education-Case-Study_Mar-2015.pdf)
2581 [of-Education-Case-Study_Mar-2015.pdf](https://s3.amazonaws.com/sitesusa/wp-content/uploads/sites/1151/2016/10/The-Data-Disclosure-Decision-Department-of-Education-Case-Study_Mar-2015.pdf).
- 2582 [42] Mark Elliot and Angela Dale. *Scenarios of attack: the data intruder’s perspective*
2583 *on statistical disclosure risk*. Spring 1999.
- 2584 [43] El Emam. *Methods for the de-identification of electronic health records for genomic*
2585 *research*. 2011. URL: [https://genomemedicine.biomedcentral.com/articles/10.](https://genomemedicine.biomedcentral.com/articles/10.1186/gm239)
2586 [1186/gm239](https://genomemedicine.biomedcentral.com/articles/10.1186/gm239).
- 2587 [44] K. El Emam and B. Malin. “Appendix B: Concepts and Methods for De-Identifying
2588 Clinical Trial Data”. In: *Sharing Clinical Trial Data: Maximizing Benefits, Mini-*
2589 *mizing Risk*. Washington, DC: Institute of Medicine of the National Academies,
2590 The National Academies Press, 2015.
- 2591 [45] Khaled El Emam and Luk Arbuckle. *Anonymizing Health Data: Case Studies and*
2592 *Methods to Get you Started*. Sebastopol, CA: O’Reilly Media, 2013.
- 2593 [46] Ali Farzanehfar, Florimond Houssiau, and Yves-Alexandre de Montjoye. “The risk
2594 of re-identification remains high even in country-scale location datasets”. In: *Pat-*
2595 *terns* 2.3 (2021), p. 100204. ISSN: 2666-3899. DOI: [https://doi.org/10.1016/j.](https://doi.org/10.1016/j.patter.2021.100204)
2596 [patter.2021.100204](https://doi.org/10.1016/j.patter.2021.100204). URL: [https://www.sciencedirect.com/science/article/pii/](https://www.sciencedirect.com/science/article/pii/S2666389921000143)
2597 [S2666389921000143](https://www.sciencedirect.com/science/article/pii/S2666389921000143).
- 2598 [47] Confidentiality and Data Access Committee. *Statistical Policy Working Paper 22:*
2599 *Report on Statistical Disclosure Limitation Methodology*. Tech. rep. Federal Com-
2600 mittee on Statistical Methodology, 2005. URL: [https://www.hhs.gov/sites/default/](https://www.hhs.gov/sites/default/files/spwp22.pdf)
2601 [files/spwp22.pdf](https://www.hhs.gov/sites/default/files/spwp22.pdf).
- 2602 [48] Federal Committee on Statistical Methodology. *Data Protection Toolkit*. Sept. 2020.
2603 URL: <https://nces.ed.gov/fcsm/dpt>.
- 2604 [49] Matthew Fredrikson et al. “Privacy in Pharmacogenetics: An End-to-End Case
2605 Study of Personalized Warfarin Dosing”. In: *23rd USENIX Security Symposium*.
2606 San Diego, CA. URL: [https://www.usenix.org/system/files/conference/usenixsecurity14/](https://www.usenix.org/system/files/conference/usenixsecurity14/sec14-paper-fredrikson-privacy.pdf)
2607 [sec14-paper-fredrikson-privacy.pdf](https://www.usenix.org/system/files/conference/usenixsecurity14/sec14-paper-fredrikson-privacy.pdf).

- [50] Simson Garfinkel. *De-Identification of Personally Identifiable Information*. Tech. rep. NIST IR 8053. National Institute of Science and Technology, Nov. 2015. URL: <http://nvlpubs.nist.gov/nistpubs/ir/2015/NIST.IR.8053.pdf>.
- [51] Simson L. Garfinkel. *De-identification of personal information*. 2015. DOI: 10.6028/NIST.IR.8053. URL: <https://nvlpubs.nist.gov/nistpubs/ir/2015/NIST.IR.8053.pdf>.
- [52] Simson L. Garfinkel. *Government Data De-Identification Stakeholder’s Meeting June 29, 2016 Meeting Report*. 2016. DOI: 10.6028/NIST.IR.8150. URL: <https://nvlpubs.nist.gov/nistpubs/ir/2016/NIST.IR.8150.pdf>.
- [53] Genetics Home Reference. *What are single nucleotide polymorphisms (SNPs)?* Last access June 16, 2022. 2022. URL: <https://ghr.nlm.nih.gov/primer/genomicresearch/snp>.
- [54] Genetics Home Reference. *What is DNA*. Last access June 16, 2022. 2022. URL: <https://ghr.nlm.nih.gov/primer/basics/dna>.
- [55] Craig Gentry. “A Fully Homomorphic Encryption Scheme”. AAI3382729. PhD thesis. Stanford, CA, USA, 2009. ISBN: 9781109444506.
- [56] Ruobin Gong, Erica L. Groshen, and Salil Vadhan. “Harnessing the Known Unknowns: Differential Privacy and the 2020 Census”. In: *Harvard Data Science Review* Special Issue 2 (June 2022). URL: <https://hdsr.mitpress.mit.edu/pub/fgyf5cne>.
- [57] Melissa Gymrek et al. “Identifying Personal Genomes by Surname Inference”. In: *Science* 339 (6117 Jan. 2013), pp. 321–329.
- [58] Michael B. Hawes. “Implementing Differential Privacy: Seven Lessons From the 2020 United States Census”. In: *Harvard Data Science Review* 2.2 (Apr. 2020). URL: <https://hdsr.mitpress.mit.edu/pub/dgg03vo6>.
- [59] TN Herzog, FJ Scheuren, and WE Winkler. *Data Quality and Record Linkage Techniques*. New York/London: Springer, 2007.
- [60] Vagelis Hristidis, ed. *Information Discovery on Electronic Health Records*. 1st. Chapman and Hall/CRC, 2009. ISBN: 1420090380. URL: <https://doi.org/10.1201/9781420090413>.
- [61] IHE IT Infrastructure Technical Committee. *IHE IT Infrastructure Handbook: De-Identification*. Integrating the Healthcare Enterprise, Mar. 2014. URL: https://ihe.net/uploadedFiles/Documents/ITI/IHE_ITI_Handbook_De-Identification_Rev1.0_2014-03-14.pdf.
- [62] Clay Johnson III. *OMB Memorandum M-07-16: Safeguarding Against and Responding to the Breach of Personally Identifiable Information*. May 2007. URL: <https://georgewbush-whitehouse.archives.gov/omb/memoranda/fy2007/m07-16.pdf>.

- [63] Information Commissioner’s Office. *Anonymisation: code of practice, managing data protection risk*. 2012. URL: <https://ico.org.uk/media/1061/anonymisation-code.pdf>.
- [64] *ISO 26324:2012, Information and documentation – Digital object identifier system*. Geneva, Switzerland, 2012. URL: <https://www.iso.org/standard/43506.html>.
- [65] *ISO/IEC 24760-1:2011, Information technology – Security techniques – A framework for identity management – Part 1: Terminology and concepts*. 2011.
- [66] *ISO/TS 25237:2008(E) Health Informatics — Pseudonymization*. Geneva, Switzerland, 2008.
- [67] Auguste Kerckhoffs. “II. Desiderata De La Cryptographie Militaire”. In: *Journal des sciences militaires* IX (Jan. 1883), pp. 5–38.
- [68] Shehab Khan. ““Human chimera”: Man fails paternity test because genes in his saliva are different to those in sperm”. In: *The Independent* (Oct. 2015).
- [69] Vladimir Kolesnikov et al. “Efficient Batched Oblivious PRF with Applications to Private Set Intersection”. In: *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*. CCS ’16. Vienna, Austria: Association for Computing Machinery, 2016, pp. 818–829. ISBN: 9781450341394. DOI: [10.1145/2976749.2978381](https://doi.org/10.1145/2976749.2978381). URL: <https://doi.org/10.1145/2976749.2978381>.
- [70] Leah Krehling. *De-Identification Guideline*. Tech. rep. WL-2020-01. Department of Electrical and Computer Engineering, Western University, 2020, p. 45.
- [71] Sandra Lechner and Winfried Pohlmeier. “To Blank or Not to Blank? A Comparison of the Effects of Disclosure Limitation Methods on Nonlinear Regression Estimates”. In: *Privacy in Statistical Databases, Lecture Notes in Computer Science* 3050 (2004), pp. 187–200.
- [72] Jaewoo Lee and Chris Clifton. “How Much Is Enough? Choosing ϵ for Differential Privacy”. In: *Information Security*. Ed. by Xuejia Lai, Jianying Zhou, and Hui Li. Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, pp. 325–340. ISBN: 978-3-642-24861-0.
- [73] Ninghui Li, Tiancheng Li, and Suresh Venkatasubramanian. “t-Closeness: Privacy Beyond k-Anonymity and l-Diversity”. In: *2007 IEEE 23rd International Conference on Data Engineering*. 2007, pp. 106–115. DOI: [10.1109/ICDE.2007.367856](https://doi.org/10.1109/ICDE.2007.367856).
- [74] Yehuda Lindell. “Secure Multiparty Computation”. In: *Commun. ACM* 64.1 (Dec. 2020), pp. 86–96. ISSN: 0001-0782. DOI: [10.1145/3387108](https://doi.org/10.1145/3387108). URL: <https://doi.org/10.1145/3387108>.
- [75] M. Altman M et al. “Towards a Modern Approach to Privacy-Aware Government Data Release”. In: *Berkeley Journal of Technology Law Internet* (2016). URL: <https://lawcat.berkeley.edu/record/1127405?ln=en>.

- [76] Ashwin Machanavajjhala et al. “l-diversity: Privacy beyond k-anonymity”. In: *Proc. 22nd Intl. Conf. Data Engg. (ICDE)*. 2006.
- [77] Sean Martin. *When De-identifying Patient Information, Follow the HITRUST Framework*. Sept. 2016. URL: <https://hitrustalliance.net/de-identifying-patient-information-follow-hitrust-framework/>.
- [78] Daniel A. Mayer et al. “Implementation and Performance Evaluation of Privacy-Preserving Fair Reconciliation Protocols on Ordered Sets”. In: *Proceedings of the First ACM Conference on Data and Application Security and Privacy*. CODASPY ’11. San Antonio, TX, USA: Association for Computing Machinery, 2011, pp. 109–120. ISBN: 9781450304665. DOI: [10.1145/1943513.1943529](https://doi.org/10.1145/1943513.1943529). URL: <https://doi.org/10.1145/1943513.1943529>.
- [79] E McCallister, T Grance, and K A Scarfone. *Guide to protecting the confidentiality of Personally Identifiable Information (PII)*. Gaithersburg, MD, 2010. DOI: [10.6028/NIST.SP.800-122](https://nvlpubs.nist.gov/nistpubs/Legacy/SP/nistspecialpublication800-122.pdf). URL: <https://nvlpubs.nist.gov/nistpubs/Legacy/SP/nistspecialpublication800-122.pdf>.
- [80] William K. Michener et al. “Participatory design of DataONE—Enabling cyber-infrastructure for the biological and environmental sciences”. In: *Ecological Informatics* 11 (2012). Data platforms in integrative biodiversity research, pp. 5–15. ISSN: 1574-9541. DOI: <https://doi.org/10.1016/j.ecoinf.2011.08.007>. URL: <https://www.sciencedirect.com/science/article/pii/S1574954111000768>.
- [81] Yves-Alexandre de Montjoye et al. “Unique in the Crowd: The Privacy Bounds of Human Mobility”. In: *Nature Scientific Reports* 3 (1376 2013).
- [82] Yves-Alexandre de Montjoye et al. “Unique in the Shopping Mall: On the Reidentifiability of Credit Card Metadata”. In: *Science* 347 (536 2015).
- [83] Arvind Narayanan and Ed Felten. *No silver bullet: De-identification still doesn’t work*. Working Paper. July 2014. URL: <http://randomwalker.info/publications/no-silver-bullet-de-identification.pdf>.
- [84] Arvind Narayanan and Vitaly Shmatikov. “Robust De-anonymization of Large Sparse Datasets”. In: *2008 IEEE Symposium on Security and Privacy (sp 2008)*. 2008, pp. 111–125. DOI: [10.1109/SP.2008.33](https://doi.org/10.1109/SP.2008.33).
- [85] *NIST Big Data Interoperability Framework: volume 1, definitions, version 2*. Gaithersburg, MD, 2018. DOI: [10.6028/NIST.SP.1500-1r1](https://nvlpubs.nist.gov/nistpubs/SpecialPublications/NIST.SP.1500-1r1.pdf). URL: <https://nvlpubs.nist.gov/nistpubs/SpecialPublications/NIST.SP.1500-1r1.pdf>.
- [86] *NIST Privacy Framework: A Tool for Improving Privacy Through Enterprise Risk Management, Version 1.0*. Gaithersburg, MD, 2022. DOI: [10.6028/NIST.CSWP.10](https://doi.org/10.6028/NIST.CSWP.10). URL: <https://nvlpubs.nist.gov/nistpubs/CSWP/NIST.CSWP.01162020.pdf>.

- [87] Christine M. O’Keefe and James O. Chipperfield. “A Summary of Attack Methods and Confidentiality Protection Measures for Fully Automated Remote Analysis Systems”. In: *International Statistical Review / Revue Internationale de Statistique* 81.3 (2013), pp. 426–455. ISSN: 03067734, 17515823. URL: <http://www.jstor.org/stable/43299645> (visited on 06/28/2022).
- [88] Barack Obama. *Executive Order 13642—Making Open and Machine Readable the New Default for Government Information*. May 2013. URL: <https://obamawhitehouse.archives.gov/the-press-office/2013/05/09/executive-order-making-open-and-machine-readable-new-default-government->.
- [89] Office of Civil Rights, US Department of Health and Human Services. *Guidance Regarding Methods for De-identification of Protected Health Information in Accordance with the Health Insurance Portability and Accountability Act (HIPAA) Privacy Rule*. Nov. 2012. URL: <http://www.hhs.gov/hipaa/for-professionals/privacy/special-topics/de-identification/>.
- [90] Office of Civil Rights, US Department of Health and Human Services. *Individuals’ Right under HIPAA to Access their Health Information 45 CFR § 164.524*. Last accessed June 17, 2022. 2022. URL: <https://www.hhs.gov/hipaa/for-professionals/privacy/guidance/access/index.html>.
- [91] Office of Management and Budget. *Circular A110 Revised 11/19/93, as further amended 9/30/99*. URL: https://obamawhitehouse.archives.gov/omb/circulars_a110/.
- [92] Office of Management and Budget. *Statistical Programs and Standards*. Last accessed July 15, 2022. 2022. URL: <https://www.whitehouse.gov/omb/information-regulatory-affairs/statistical-programs-standards/>.
- [93] Office of Safeguards, US Internal Revenue Service. *Publication 1075: Tax Information Security Guidelines For Federal, State and Local Agencies*. 2021. URL: <https://www.irs.gov/pub/irs-pdf/p1075.pdf>.
- [94] Paul Ohm. “Broken Promises of Privacy: Responding to the Surprising Failure of Anonymization”. In: *UCLA Law Review* 57 (July 2012), pp. 1701–1778.
- [95] *OHRP-Guidance on Research Involving Private Information or Biological Specimens*. Aug. 2008. URL: <http://www.hhs.gov/ohrp/policy/cdebiol.html>.
- [96] Joanne Pascale et al. *Issue Paper on Disclosure Review for Information Products with Qualitative Research Findings*. Mar. 2020. URL: <https://www.census.gov/library/working-papers/2020/adrm/rsm2020-01.html>.
- [97] Joanne Pascale et al. “Protecting the Identity of Participants in Qualitative Research”. In: *Journal of Survey Statistics and Methodology* 10.3 (Jan. 2022), pp. 549–567. ISSN: 2325-0984. DOI: 10.1093/jssam/smab048. eprint: <https://academic.oup.com/jssam/article-pdf/10/3/549/44275508/smab048.pdf>. URL: <https://doi.org/10.1093/jssam/smab048>.

- 2757 [98] Andrew Peterson. “Why the names of six people who complained of sexual assault
2758 were published online by Dallas police”. In: *The Washington Post* (Apr. 2016).
2759 URL: [https://www.washingtonpost.com/news/the-switch/wp/2016/04/29/why-the-
2760 names-of-six-people-who-complained-of-sexual-assault-were-published-online-
2761 by-dallas-police/](https://www.washingtonpost.com/news/the-switch/wp/2016/04/29/why-the-names-of-six-people-who-complained-of-sexual-assault-were-published-online-by-dallas-police/).
- 2762 [99] Thomas Piketty and Emmanuel Saez. “Income Inequality in the United States 1913-
2763 1998”. In: *Quarterly Journal of Economics* 118 (1 2003), pp. 1–41.
- 2764 [100] “Pillar Investigates: USCCB gen sec Burrill resigns after sexual misconduct alle-
2765 gations”. In: *The Pillar* (July 2021). URL: [https://www.pillarcatholic.com/p/pillar-
2766 investigates-usccb-gen-sec](https://www.pillarcatholic.com/p/pillar-investigates-usccb-gen-sec).
- 2767 [101] Sandro Pinto and Nuno Santos. “Demystifying Arm TrustZone: A Comprehensive
2768 Survey”. In: *ACM Comput. Surv.* 51.6 (Jan. 2019). ISSN: 0360-0300. DOI: [10.1145/
2769 3291047](https://doi.org/10.1145/3291047). URL: <https://doi.org/10.1145/3291047>.
- 2770 [102] *Private Lives and Public Policies: Confidentiality and Accessibility of Government
2771 Statistics*. Panel on Confidentiality and Data Access, National Research Council,
2772 p. 288. ISBN: 0-309-57611-3. URL: <http://www.nap.edu/catalog/2122/>.
- 2773 [103] *Public Law 93-579: The Privacy Act*. 88 Stat. 1896, 5 U.S.C. § 552a.
- 2774 [104] Balaji Raghunathan. *The Complete Book of Data Anonymization: From Planning
2775 to Implementation*. USA: Auerbach Publications, 2013. ISBN: 1439877300.
- 2776 [105] William H. Rehnquist. *Department of State v. Washington Post Co.*, 456 U.S. 595
2777 (1982). 1982. URL: <https://www.loc.gov/item/usrep456595/>.
- 2778 [106] *Report 08-536, Privacy: Alternatives Exist for Enhancing Protection of Personally
2779 Identifiable Information*. May 2008. URL: [http://www.gao.gov/new.items/d08536.
2780 pdf](http://www.gao.gov/new.items/d08536.pdf).
- 2781 [107] Diane Ridgeway et al. *Challenge Design and Lessons Learned from the 2018 Dif-
2782 ferential Privacy Challenges*. 2021. DOI: [10.6028/NIST.TN.2151](https://doi.org/10.6028/NIST.TN.2151). URL: [https:
2783 //nvlpubs.nist.gov/nistpubs/TechnicalNotes/NIST.TN.2151.pdf](https://nvlpubs.nist.gov/nistpubs/TechnicalNotes/NIST.TN.2151.pdf).
- 2784 [108] Pierangela Samarati and Latanya Sweeney. “Protecting privacy when disclosing
2785 information: k-anonymity and its enforcement through generalization and suppres-
2786 sion”. In: *Proceedings of the IEEE Symposium on Research in Security and Privacy*
2787 (May 1998).
- 2788 [109] Pierangela Samarti. “Protecting Respondents’ Identities in Microdata Release”.
2789 In: *IEEE Transactions on Knowledge and Data Engineering* 13 (6 Nov. 2001),
2790 pp. 1010–1027.
- 2791 [110] Josep Sanz and Josep Domingo-Ferrer. “A Comparative Study of Microaggrega-
2792 tion Methods”. In: *Questiio: Quaderns d’Estadística, Sistemes, Informàtica i In-
2793 vestigació Operativa* 22 (3 Aug. 2000), pp. 511–526.

- 2794 [111] M Scaiano et al. “unified framework for evaluating the risk of re-identification of
2795 text de-identification tools”. In: *Journal of Biomedical Informatics* 63 (Oct. 2016),
2796 pp. 174–183.
- 2797 [112] Matthias Schunter. “Intel Software Guard Extensions: Introduction and Open Re-
2798 search Challenges”. In: *Proceedings of the 2016 ACM Workshop on Software PRO-*
2799 *tection*. SPRO ’16. Vienna, Austria: Association for Computing Machinery, 2016,
2800 p. 1. ISBN: 9781450345767. DOI: [10.1145/2995306.2995307](https://doi.org/10.1145/2995306.2995307). URL: [https://doi.org/](https://doi.org/10.1145/2995306.2995307)
2801 [10.1145/2995306.2995307](https://doi.org/10.1145/2995306.2995307).
- 2802 [113] M Seastrom. “Licensing”. In: *Confidentiality, Disclosure and Data Access: Theory*
2803 *and Practical Application for Statistical Agencies*. Ed. by P. Doyle et al. Elsevier
2804 Science, 2001.
- 2805 [114] Jordi Soria-Comas and Josep Domingo-Ferrer. “Connecting privacy models: syn-
2806 ergies between k-anonymity, t-closeness, and differential privacy”. In: Working
2807 Paper (English Only). Ottawa, Canada, Oct. 2013. URL: [https://www.unece.](https://www.unece.org/fileadmin/DAM/stats/documents/ece/ces/ge.46/2013/Topic_2_soria-comas_domingo-ferrer.pdf)
2808 [org/fileadmin/DAM/stats/documents/ece/ces/ge.46/2013/Topic_2_soria-](https://www.unece.org/fileadmin/DAM/stats/documents/ece/ces/ge.46/2013/Topic_2_soria-comas_domingo-ferrer.pdf)
2809 [comas_domingo-ferrer.pdf](https://www.unece.org/fileadmin/DAM/stats/documents/ece/ces/ge.46/2013/Topic_2_soria-comas_domingo-ferrer.pdf).
- 2810 [115] Philip Steel and Jon Sperling. *The Impact of Multiple Geographies and Geographic*
2811 *Detail on Disclosure Risk: Interactions between Census Tract and ZIP Code Tab-*
2812 *ulation Geography*. 2001. URL: [https://www.census.gov/content/dam/Census/](https://www.census.gov/content/dam/Census/library/working-papers/2001/adrm/steel-sperling-2001.pdf)
2813 [library/working-papers/2001/adrm/steel-sperling-2001.pdf](https://www.census.gov/content/dam/Census/library/working-papers/2001/adrm/steel-sperling-2001.pdf).
- 2814 [116] Jackson M. Steinkamp et al. “Evaluation of Automated Public De-Identification
2815 Tools on a Corpus of Radiology Reports”. In: *Radiol Artif Intell* 2 (6 Oct. 2020).
2816 URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8082401/>.
- 2817 [117] John Paul Stevens. *U.S. Dept. of Justice v. Reporters Committee For Freedom of*
2818 *Press*, 489 U.S. 749 (1989). 1988. URL: <https://www.loc.gov/item/usrep489749/>.
- 2819 [118] John Paul Stevens. *United States Department of State v. Ray et al.*, 502 U.S. 164
2820 (1991). 1991. URL: <https://www.loc.gov/item/usrep502164/>.
- 2821 [119] Kevin Stine et al. *Volume I: guide for mapping types of information and infor-*
2822 *mation systems to security categories*. Gaithersburg, MD, 2008. DOI: [10.6028/](https://nvlpubs.nist.gov/nistpubs/Legacy/SP/nistspecialpublication800-60v1r1.pdf)
2823 [NIST.SP.800-60v1r1](https://nvlpubs.nist.gov/nistpubs/Legacy/SP/nistspecialpublication800-60v1r1.pdf). URL: [https://nvlpubs.nist.gov/nistpubs/Legacy/SP/](https://nvlpubs.nist.gov/nistpubs/Legacy/SP/nistspecialpublication800-60v1r1.pdf)
2824 [nistspecialpublication800-60v1r1.pdf](https://nvlpubs.nist.gov/nistpubs/Legacy/SP/nistspecialpublication800-60v1r1.pdf).
- 2825 [120] Tim Stobierski. In: *Business Insights* (Feb. 2021). URL: [https://online.hbs.edu/](https://online.hbs.edu/blog/post/data-life-cycle)
2826 [blog/post/data-life-cycle](https://online.hbs.edu/blog/post/data-life-cycle).
- 2827 [121] Teresa A. Sullivan. “Coming to Our Census: How Social Statistics Underpin Our
2828 Democracy (and Republic)”. In: *Harvard Data Science Review* 2.1 (Jan. 2020).
2829 URL: <https://hdsr.mitpress.mit.edu/pub/1g1cbvkx>.
- 2830 [122] Latanya Sweeney. “k-anonymity: a model for protecting privacy”. In: *International*
2831 *Journal on Uncertainty, Fuzziness and Knowledge-based Systems* 10 (5 2002),
2832 pp. 557–570.

- 2833 [123] Latanya Sweeney. “*k*-anonymity: a model for protecting privacy”. In: *Int. J. Un-*
2834 *certain. Fuzziness Knowl.-Based Syst.* 10 (5 Oct. 2002), pp. 557–570. URL: [http:](http://dx.doi.org/10.1142/S0218488502001648)
2835 [//dx.doi.org/10.1142/S0218488502001648](http://dx.doi.org/10.1142/S0218488502001648).
- 2836 [124] Latanya Sweeney. “Weaving Technology and Policy Together to Maintain Confi-
- 2837 dentiality”. In: *Journal of Law, Medicine and Ethics* 25 (1997), pp. 98–110.
- 2838 [125] “The Debate Over ‘Re-Identification’ Of Health Information: What Do We Risk?”
- 2839 In: *Health Affairs Blog* (Aug. 2012). DOI: [10.1377/hblog20120810.021952](https://doi.org/10.1377/hblog20120810.021952). URL:
- 2840 <https://www.healthaffairs.org/doi/10.1377/forefront.20120810.021952>.
- 2841 [126] *Title V of the E-Government Act of 2002: Confidential Information Protection and*
- 2842 *Statistical Efficiency Act (CIPSEA) PL 107–347, 116 Stat. 2899, 44 USC § 101*
- 2843 *Section 502(8)*.
- 2844 [127] TransCelerate Biopharma, Inc. *Data De-identification and Anonymization of Indi-*
- 2845 *vidual Patient Data in Clinical Studies—A Model Approach*. 2013.
- 2846 [128] Michael Carl Tschantz and Jeannette M. Wing. *Formal Methods for Privacy*. Tech.
- 2847 rep. CMU-CS-09-154. Pittsburgh, PA: Carnegie Mellon University, Aug. 2009. URL:
- 2848 <http://reports-archive.adm.cs.cmu.edu/anon/2009/CMU-CS-09-154.pdf>.
- 2849 [129] A. M. Turing. “On Computable Numbers, with an Application to the Entschei-
- 2850 dungsproblem”. In: *Proceedings of the London Mathematical Society, Series 2* (42
- 2851 1936–37), pp. 230–265.
- 2852 [130] US Census Bureau. *Census Confidentiality and Privacy: 1790-2002*. 2003. URL:
- 2853 <https://www.census.gov/prod/2003pubs/conmono2.pdf>.
- 2854 [131] US Census Bureau. *The “72-Year Rule”*. Jan. 2022. URL: [https://www.census.gov/](https://www.census.gov/history/www/genealogy/decennial_census_records/the_72_year_rule_1.html)
- 2855 [history/www/genealogy/decennial_census_records/the_72_year_rule_1.html](https://www.census.gov/history/www/genealogy/decennial_census_records/the_72_year_rule_1.html).
- 2856 [132] US Congress. *Public Law 104-191: Health Insurance Portability and Accountabil-*
- 2857 *ity Act of 1996 (HIPAA)*. Aug. 1996. URL: [https://www.congress.gov/bill/104th-](https://www.congress.gov/bill/104th-congress/house-bill/3103)
- 2858 [congress/house-bill/3103](https://www.congress.gov/bill/104th-congress/house-bill/3103).
- 2859 [133] US Congress. *Public Law 114-185: FOIA Improvement Act of 2016*. 2016. URL:
- 2860 <https://www.congress.gov/114/plaws/publ185/PLAW-114publ185.pdf>.
- 2861 [134] US Congress. *The Freedom Of Information Act, 5 U.S.C. § 552*. 2022. URL: [https:](https://www.justice.gov/oip/freedom-information-act-5-usc-552)
- 2862 [//www.justice.gov/oip/freedom-information-act-5-usc-552](https://www.justice.gov/oip/freedom-information-act-5-usc-552).
- 2863 [135] US Department of Health and Human Services. *45 CFR Part 46: Federal Policy*
- 2864 *for the Protection of Human Subjects*. Jan. 2017. URL: [https://www.govinfo.gov/](https://www.govinfo.gov/content/pkg/FR-2017-01-19/pdf/2017-01058.pdf)
- 2865 [content/pkg/FR-2017-01-19/pdf/2017-01058.pdf](https://www.govinfo.gov/content/pkg/FR-2017-01-19/pdf/2017-01058.pdf).
- 2866 [136] US Department of Health and Human Services. *Guidance Regarding Methods for*
- 2867 *De-identification of Protected Health Information in Accordance with the Health*
- 2868 *Insurance Portability and Accountability Act (HIPAA) Privacy Rule*. 2012. URL:
- 2869 [https://www.hhs.gov/hipaa/for-professionals/privacy/special-topics/de-](https://www.hhs.gov/hipaa/for-professionals/privacy/special-topics/de-identification/index.html)
- 2870 [identification/index.html](https://www.hhs.gov/hipaa/for-professionals/privacy/special-topics/de-identification/index.html).

- 2871 [137] *Joel Havermann, plaintiff—Appellant v. Carolyn W. Colvin, Acting Commissioner*
2872 *of the Social Security Administration, Defendant—Appellee, No. 12-2453, US Court*
2873 *of Appeals for the Fourth Circuit, 537 Fed. Appx. 142; 2013 US App. Aug 1,*
2874 *2013. Joel Havemann v. Carolyn W. Colvin, Civil No. JFM-12-1325. 2015 US Dist.*
2875 *LEXIS 27560. Mar. 2015.*
- 2876 [138] “Utility”. In: *Glossary of Statistical Terms* (Aug. 2002). Last accessed June 23,
2877 2022. URL: <https://stats.oecd.org/glossary/detail.asp?ID=4884>.
- 2878 [139] Russell T. Vought. *Phase 1 Implementation of the Foundations for Evidence-Based*
2879 *Policymaking Act of 2018: Learning Agendas, Personnel, and Planning Guidance.*
2880 URL: <https://www.whitehouse.gov/wp-content/uploads/2019/07/M-19-23.pdf>.
- 2881 [140] Charlie Warzel and Stuart A. Thompson. “How Your Phone Betrays Democracy”.
2882 In: *The New York Times* (Dec. 2019). URL: [https://www.nytimes.com/interactive/](https://www.nytimes.com/interactive/2019/12/21/opinion/location-data-democracy-protests.html)
2883 [2019/12/21/opinion/location-data-democracy-protests.html](https://www.nytimes.com/interactive/2019/12/21/opinion/location-data-democracy-protests.html).
- 2884 [141] Cathy Wasserman and Eric Ossiander. *Department of Health Agency Standards for*
2885 *Reporting Data with Small Numbers.* May 2018. URL: [https://doh.wa.gov/sites/](https://doh.wa.gov/sites/default/files/legacy/Documents/1500/SmallNumbers.pdf)
2886 [default/files/legacy/Documents/1500/SmallNumbers.pdf](https://doh.wa.gov/sites/default/files/legacy/Documents/1500/SmallNumbers.pdf).
- 2887 [142] Leon Willenborg and Ton de Waal. “Chapter 3 Data Analytic Impact of SDC
2888 Techniques on Microdata”. In: *Elements of Statistical Disclosure Control* (2012),
2889 pp. 72–92.
- 2890 [143] Li Xiong et al. “Privacy-Preserving Information Discovery on EHRs”. In: *Informa-*
2891 *tion Discovery on Electronic Health Records.* Ed. by Vagelis Hristidis. 1st. Chap-
2892 man and Hall/CRC, 2009. ISBN: 1420090380. URL: [https://doi.org/10.1201/](https://doi.org/10.1201/9781420090413)
2893 [9781420090413](https://doi.org/10.1201/9781420090413).

Appendix A. Standards

- ASTM E1869-04(2014) Standard Guide for Confidentiality, Privacy, Access, and Data Security Principles for Health Information Including Electronic Health Records.
- DICOM PS3.15 2016d – Security and System Management Profiles Chapter E Attribute Confidentiality Profiles, DICOM Standards Committee, NEMA 2016. http://dicom.nema.org/medical/dicom/current/output/html/part15.html#chapter_E
- HITRUST De-Identification Working Group (2015, March). De-Identification Framework: A Consistent, Managed Methodology for the De-Identification of Personal Data and the Sharing of Compliance and Risk Information. Frisco, TX: HITRUST. Retrieved from <https://hitrustalliance.net/de-identification-license-agreement/>
- ISO 8000-2:2012(E) Data quality – Part 2: Vocabulary, 2012. ISO, Geneva, Switzerland. 2012.
- ISO/IEC 27000:2014 Information technology -- Security techniques -- Information security management systems -- Overview and vocabulary. ISO, Geneva, Switzerland. 2012.
- ISO/IEC 24760-1:2011 Information technology -- Security techniques -- A framework for identity management -- Part 1: Terminology and concepts. ISO, Geneva, Switzerland. 2011.
- ISO/TS 25237:2008(E) Health Informatics – Pseudonymization. ISO, Geneva, Switzerland. 2008.
- ISO/IEC 20889 WORKING DRAFT 2016-05-30, Information technology – Security techniques – Privacy enhancing data de-identification techniques. ISO, Geneva, Switzerland. 2016.
- IHE IT Infrastructure Handbook, De-Identification, Integrating the Healthcare Enterprise, June 6, 2014. http://www.ihe.net/User_Handbooks/

Appendix A.1. NIST Publications

- *NIST Privacy Framework: A Tool for Improving Privacy Through Enterprise Risk Management, Version 1.0*. Gaithersburg, MD, 2022. DOI: [10.6028/NIST.CSWP.10](https://doi.org/10.6028/NIST.CSWP.10). URL: <https://nvlpubs.nist.gov/nistpubs/CSWP/NIST.CSWP.01162020.pdf>
- Kevin Stine et al. *Volume I: guide for mapping types of information and information systems to security categories*. Gaithersburg, MD, 2008. DOI: [10.6028/NIST.SP.800-60v1r1](https://doi.org/10.6028/NIST.SP.800-60v1r1). URL: <https://nvlpubs.nist.gov/nistpubs/Legacy/SP/nistspecialpublication800-60v1r1.pdf>
- Simson L. Garfinkel. *De-identification of personal information*. 2015. DOI: [10.6028/NIST.IR.8053](https://doi.org/10.6028/NIST.IR.8053). URL: <https://nvlpubs.nist.gov/nistpubs/ir/2015/NIST.IR.8053.pdf>

- 2929 • Simson L. Garfinkel. *Government Data De-Identification Stakeholder's Meeting*
2930 *June 29, 2016 Meeting Report*. 2016. DOI: [10.6028/NIST.IR.8150](https://doi.org/10.6028/NIST.IR.8150). URL: [https:](https://nvlpubs.nist.gov/nistpubs/ir/2016/NIST.IR.8150.pdf)
2931 [//nvlpubs.nist.gov/nistpubs/ir/2016/NIST.IR.8150.pdf](https://nvlpubs.nist.gov/nistpubs/ir/2016/NIST.IR.8150.pdf)

2932 **Appendix A.2. Other U.S. Government Publications**

- 2933 • *Census Confidentiality and Privacy: 1790-2002*, US Census Bureau, 2003. [https:](https://www.census.gov/prod/2003pubs/conmono2.pdf)
2934 [//www.census.gov/prod/2003pubs/conmono2.pdf](https://www.census.gov/prod/2003pubs/conmono2.pdf)
- 2935 • *Data De-identification: An Overview of Basic Terms*, Privacy Technical Assistance
2936 Center, US Department of Education. May 2013. [http://ptac.ed.gov/sites/default/](http://ptac.ed.gov/sites/default/files/data_deidentification_terms.pdf)
2937 [files/data_deidentification_terms.pdf](http://ptac.ed.gov/sites/default/files/data_deidentification_terms.pdf)
- 2938 • *Data Disclosure Decision*, Department of Education (ED) Disclosure Review Board
2939 (DRB), A Product of the Federal CIO Council Innovation Committee. Version 1.0,
2940 2015.
- 2941 • *Disclosure Avoidance Techniques at the US Census Bureau: Current Practices and*
2942 *Research*, Research Report Series (Disclosure Avoidance #2014-02), Amy Lauger,
2943 Billy Wisniewski, and Laura McKenna, Center for Disclosure Avoidance Research,
2944 US Census. Bureau, September 26, 2014. [https://www.census.gov/srd/CDAR/cdar2014-02_](https://www.census.gov/srd/CDAR/cdar2014-02_Discl_Avoid_Techniques.pdf)
2945 [Discl_Avoid_Techniques.pdf](https://www.census.gov/srd/CDAR/cdar2014-02_Discl_Avoid_Techniques.pdf)
- 2946 • *Frequently Asked Questions – Disclosure Avoidance, Privacy Technical Assistance*
2947 *Center*, U.S. Department of Education. October 2012 (revised July 2015). [http:](http://ptac.ed.gov/sites/default/files/FAQ_Disclosure_Avoidance.pdf)
2948 [//ptac.ed.gov/sites/default/files/FAQ_Disclosure_Avoidance.pdf](http://ptac.ed.gov/sites/default/files/FAQ_Disclosure_Avoidance.pdf)
- 2949 • *Guidance Regarding Methods for De-identification of Protected Health Information*
2950 *in Accordance with the Health Insurance Portability and Accountability Act (HIPAA)*
2951 *Privacy Rule*, U.S. Department of Health & Human Services, Office for Civil Rights,
2952 November 26, 2012. [http://www.hhs.gov/ocr/privacy/hipaa/understanding/coveredentities/](http://www.hhs.gov/ocr/privacy/hipaa/understanding/coveredentities/De-identification/hhs_deid_guidance.pdf)
2953 [De-identification/hhs_deid_guidance.pdf](http://www.hhs.gov/ocr/privacy/hipaa/understanding/coveredentities/De-identification/hhs_deid_guidance.pdf)
- 2954 • <http://www.hhs.gov/ohrp/policy/cdebiol.html>
- 2955 • http://ptac.ed.gov/sites/default/files/data_deidentification_terms.pdf
- 2956 • The Data Disclosure Decision, Department of Education (ED) Disclosure Review
2957 Board (DRB), A Product of the Federal CIO Council Innovation Committee.
- 2958 • https://www.cdc.gov/nchs/data/nchs_microdata_release_policy_4-02a.pdf
- 2959 • http://www.cdc.gov/nchs/nvss/dvs_data_release.htm
- 2960 • *Linking Data for Health Services Research: A Framework and Instructional Guide.*,
2961 *Dusetzina SB, Tyree S, Meyer AM, Meyer A, Green L, Carpenter WR. (Prepared*
2962 *by the University of North Carolina at Chapel Hill under Contract No. 290-2010-*

- 2963 000141.) *AHRQ Publication No. 14-EHC033-EF*. Rockville, MD: Agency for Health-
2964 care Research and Quality; September 2014.
- 2965 • *National Center for Health Statistics Data Release and Access Policy for Micro-data*
2966 *and Compressed Vital Statistics File*, Centers for Disease Control, April 26, 2011.
2967 http://www.cdc.gov/nchs/nvss/dvs_data_release.htm
 - 2968 • *National Center for Health Statistics Policy on Micro-Data Dissemination*, Centers
2969 for Disease Control, July 2002. [https://www.cdc.gov/nchs/data/nchs_microdata_release_policy_4-](https://www.cdc.gov/nchs/data/nchs_microdata_release_policy_4-02a.pdf)
2970 [02a.pdf](https://www.cdc.gov/nchs/data/nchs_microdata_release_policy_4-02a.pdf)
 - 2971 • *OHRP-Guidance on Research Involving Private Information or Biological Specimens* (2008), Department of Health & Human Services, Office of Human Research
2972 Protections (OHRP), August 16, 2008. <http://www.hhs.gov/ohrp/policy/cdebiol.html>
 - 2973 • OMB Circular A-130, *Managing Information as a Strategic Resource*, July 2016.
 - 2974 • *Privacy and Confidentiality Research and the U.S. Census Bureau, Recommendations Based on a Review of the Literature*, Thomas S. Mayer, Statistical Research Di-
2975 vision, US Bureau of the Census. February 7, 2002. [https://www.census.gov/srd/papers/pdf/rsm2002-](https://www.census.gov/srd/papers/pdf/rsm2002-01.pdf)
2976 [01.pdf](https://www.census.gov/srd/papers/pdf/rsm2002-01.pdf)
 - 2977 • *Statistical Policy Working Paper 22 (Second version, 2005)*, Report on Statistical
2978 Disclosure Limitation Methodology, Federal Committee on Statistical Methodology,
2979 December 2005.

2982 Selected Publications by Other Governments

- 2983 • *Privacy business resource 4: De-identification of data and information*, Office of the
2984 Australian Information Commissioner, Australian Government, April 2014. [http://www.oaic.gov.au/ima-](http://www.oaic.gov.au/ima-resources/privacy-business-resources/privacy_business_resource_4.pdf)
2985 [resources/privacy-business-resources/privacy_business_resource_4.pdf](http://www.oaic.gov.au/ima-resources/privacy-business-resources/privacy_business_resource_4.pdf)
- 2986 • *Opinion 05/2014 on Anonymisation Techniques*, Article 29 Data Protection Working
2987 Party, 0829/14/EN WP216, Adopted on 10 April 2014.
- 2988 • *Anonymisation: Managing data protection risk, Code of Practice 2012*, Information
2989 Commissioner's Office. [https://ico.org.uk/media/for-organisations/documents/1061/anonymisation-](https://ico.org.uk/media/for-organisations/documents/1061/anonymisation-code.pdf)
2990 [code.pdf](https://ico.org.uk/media/for-organisations/documents/1061/anonymisation-code.pdf). 108 pages.
- 2991 • *The Anonymisation Decision-Making Framework*, Mark Elliot, Elaine Mackey, Kieron
2992 O'Hara and Caroline Tudor, UKAN, University of Manchester, July 2016. [http:](http://ukanon.net/ukan-resources/ukan-decision-making-framework/)
2993 [//ukanon.net/ukan-resources/ukan-decision-making-framework/](http://ukanon.net/ukan-resources/ukan-decision-making-framework/)

2994 Reports and Books

- 2995 • *Private Lives and Public Policies: Confidentiality and Accessibility of Government*
2996 *Statistics (1993)*, George T. Duncan, Thomas B. Jabine, and Virginia A. de Wolf,

- 2997 Editors; Panel on Confidentiality and Data Access; *Commission on Behavioral and*
2998 *Social Sciences and Education*; *Division of Behavioral and Social Sciences and Ed-*
2999 *ucation*; National Research Council, 1993. <http://dx.doi.org/10.17226/2122>
- 3000 • *Sharing Clinical Trial Data: Maximizing Benefits, Minimizing Risk*, Committee on
3001 Strategies for Responsible Sharing of Clinical Trial Data, Board on Health Sciences
3002 Policy, Institute of Medicine of the National Academies, The National Academies
3003 Press, Washington, DC. 2015.
 - 3004 • P. Doyle and J. Lane, *Confidentiality, Disclosure and Data Access: Theory and Prac-*
3005 *tical Applications for Statistical Agencies*, North-Holland Publishing, Dec 31, 2001.
 - 3006 • George T. Duncan, Mark Elliot, Juan-José Salazar-Gonzalez, *Statistical Confiden-*
3007 *tiality: Principles and Practice*, Springer, 2011.
 - 3008 • Cynthia Dwork and Aaron Roth, *The Algorithmic Foundations of Differential Pri-*
3009 *vac*y (Foundations and Trends in Theoretical Computer Science). Now Publishers,
3010 August 11, 2014. <http://www.cis.upenn.edu/~aaroht/privacybook.html>
 - 3011 • Khaled El Emam, *Guide to the De-Identification of Personal Health Information*,
3012 CRC Press, 2013.
 - 3013 • Khaled El Emam and Luk Arbuckle, *Anonymizing Health Data*, O'Reilly, Cam-
3014 bridge, MA. 2013.
 - 3015 • K El Emam and B Malin, "Appendix B: Concepts and Methods for De-Identifying
3016 Clinical Trial Data," in *Sharing Clinical Trial Data: Maximizing Benefits, Minimiz-*
3017 *ing Risk*, Institute of Medicine of the National Academies, The National Academies
3018 Press, Washington, DC. 2015.
 - 3019 • Anco Hundepool, Josep Domingo-Ferrer, Luisa Franconi, Sarah Giessing, Eric Schulte
3020 Nordholt, Keith Spicer, Peter-Paul de Wolf, *Statistical Disclosure Control*, Wiley,
3021 September 2012.

3022 How-To Articles

- 3023 • Leah Krehling, *De-Identification Guideline*, WHISPERLAB, Technical Report WL-
3024 2020-01, Department of Electrical and Computer Engineering, Western University,
3025 2020.
- 3026 • Olivia Angiuli, Joe Blitstein, and Jim Waldo, *How to De-Identify Your Data*, Com-
3027 munications of the ACM, December 2015.
- 3028 • Jörg Drechsler, Stefan Bender, Susanne Rässler, *Comparing fully and partially syn-*
3029 *thetic datasets for statistical disclosure control in the German IAB Establishment*
3030 *Panel*. 2007, United Nations, Economic Commission for Europe. Working paper,
3031 11, New York, 8 p. <http://fdz.iab.de/342/section.aspx/Publikation/k080530j05>

- 3032 • Ebaa Fayyumi and B. John Oommen, A survey on statistical disclosure control and
3033 micro-aggregation techniques for secure statistical databases. 2010, *Software Prac-*
3034 *tice and Experience*. 40, 12 (November 2010), 1161-1188. DOI=10.1002/spe.v40:12
3035 <http://dx.doi.org/10.1002/spe.v40:12><http://dx.doi.org/10.1002/spe.v40:12>
- 3036 • Jingchen Hu, Jerome P. Reiter, and Quanli Wang, Disclosure Risk Evaluation for
3037 Fully Synthetic Categorical Data, *Privacy in Statistical Databases*, pp. 185-199,
3038 2014. https://link.springer.com/chapter/10.1007/978-3-319-11257-2_15
- 3039 • Matthias Templ, Bernhard Meindl, Alexander Kowarik and Shuang Chen, Introduc-
3040 tion to Statistical Disclosure Control (SDC), IHSN Working Paper No. 007, Inter-
3041 national Household Survey Network, August 2014. [http://www.ihsn.org/home/sites/](http://www.ihsn.org/home/sites/default/files/resources/ihsn-working-paper-007-Oct27.pdf)
3042 [default/files/resources/ihsn-working-paper-007-Oct27.pdf](http://www.ihsn.org/home/sites/default/files/resources/ihsn-working-paper-007-Oct27.pdf)
- 3043 • Natalie Shlomo, Statistical Disclosure Control Methods for Census Frequency Ta-
3044 bles, *International Statistical Review* (2007), 75, 2, 199-217. [https://www.jstor.org/](https://www.jstor.org/stable/41508461)
3045 [stable/41508461](https://www.jstor.org/stable/41508461)

3046 **Appendix B. List of Symbols, Abbreviations, and Acronyms**

3047 Selected acronyms and abbreviations used in this paper are defined below.

3048 **ACM** Association for Computing Machinery

3049 **AHRQ** Agency for Healthcare Research and Quality

3050 **AMD** Advanced Micro Devices

3051 **ARM** Advanced RISC Machines (formerly Acron RISC Machine)

3052 **ARMP** average record matching probability

3053 **ASTM** ASTM (formerly the American Society for Testing and Materials)

3054 **CED-DA** Center for Enterprise Dissemination-Disclosure Avoidance

3055 **CFR** Code of Federal Regulations

3056 **CIO** chief information officer

3057 **CIPSEA** The Confidential Information Protection and Statistical Efficiency Act of 2002

3058 **CNSS** Committee on National Security Systems

3059 **CNSSI** Committee on National Security Systems instruction

3060 **CPU** central processing unit

3061 **CRC** (formerly the Chemical Rubber Company)

3062 **DC** District of Columbia

3063 **DCMA** Defense Contract Management Agency

3064 **DICOM** Digital Imaging and Communications in Medicine

3065 **DNA** deoxyribonucleic acid

3066 **DOI** digital object identifier

3067 **DRB** disclosure review board

3068 **DUA** data use agreement

3069 **EDDRB** Department of Education disclosure review board

3070 **FCSM** Federal Committee on Statistical Methodology

3071 **FHE** Fully-homomorphic encryption

3072 **FISMA** Federal Information Security Modernization Act

3073 **FOIA** Freedom of Information Act

3074	HHS Health and Human Services
3075	HIPAA Health Insurance Portability and Accountability Act
3076	HITRUST (formerly the Health Industry Trust Alliance)
3077	IAB Institut für Arbeitsmarkt-und Berufsforschung (Germany's Institute for Employment
3078	and Research)
3079	ICSP Interagency Council on Statistical Policy
3080	ID Identification number
3081	IEC International Electrotechnical Commission
3082	IHE Integrating the Healthcare Enterprise
3083	IHSN International Household Survey Network
3084	IP internet protocol
3085	IR inter-agency report
3086	IRB institutional review board
3087	IRS Internal Revenue Service
3088	ISO (formerly International Organization for Standardization)
3089	ISO/TS ISO Technical Standard
3090	IT information technology
3091	ITL Information Technology Laboratory
3092	KIRP Known inclusion re-identification probability
3093	MA Massachusetts
3094	MCC Millennium Challenge Corporation
3095	MD Maryland
3096	MIT Massachusetts Institute of Technology
3097	MPC multi-party computation
3098	NEMA National Electrical Manufacturers Association
3099	NIST National Institute of Standards and Technology
3100	NISTIR National Institute of Standards and Technology interagency report
3101	OECD Organisation for Economic Co-operation and Development
3102	OHRP Office for Human Research Protections

3103	OMB Office of Management and Budget
3104	OPRE Office of Planning, Research and Evaluation
3105	PDF portable document file
3106	PEC privacy enhancing cryptography
3107	PHI protected health information
3108	PII personally identifiable information
3109	PL public law
3110	PUF public use file
3111	RMP record matching probability
3112	SDC statistical disclosure control
3113	SDL statistical disclosure limitation
3114	SHA secure hash algorithm
3115	SLA service-level agreement
3116	SP special publication
3117	TEE trusted execution environments
3118	TX Texas
3119	UIRP Unknown inclusion re-identification probability
3120	UK United kingdom
3121	UKAN United Kingdom Advocacy Network
3122	US United States
3123	USC United States Code
3124	WHISPERLAB Western Information Security and Privacy Research Laboratory
3125	WP working paper

3126 **Appendix C. Glossary**

3127 Selected terms used in the publication are defined below. Where noted, the definition is
3128 sourced from another publication.

3129 **anonymization** A process that removes the association between the identifying dataset
3130 and the data subject. (ISO 25237-2008)

3131 **attribute** An inherent characteristic. (ISO 9241-302:2008)

3132 **attribute disclosure** Re-identification event in which an entity learns confidential infor-
3133 mation about a data principal, without necessarily identifying the data principal.
3134 (ISO/IEC 20889 WORKING DRAFT 2 2016-05-27)

3135 **anonymity** Condition in identification whereby an entity can be recognized as distinct,
3136 without sufficient identity information to establish a link to a known identity. (ISO/IEC
3137 24760-1:2011)

3138 **anticipated re-identification rate** When an organization contemplates performing re-identification,
3139 the re-identification rate that the resulting de-identified data are likely to have.

3140 **attacker** A person who seeks to exploit potential vulnerabilities of a system.

3141 **attribute** Characteristic or property of an entity that can be used to describe its state, ap-
3142 pearance, or other aspect. (ISO/IEC 24760-1:2011)[65]

3143 **brute force attack** In cryptography, an attack that involves trying all possible combina-
3144 tions to find a match.

3145 **characteristic** Distinguishing feature. (ISO 8000-2:2012(E))

3146 **coded** 1. Identifying information (such as name or social security number) that would
3147 enable the investigator to readily ascertain the identity of the individual to whom
3148 the private information or specimens pertain has been replaced with a number, let-
3149 ter, symbol, or combination thereof (i.e., the code); 2. A key to decipher the code
3150 exists, enabling linkage of the identifying information to the private information or
3151 specimens. [95]

3152 **control** Measure that is modifying risk. Note: controls include any process, policy, device,
3153 practice, or other actions which modify risk. (ISO/IEC 27000:2014)

3154 **covered entity** Under HIPAA, a health plan, a health care clearinghouse, or a health care
3155 provider that conducts certain health care transactions electronically (e.g., billing).
3156 (HIPAA Privacy Rule)

3157 **data** Re-interpretable representation of information in a formalized manner suitable for
3158 communication, interpretation, or processing. (ISO 8000-2:2012(E))

- 3159 **data accuracy** Closeness of agreement between a property value and the true value. (ISO
3160 8000-2:2012(E))
- 3161 **data dictionary** collection of data dictionary entries that allows lookup by entity identifier.
3162 (ISO 8000-2:2012(E))
- 3163 **data dictionary entry** Description of an entity type containing, at a minimum, an unam-
3164 biguous identifier, a term, and a definition. (ISO 8000-2:2012(E))
- 3165 **data intruder** A data user who attempts to disclose information about a population through
3166 identification or attribution. (OECD Glossary of Statistical Terms)
- 3167 **data life cycle** The set of processes in an application that transform raw data into action-
3168 able knowledge. (NIST SP 1500-1)
- 3169 **data subjects** Persons to whom data refer. (ISO/TS 25237:2008)
- 3170 **data use agreement** Executed agreement between a data provider and a data recipient that
3171 specifies the terms under which the data can be used.
- 3172 **data universe** All possible data within a specified domain.
- 3173 **dataset** A collection of data.
- 3174 **dataset with identifiers** A dataset that contains information that directly identifies indi-
3175 viduals.
- 3176 **dataset without identifiers** A dataset that does not contain direct identifiers.
- 3177 **de-identification** A process that is applied to a dataset with the goal of preventing or lim-
3178 iting informational risks to individuals, protected groups, and establishments, while
3179 still allowing for the production of aggregate statistics.²⁸
- 3180 **de-identification model** An approach to the application of data de-identification tech-
3181 niques that enables the calculation of re-identification risk. (ISO/IEC 20889 WORK-
3182 ING DRAFT 2 2016-05-27)
- 3183 **de-identification process** A general term for any process of removing the association be-
3184 tween a set of identifying data and the data principal. (ISO/TS 25237:2008)
- 3185 **de-identified information** Records that have had enough PII removed or obscured such
3186 that the remaining information does not identify an individual, and there is no rea-
3187 sonable basis to believe that the information can be used to identify an individual.
3188 (SP800-122)
- 3189 **direct identifying data** Data that directly identify a single individual. (ISO/TS 25237:2008)

²⁸ISO/TS 25237:2008 defines de-identification as the “general term for any process of removing the association between a set of identifying data and the data subject” [66, p.3]. This document intentionally adopts a broader definition for de-identification that allows for noise-introducing techniques, such as differential privacy and the creation of synthetic datasets that are based on privacy-preserving models.

- 3190 **disclosure** Divulging of, or provision of access to, data. (ISO/TS 25237:2008)
- 3191 **disclosure limitation** Statistical methods used to hinder anyone from identifying an indi-
3192 vidual respondent or establishment by analyzing published data, especially by ma-
3193 nipulating mathematical and arithmetical relationships among the data. [p.21][130]
- 3194 **effectiveness** The extent to which planned activities are realized and planned results achieved.
3195 (ISO/IEC 27000:2014)
- 3196 **entity** An item inside or outside an information and communication technology system,
3197 such as a person, an organization, a device, a subsystem, or a group of such items
3198 that has recognizably distinct existence. (ISO/IEC 24760-1:2011)
- 3199 **expert determination** Within the context of de-identification, refers to the Expert Deter-
3200 mination method for de-identifying protected health information in accordance with
3201 the HIPAA Privacy Rule de-identification standard.
- 3202 **Federal Committee on Statistical Methodology (FCSM)** An interagency committee ded-
3203 icated to improving the quality of Federal statistics. The FCSM was created by the
3204 Office of Management and Budget (OMB) to inform and advise OMB and the Inter-
3205 agency Council on Statistical Policy (ICSP) on methodological and statistical issues
3206 that affect the quality of Federal data. (fscm.sites.usa.gov)
- 3207 **genomic information** Information based on an individual's genome, such as a sequence
3208 of DNA or the results of genetic testing.
- 3209 **harm** Any adverse effects that would be experienced by an individual (i.e., that may be
3210 socially, physically, or financially damaging) or an organization if the confidentiality
3211 of PII were breached. (SP 800-122)
- 3212 **Health Insurance Portability and Accountability Act of 1996 (HIPAA)** A federal statute
3213 that called on the federal Department of Health and Human Services to establish reg-
3214 ulatory standards to protect the privacy and security of individually identifiable health
3215 information. See <https://www.hhs.gov/hipaa/for-professionals/index.html>.
- 3216 **HIPAA** See *Health Insurance Portability and Accountability Act of 1996*.
- 3217 **HIPAA Privacy Rule** Establishes national standards to protect individuals' medical records
3218 and other personal health information and applies to health plans, health care clear-
3219 inghouses, and those health care providers that conduct certain health care transac-
3220 tions electronically. (HIPAA Privacy Rule, 45 CFR 160, 162, 164). See [https://www.hhs.gov/hipaa/for-](https://www.hhs.gov/hipaa/for-professionals/privacy/index.html)
3221 [professionals/privacy/index.html](https://www.hhs.gov/hipaa/for-professionals/privacy/index.html).
- 3222 **identification** The process of using claimed or observed attributes of an entity to single
3223 out the entity among other entities in a set of identities. (ISO/TS 25237:2008)
- 3224 **identifying information** Information that can be used to distinguish or trace an individ-
3225 ual's identity (e.g., their name, social security number, biometric records, etc.) alone

3226 or when combined with other personal or identifying information that is linked or
3227 linkable to a specific individual (e.g., date and place of birth, mother’s maiden name,
3228 etc.). (OMB M-07-16)

3229 **identifier** Information used to claim an identity, before a potential corroboration by a cor-
3230 responding authenticator. (ISO/TS 25237:2008)

3231 **imputation** A procedure for entering a value for a specific data item where the response is
3232 missing or unusable. (OECD Glossary of Statistical Terms)

3233 **inference** Refers to the ability to deduce the identity of a person associated with a set of
3234 data through “clues” contained in that information. This analysis permits determi-
3235 nation of the individual’s identity based on a combination of facts associated with
3236 that person even though specific identifiers have been removed, like name and social
3237 security number. (ASTM E1869-04)[15]

3238 **information** Knowledge concerning objects, such as facts, events, things, processes, or
3239 ideas, including concepts, that within a certain context has a particular meaning.
3240 (ISO 8000-2:2012(E))

3241 **k-anonymity** A technique “to release person-specific data such that the ability to link to
3242 other information using the quasi-identifier is limited” [122]. *k*-anonymity achieves
3243 this through suppression of identifiers and output perturbation.

3244 **l-diversity** A refinement to the *k*-anonymity approach that assures that groups of records
3245 specified by the same identifiers have sufficient diversity to prevent inferential dis-
3246 closure. [76]

3247 **masking** The process of systematically removing a field or replacing it with a value in a
3248 way that does not preserve the analytic utility of the value, such as replacing a phone
3249 number with asterisks or a randomly generated pseudonym. [45]

3250 **motivated intruder test** The ‘motivated intruder’ is taken to be a person who starts with-
3251 out any prior knowledge but who wishes to identify the individual from whose per-
3252 sonal data the anonymised data has been derived. This test is meant to assess whether
3253 the motivated intruder would be successful. [63]

3254 **noise** A convenient term for a series of random disturbances borrowed through communi-
3255 cation engineering, from the theory of sound. In communication theory, noise results
3256 in the possibility of a signal sent, *x*, being different from the signal received, *y*, and
3257 the latter has a probability distribution conditional upon *x*. If the disturbances con-
3258 sist of impulses at random intervals, it is sometimes known as “shot noise.” (OECD
3259 Glossary of Statistical Terms)

3260 **non-deterministic noise** A random value that cannot be predicted.

3261 **non-ignorable bias** A bias introduced into data or an analytics procedure that results in a
3262 change that cannot be ignored.

- 3263 **non-public personal information** Information about a person that is not publicly known;
3264 called “private information” in some other publications.
- 3265 **personal identifier** Information with the purpose of uniquely identifying a person within
3266 a given context. (ISO/TS 25237:2008)
- 3267 **personal data** Any information relating to an identified or identifiable natural person (*data*
3268 *subject*). (ISO/TS 25237:2008)
- 3269 **personal information** See *personal data*.
- 3270 **personally identifiable information (PII)** Any information about an individual maintained
3271 by an agency, including (1) any information that can be used to distinguish or trace
3272 an individual’s identity, such as name, social security number, date and place of birth,
3273 mother’s maiden name, or biometric records; and (2) any other information that is
3274 linked or linkable to an individual, such as medical, educational, financial, and em-
3275 ployment information. [106](SP 800-122)
- 3276 **perturbation-based methods** Perturbation-based methods falsify the data before publica-
3277 tion by introducing an element of error purposely for confidentiality reasons. This
3278 error can be inserted in the cell values after the table is created, which means the
3279 error is introduced to the output of the data and will therefore be referred to as output
3280 perturbation, or the error can be inserted in the original data on the microdata level,
3281 which is the input of the tables one wants to create; the method with then be referred
3282 to as data perturbation—input perturbation being the better but uncommonly used
3283 expression. Possible methods are: rounding; random perturbation; [and] disclosure
3284 control methods for microstatistics applied to macrostatistics. (OECD Glossary of
3285 Statistical Terms)
- 3286 **privacy** Freedom from intrusion into the private life or affairs of an individual when that
3287 intrusion results from undue or illegal gathering and use of data about that individual.
3288 (ISO/IEC 2382-8:1998, definition 08-01-23)
- 3289 **privacy risk**
- 3290 **privacy loss** A measure of the extent to which a data release may reveal information that
3291 is specific to an individual.
- 3292 **privacy loss budget** An upper bound on the cumulative total privacy loss for individuals.
- 3293 **property value** Instance of a specific value together with an identifier for a data dictionary
3294 entry that defines a property. (ISO 8000-2:2012(E))
- 3295 **protected health information (PHI)** Individually identifiable health information: (1) Ex-
3296 cept as provided in paragraph (2) of this definition, that is: (i) Transmitted by elec-
3297 tronic media; (ii) Maintained in electronic media; or (iii) Transmitted or maintained
3298 in any other form or medium. (2) *Protected health information* excludes individu-
3299 ally identifiable health information in: (i) Education records covered by the Fam-

3300 ily Educational Rights and Privacy Act, as amended, [20 USC. 1232g](#); (ii) Records
3301 described at [20 USC. 1232g\(a\)\(4\)\(B\)\(iv\)](#); and (iii) Employment records held by a
3302 covered entity in its role as employer. (HIPAA Privacy Rule, 45 CFR 160.103). See
3303 <https://www.hhs.gov/hipaa/for-professionals/privacy/laws-regulations/index.html>.

3304 **pseudonymization** A particular type of de-identification that both removes the association
3305 with a data subject and adds an association between a particular set of characteristics
3306 related to the data subject and one or more pseudonyms.²⁹ Typically, pseudonymiza-
3307 tion is implemented by replacing direct identifiers with a pseudonym, such as a ran-
3308 domly generated value.

3309 **pseudonym** Personal identifier that is different from the normally used personal identifier.
3310 (ISO/TS 25237:2008)

3311 **quality** Degree to which a set of inherent characteristics fulfils requirements. (ISO 8000-
3312 2:2012(E))

3313 **quasi-identifier** A variable that can be used to identify an individual through association
3314 with another variable.

3315 **recipient** Natural or legal person, public authority, agency, or any other body to whom
3316 data are disclosed. (ISO/TS25237:2008)

3317 **redaction** The removal of information from a document or dataset for legal or security
3318 purposes.

3319 **re-identification** A general term for any process that restores the association between a
3320 set of de-identified data and a data subject.

3321 **re-identification risk** The likelihood that a third party can re-identify data subjects in a
3322 de-identified dataset.

3323 **re-identification rate** The percentage of records in a dataset that can be re-identified.

3324 **re-identificaiton probability** TBD

3325 **requirement** A need or expectation that is stated, generally implied or obligatory. (ISO
3326 8000-2:2012(E))

3327 **risk** A measure of the extent to which an entity is threatened by a potential circumstance
3328 or event, and typically a function of: (i) the adverse impacts that would arise if the
3329 circumstance or event occurs; and (ii) the likelihood of occurrence. (CNSSI No.
3330 4009)

3331 **risk assessment** The process of identifying, estimating, and prioritizing risks to organi-
3332 zational operations (including mission, functions, image, reputation), organizational

²⁹This definition is the same as the definition in ISO/TS 25237:2008, except that the word “anonymization” is replaced with the word “de-identification.”

3333 assets, individuals, other organizations, and the Nation, resulting from the operation
3334 of an information system. Part of risk management, incorporates threat and vulner-
3335 ability analyses, and considers mitigations provided by security controls planned or
3336 in place. Synonymous with risk analysis. (NIST SP 800-39)

3337 **safe harbor** Within the context of de-identification, refers to the Safe Harbor method for
3338 de-identifying protected health information in accordance with the Health Insurance
3339 Portability and Accountability Act (HIPAA) Privacy Rule. See [https://www.hhs.gov/hipaa/for-](https://www.hhs.gov/hipaa/for-professionals/privacy/special-topics/de-identification/index.html)
3340 [professionals/privacy/special-topics/de-identification/index.html](https://www.hhs.gov/hipaa/for-professionals/privacy/special-topics/de-identification/index.html).

3341 **statistical disclosure control** The set of methods to reduce the risk of disclosing informa-
3342 tion on individuals, businesses or other organizations. Such methods are only related
3343 to the dissemination step and are usually based on restricting the amount of or modi-
3344 fying the data released. (OECD Glossary of Statistical Terms)

3345 **suppression** One of the most commonly used ways of protecting sensitive cells in a table is
3346 via suppression. It is obvious that in a row or column with a suppressed sensitive cell,
3347 at least one additional cell must be suppressed, or the value in the sensitive cell could
3348 be calculated exactly by subtraction from the marginal total. For this reason, certain
3349 other cells must also be suppressed. These are referred to as secondary suppressions.
3350 (OECD Glossary of Statistical Terms)

3351 **synthetic data generation** A process in which seed data are used to create artificial data
3352 that have some of the statistical characteristics as the seed data.