

1 **Sharing of data from clinical** 2 **research projects – guidance from** 3 **the Swiss CTU network**

4 Authors (there might be changes in sequence)

5 Brigitta Gahl, PhD¹, Alan Haynes, PhD¹, Constantin Sluka, PhD², Elise Dupuis-Lozeron, PhD³,
6 Francisca Jörger, PhD⁴, Renate Schur, MSc⁴, Andri Christen, PhD⁵, Sven Trelle, MD¹

7

8

9

10 ¹ CTU Bern, University of Bern, Switzerland

11 ²Clinical Trial Unit, University of Basel, Department of Clinical Research, University Hospital,
12 Switzerland

13 ³ Clinical Research Centre, Department of Health and Community Medicine, University of
14 Geneva & Geneva University Hospitals, Geneva, Switzerland.

15 ⁴ Clinical Trials Center, University Hospital Zurich, Zürich, Switzerland

16 ⁵ Swiss Clinical Trial Organisation SCTO, Switzerland

17

18

19 Table of Contents

20

| | | |
|----|---|----|
| 21 | 1. Abstract..... | 4 |
| 22 | 2. Abbreviations..... | 5 |
| 23 | 3. Boxes with recommendations..... | 6 |
| 24 | 4. Introduction..... | 7 |
| 25 | 5. Legal basis in Switzerland..... | 9 |
| 26 | 6. Informed consent..... | 11 |
| 27 | 7. Data management plan..... | 13 |
| 28 | 8. Anonymization..... | 15 |
| 29 | 8.1. Goal..... | 15 |
| 30 | 8.2. Identifying variables..... | 15 |
| 31 | 8.3. The process of anonymizing data..... | 16 |
| 32 | 8.3.1 Assessment of the data..... | 16 |
| 33 | 8.3.2 Detailed specification of required data processing steps..... | 17 |
| 34 | 8.3.3 Data processing..... | 17 |
| 35 | 8.3.4 Quality control..... | 17 |
| 36 | 9. Data structure and format..... | 20 |
| 37 | 9.1. Data structure..... | 20 |
| 38 | 9.2. Data format..... | 20 |
| 39 | 9.3. Character encoding..... | 21 |
| 40 | 10. Coding of variables..... | 22 |
| 41 | 10.1. Variable types..... | 22 |
| 42 | 10.2. Variable labels..... | 23 |
| 43 | 10.3. Time structures in the data sampling..... | 24 |
| 44 | 11. Metadata and documentation..... | 25 |
| 45 | 11.1. Metadata schemes..... | 25 |
| 46 | 11.2. Additional documentation..... | 26 |
| 47 | 11.3. Is statistical analysis code needed for data sharing?..... | 27 |
| 48 | 12. Version control..... | 28 |
| 49 | 13. Selection of repository..... | 29 |

| | | | |
|----|--------|--|----|
| 50 | 13.1. | Principles..... | 29 |
| 51 | 13.2. | Time point..... | 30 |
| 52 | 13.3. | Identifying potential repositories | 30 |
| 53 | 13.4. | Selection criteria | 30 |
| 54 | 14. | Requesting and use of data..... | 32 |
| 55 | 15. | References | 34 |
| 56 | 16. | Glossary..... | 37 |
| 57 | 17. | Appendix | 44 |
| 58 | 17.1. | Further detailed specification of required data processing steps | 44 |
| 59 | 17.1.1 | Example data to be considered for deletion | 44 |
| 60 | 17.1.2 | Examples and details on manipulations to decrease precision..... | 44 |
| 61 | 17.2. | Further details on coding of variables..... | 46 |
| 62 | 17.2.1 | Formatting of date and time variables..... | 46 |
| 63 | 17.2.2 | Examples for further documentation of the dataset | 47 |
| 64 | 17.3. | Meta data scheme from ISRCTN..... | 51 |
| 65 | 17.4. | Information required for additional documentation..... | 53 |
| 66 | 17.5. | Checklist for selecting a data repository | 54 |
| 67 | | | |
| 68 | | | |

69 1. Abstract

70 **Objectives:** Data sharing has become a requirement of many funding bodies and is becoming
71 a scientific standard in many disciplines. In medical research, however, data sharing can conflict
72 with clinicians' obligation to protect patients' privacy. General recommendations on data
73 sharing exist also for clinical research, but so far lack practical and Swiss-specific aspects. The
74 objective of this document is to provide practical recommendations for all relevant aspects of
75 data sharing in agreement with legislation in Switzerland.

76 **Methods:** This document was written by members of the Swiss CTU Network, a network of
77 academic clinical trial units. The process did not follow a formalized Delphi process. After an
78 internal consensus round, this report is now published as pre-print for external review. A second
79 version will incorporate external comments.

80 We plan to publish this document as a text in progress, as we expect relevant changes in related
81 fields such as the development of further dedicated medical repositories or methodological
82 advances in anonymization techniques.

83 **Results:** We developed principles and practical recommendations with respect to informed
84 consent, data management plan, anonymization, data structure and format, coding of variables,
85 metadata and documentation, version control, selection of repository, requesting and use of data.
86 We also provide a summary of legal aspects relevant for the Swiss context.

87 **Conclusions:** The intension to share data has an impact not only after a clinical trial or an
88 observational study is completed, but also during the planning period, the conduct and the
89 analysis phase. Clinical researchers need to be aware at the beginning of a study on how to
90 inform patients and at least the amount of work related to preparing data for sharing, metadata,
91 and any further documentation. This report provides details of aspects to be considered,
92 suggests decision criteria, and provides examples and checklists, in order to support data
93 sharing in practice.

94

95 2. Abbreviations

| | | |
|-----|--------------|---|
| 96 | ADaM..... | Analysis Data Model |
| 97 | AHV | Alters- und Hinterlassenversicherung (<i>Old Age Insurance</i>) |
| 98 | API | Application Protocol Interface |
| 99 | CDASH | Clinical Data Acquisition Standards Harmonization |
| 100 | CDISC | Clinical Data Interchange Standards Consortium |
| 101 | CERN | Conseil Européen pour la Recherche Nucléaire (<i>European Organization for</i> |
| 102 | | <i>Nuclear Research</i>) |
| 103 | ClinO | Clinical Trials Ordinance |
| 104 | CRF | Case Report Form |
| 105 | CSV | Comma Separated Value |
| 106 | CTU | Clinical Trials Unit |
| 107 | DAC | Data Access Committee |
| 108 | DMP | Data Management Plan |
| 109 | DOI..... | Digital Object Identifier |
| 110 | ECRIN..... | European Clinical Research Infrastructure Network |
| 111 | ELSI | Ethical, Legal and Social Implications |
| 112 | EOSC..... | European Open Science Cloud |
| 113 | EU..... | European Union |
| 114 | FAIR..... | Findable, Accessible, Interoperable, Reuseable |
| 115 | FADP..... | Federal Act of Data Protection |
| 116 | GCP | Good Clinical Practice |
| 117 | HIPAA..... | U.S. Health Insurance Portability and Accountability Act |
| 118 | HRA | Human Research Act |
| 119 | HRO | Human Research Ordinance |
| 120 | ICMJE | International Committee of Medical Journal Editors |
| 121 | ICPSR..... | Inter-university Consortium for Political and Social Research |
| 122 | ICTRP..... | International Clinical Trials Registry Platform |
| 123 | ID..... | IDentificator |
| 124 | IP | Intellectual Properties |
| 125 | IPD | Individual Participant Data |
| 126 | ISRCTN..... | International Standard Randomised Controlled Trials Number |
| 127 | KDSG..... | Kantonales Datenschutzgesetz (<i>Cantonal Act of Data Protection</i>) |
| 128 | MedDRA | Medical Dictionary for Regulatory Activities |
| 129 | PID | Patient IDentificator |
| 130 | SAP..... | Statistical Analysis Plan |
| 131 | SDTM..... | Study Data Tabulation Model |
| 132 | SNOMED | Systematized Nomenclature of Human and Veterinary Medicine |
| 133 | SPHN..... | Swiss Personalized Health Network |
| 134 | TSV | Tab Separated Value |
| 135 | UK..... | United Kingdom |
| 136 | URL..... | Uniform Resource Locator |
| 137 | US..... | United States |
| 138 | UTF | Unicode Transformation Format |
| 139 | WHO | World Health Organisation |
| 140 | XML..... | eXtensible Markup Language |
| 141 | ZIP..... | Zone Improvement Plan |

142 **3. Boxes with recommendations**

| | | |
|-----|---|----|
| 143 | Box 1: Recommendations concerning consent | 12 |
| 144 | Box 2: Recommendations concerning the data management plan | 13 |
| 145 | Box 3: Recommendations concerning anonymization..... | 18 |
| 146 | Box 4: Recommendations on data structure and format | 21 |
| 147 | Box 5: Recommendation concerning variables within a shared dataset | 24 |
| 148 | Box 6: Recommendation for metadata and additional documentation | 26 |
| 149 | Box 7: Recommendation regarding availability of analysis code..... | 27 |
| 150 | Box 8: Recommendations selection of repository | 31 |
| 151 | | |
| 152 | | |

153 4. Introduction

154 Sharing of research data^{Glossary}¹ has evolved as standard practice in many disciplines. There are
155 two main drivers for the current data sharing movement: one is to enable reproducibility of
156 research results with the goal of increasing transparency and credibility of science (1); the other
157 is to enable reuse of data for new research questions. The International Committee of Medical
158 Journal Editors (ICMJE) considers it an ethical obligation to responsibly share data generated
159 by clinical trials (2). The main reason for this opinion is that trial participants have put
160 themselves at risk by accepting to receive a treatment under study. Slightly different means are
161 needed to achieve these two goals with respect to the amount of data to be shared and with
162 respect to additional documentation that comes with the shared data.

163 The purpose of this document is to give specific recommendations for each decision that has to
164 be made when sharing data from a clinical trial or an observational study, be it individual
165 participant data (IPD) or aggregate data. A principal investigator (PI) or other study team
166 members who intend to share data should find answers to their questions in this document and
167 whatever else there is to consider. We provide minimal options to fulfil current requirements to
168 still allow for the culture of data sharing to evolve.

169 Important papers on sharing individual patient data were published. Ohmann et al. developed a
170 consensus document for sharing individual patient data using a consensus-building process
171 among an interdisciplinary task force of research professionals as part of an European project
172 (3). The document provides 10 principles and 50 recommendations to support data sharing and
173 remove obstacles on many different levels such as collaboration culture and incentives, but also
174 on technical and organizational aspects for “making data sharing a reality” (3). Our own
175 statement is rather dedicated to the reality of data sharing clinical researchers are facing. It is
176 our conviction that the establishment of data sharing will affect collaboration among clinical
177 researchers and involvement of the research community.

178 The FAIR principles (4–6) provide guidelines to improve the Findability, Accessibility,
179 Interoperability, and Reuse of data. They were developed for scientific data in general and focus
180 on machine-operability and the order of the letters represent the dependency of the principles,
181 e.g. data must be Findable to be Accessible, must be Findable and Accessible to be
182 Interoperable and must be Findable, Accessible and Interoperable to be Reuseable. Even though
183 it is very normal for us all to search for digital objects such as scientific papers in a database,
184 this is more complicated when it comes to data objects^{Glossary}/artifacts^{Glossary}. When retrieving
185 data, a whole package of related descriptions and documentation is needed to understand the
186 data and allow its reuse. As a consequence, structure is needed and thus more rules for data
187 provision. The choice of repository already determines many aspects of findability and
188 accessibility. Usually, a repository has a metadata^{Glossary} scheme (see sections 11 and 13) that
189 might be specific to the field and hence allows for specific searches. The repository might be
190 linked to other systems to allow for parallel searches in several repositories (see section 11).

¹ Terms that are defined and further explained in the glossary, are marked with ^{Glossary} whenever they appear for the first time in the text.

191 Accessibility follows from the data requesting process as defined by the repository.
192 Interoperability of data basically relates to the format, structure, coding, and documentation and
193 is covered in sections 9 and 10. Reuse of data depends on many factors and, particularly, of the
194 three other principles. We want to point out that *reproducibility* in the context of shared data
195 might not mean that researchers redoing exactly the same analyses will end up with exactly the
196 same result for each estimate. If precise data such as biomarkers have been jittered or grouped
197 as a means of anonymization^{Glossary}, derived and estimated values might differ from published
198 values. It is important that this difference is mentioned and quantified in the documentation.

199 This document was written by Swiss professionals in the field of academic clinical research.
200 Involved persons were identified within the Swiss CTU network and delegated from each
201 clinical trials unit participating in the network. The authors identified relevant topics to be
202 covered and assigned each topic to an individual author. During the writing period, further
203 topics were identified and added. The document was merged and the different parts were
204 consolidated by three members of CTU Bern. Then all authors were asked for feedback to the
205 entire document. The three members of CTU Bern incorporated all feedback reaching
206 consensus among all authors in most aspects (unless explicitly indicated). Afterwards this
207 document underwent language review and is now to be published on <https://www.medrxiv.org/>
208 for invited expert review, after another co-authors round. The whole process did not follow a
209 structured Delphi process.

210 We plan to publish this document as a text in progress, as we expect relevant changes in related
211 fields such as the development of further dedicated medical repositories or methodological
212 advances in anonymization techniques.

213

214 5. Legal basis in Switzerland

215 Health-related personal data are considered sensitive data in Switzerland. According to Article
216 4 paragraph 3 of the Federal Act of Data Protection (FADP) (7), personal data may only be
217 used for the purpose a) indicated to subjects at the time their data are collected, b) that is evident
218 from the circumstances, or c) that is required by law. The use of health-related personal data
219 for research purposes is specifically laid down in a so-called special law, the Human Research
220 Act (HRA) (8).

221 The Act regulates biomedical research on human subjects at the federal level and is based on
222 internationally recognized principles. Sharing health-related data fulfills criteria for Further
223 Use² (9) according to the Act and is regulated by Chapter 4 of the HRA (Art. 32-35). Further
224 Use presupposes that the data are already available, i.e. collected with the necessary justification
225 for another purpose and stored and made available (Art. 24 Human Research Ordinance, HRO,
226 CC 810.301). If data sharing is planned at the time of data collection, e.g. for a clinical trial, the
227 participants must be informed and consent obtained about the intended reuse of the collected
228 data and their right to dissent to that at the time of collection. Article 17 of the HRA states: "If
229 the intention exists to make further use for research of ... health-related personal data collected,
230 the consent of the persons concerned must be obtained at the time of such sampling or collection,
231 or they must be informed of their right to dissent." However, consent for further use/sharing of
232 data should not be an inclusion criterion for a trial; individuals must be given the possibility to
233 participate in a trial without giving consent for data sharing later. In exceptional cases and under
234 given circumstances (e.g. approval by an ethics committee), the law allows the reuse of health-
235 related data for research that was collected without explicit consent provided it is impossible or
236 very difficult to obtain consent or to provide information on the right to dissent, or this would
237 impose an undue burden on the person concerned. In addition, the privacy and fundamental
238 rights of the individuals must always be ensured (Art. 34 HRA). If the intention is to share only
239 coded data and the data do not contain genetic data, information about potential further use is
240 sufficient unless a participant explicitly disagrees. Explicit consent is not required (Art. 33
241 HRA).

242 If personal data are disclosed abroad, adequate data protection must be ensured (Art. 6 Federal
243 Act of Data Protection). Adequate data protection should be part of any data use agreement (see
244 section 12).

245 Anonymous data, which are not personal and cannot harm persons, are subject to neither FADP
246 nor HRA, may be freely shared. However, as described below it is typically not possible to
247 ensure that individual patient data will remain anonymous for all times to come (see section 8).

Infobox 1: Swiss legal basis in a nutshell

1. Data sharing is considered further use.

² The concept of Further Use also applies to biological material but this is not discussed in this statement. The statement is specific for data sharing aspects and does not cover other aspects of Further Use.

2. Consent for data sharing should preferably be obtained at enrolment.
3. Anonymous data does not fall under FADP nor HRA. It is unlikely that individual patient data of a clinical study can be anonymized^{Glossary}.

249 6. Informed consent

250 The sharing and use of personal health data for research has implications for patients' rights
251 and interests. The legal requirements for patient information and consent are laid down in the
252 Swiss Federal Act on Data Protection (FADP) and the Human Research Act (HRA) (see section
253 5).

254 The Ethical, Legal and Social Implications (ELSI) advisory group, which is part of the Swiss
255 Personalized Health Network (SPHN) (10) initiative, published a framework providing ethical
256 guidance on processing and sharing personal data within SPHN hereafter referred to as the ELSI
257 framework (11). The document takes into account both international guidelines and national
258 law including the HRA with a specific focus on aspects of general consent: "[The] Framework
259 refers to all data types ... that can be employed in the context of health research. This includes
260 health-related personal data ... that were not originally collected for research purposes, ...".
261 The ELSI advisory group considers a general consent (Broad Consent) sufficient for further use
262 of encoded data outside the institution regardless of the original collection purpose and whether
263 data are genetic or nongenetic (*ELSI framework III-1, Guidelines point b*). It is important in this
264 context to have an unambiguous understanding of the term *general consent*. This term is often
265 used in biobanking and related to further use of health-related data and samples collected in
266 routine medical care (12). As described in section 5, sharing data from clinical research projects
267 requires explicit informed consent because the consent given by the patient allows the use of
268 the data to answer the questions/objectives of the project and does not extend to other research
269 purposes. A general consent that was given in the context of routine medical care, for example
270 at time of admission to a hospital, is insufficient for the purpose of sharing clinical trial data.
271 The ELSI Advisory Group provides a broader definition of the term, and states that general
272 consent means "informed consent of a research participant to unspecified further research uses
273 of his or her health-related personal data or human biological material" (in the international
274 academic literature, the closest term to general consent is broad consent). See (9) 3.3.1 for a
275 definition of "further use." In this sense, the framework is applicable to the sharing of clinical
276 trial data. As described in section 5, information and consent about possible data sharing should
277 be done at project enrolment.

278 Sharing of coded or personal health-related data requires that the transfer of data is traceable at
279 any time to ensure the patients' personal rights to provide information on the type, storage, and
280 reuse (sharing) of her/his data on request or to ensure that data will no longer be available for
281 research if the consent for reuse is revoked (*ELSI framework III-1, III-4*). This is only feasible
282 if the data are either anonymized (which is in general not achievable, see section 6) or if data
283 are shared on the basis of a contract that we consider the default (see section 14). The sponsor-
284 investigator providing data needs adequate governance structures in place to maintain control
285 over the data such as data sharing agreements^{Glossary} specifying the intended use, confidentiality,
286 and the obligation to delete data of persons revoking consent and compliance with data
287 protection. As in all situations, revoked consent has to be immediately addressed (*ELSI*
288 *framework* says "revocations [...] are swiftly acted upon"), but not retroactively. Specifically,

289 the patient consent status at the moment of database export is relevant. If a patient does not give
290 consent, it should be documented when the patient was asked and what he or she was informed
291 about.

Box 1: Recommendations concerning consent

- R1. Sponsor(-investigator)s must ensure that participants are informed about potential data sharing and further use of their data at the time of enrollment in a clinical research project including anonymization of their data.
- R2. If sharing of coded data is planned (the usual case):
- a. Sponsor(-investigator)s must ensure that potential participants are informed about the potential sharing of their data. Explicit consent is not needed but the possibility to disagree must be ensured.
 - b. Sponsors/sponsor-investigators should ensure that a system is in place that allows access to this information centrally, e.g., by recording disagreements in the study database.
 - c. When sharing coded data, the sponsor(-investigator) must have a system in place to ensure participants' rights especially with regard to the use of their data and deletion of the data wherever the data were transferred to.
- R3. If sharing of uncoded personal data is planned (not recommended):
- a. Sponsor(-investigator)s must ensure that potential participants are informed about the potential sharing of personal data and the potential anonymization of their data for this purpose. Explicit written informed consent should be sought.
 - b. Sponsors/sponsor-investigators should ensure that a system is in place that allows access to consent status of each patient centrally, e.g., by recording the information in the study database.
 - c. When sharing personal data, the sponsor(-investigator) must have a system in place to ensure participants' rights especially with regard to the use of their data and deletion of the data wherever the data were transferred to.
- R4. For sharing data collected in the setting of clinical routine a general consent of a patient is sufficient unless it explicitly excludes data sharing, the general consent used in the hospital has to be carefully checked.
- R5. It is imperative to take into account the consent status of patients. If a patient withdraws consent, data of this patients have to be ignored immediately from the moment of withdrawal on, but analysis already done or data files already provided do not have to be changed.

292

293

294 7. Data management plan

295 According to (13)(v3.1.0), a Data Management Plan^{Glossary} is a document "to identify the overall
 296 strategy for data management processes for the trial; a compilation of documents that may
 297 include amendments/appendices but are not limited to: Completion Guidelines, Data Quality
 298 Plan, CRF Design Document, Database (build) Specification, Entry Guidelines, Database
 299 Testing". The Data Management Plan therefore provides an overview of all aspects related to
 300 data (management) in a clinical research project. Depending on the details provided in the study
 301 protocol, a Data Management Plan might not be needed. However, we recommend that all
 302 studies have a Data Management Plan because this supports and facilitates later data sharing
 303 activities. Several templates for such a document are freely available over the Internet. We do
 304 not recommend a particular one. However, the plan should cover the aspects relevant for data
 305 collection, handling, and storage during study conduct (and implementation/conclusion) as well
 306 as for data sharing. A possible structure and description of content is shown below. It should
 307 be noted that there are now specific journals that specialize in publishing articles on description
 308 of datasets and aspects of data management. We make no specific recommendations on this.

Box 2: Recommendations concerning the data management plan

- R6. All aspects related to data management including data sharing should be documented before conducting a clinical research project. The document should be considered a living document and regularly updated using a version control system. It might be called Data Management Plan.
- R7. Possible structure and content of a data management plan. Not all sections will be relevant to all research projects:
1. Introduction
 2. Responsibilities
 3. Description of collected/generated data
 4. Case Report Form^{Glossary} development
 5. Clinical Data Management System – study specific implementation
 - 5.1. Implementation of the study database in the Clinical Data Management System
 - 5.1.1. Codebook development
 - 5.1.2. Clinical Data Management System implementation
 - 5.1.3. Medical coding
 - 5.1.4. Data import
 - 5.2. Verification of Clinical Data Management System setup and deployment
 - 5.3. Change management
 6. Clinical Data Management System – infrastructure
 - 6.1. Data storage
 - 6.2. Data back-up
 - 6.3. Access to the data

- 6.4. Granting access to the productive version of the Clinical Data Management System and database
7. Data collection
 - 7.1. Pre-requisites for data entry
 - 7.1.1. Data entry guidelines
 - 7.1.2. Training of users and training documentation
 - 7.2. Entering data
8. Quality control procedures
 - 8.1. Real-time data validation
 - 8.2. External data validation (offline checks)
 - 8.3. Central data monitoring
 - 8.3.1. Definition of Key Performance Indicators (KPIs)
 - 8.3.2. Frequency
 - 8.3.3. Reporting
 - 8.3.4. Clinical Data Management System generated, automatic queries
 - 8.3.5. Manual queries
 - 8.3.6. Follow-up on (persisting) data discrepancies
9. Database closure
 - 9.1. Pre-closure data checks
 - 9.2. Quality assurance audit and database lock
 - 9.3. Database unlock
10. Data transfer and exports
 - 10.1. Data requests and transfer
 - 10.2. Data exports
 - 10.3. Export validation
 - 10.4. Adverse event data reconciliation
11. Clinical Data Management System archiving and provision of final materials to the sponsor
12. Data preservation
13. FAIR data sharing
 - 13.1. Repository ^{Glossary}
 - 13.1.1. Shared artifacts
 - 13.2. Data request process
 - 13.3. Ethics, legal and security issues
 - 13.3.1. Data protection
 - 13.3.2. Copyright and intellectual property

311 8. Anonymization

312 8.1. Goal

313 Anonymization is the process of handling personal data in such a way that identification of
314 individual persons is impossible or possible only with disproportionate effort. Further data
315 sources and technologies for data linkage might become available at some point, thus the effort
316 needed to identify persons is not known for all times to come (9). As a consequence, data that
317 are anonymized today might not remain anonymized according to this definition. Anonymized
318 data in the strictest sense can be shared openly, but the claim that individual patient data are
319 anonymized is not realistic (14,15). The process we describe below aims at sharing data with
320 researchers in a standardized and institutionalized way based on a standard contract or license
321 in which the data requester agrees not to try to identify patients, not to give the data to other
322 persons, and to maintain data security. In this setting and with these restrictions, present-day
323 anonymization can be considered acceptable. The goal of the anonymization process we
324 describe is to protect participants' privacy to a degree that criminal acts would be necessary for
325 identifying patients at time of sharing the data to identify participants.

326 Obviously, the anonymization process consists of manipulations that *change* the data. In order
327 not to spoil the benefit of data sharing, it is important to consider the goal of anonymization:
328 protection of patients' privacy, while also considering the usefulness of the data. Of note,
329 anonymizing data needs a lot of work.

330 8.2. Identifying variables

331 Variables^{Glossary} are called directly identifying if they contain personal information by which a
332 participant can be identified with little or no effort, and should in general not be stored within
333 the study database or, if stored, not be possible to export. The Human Research Ordinance
334 mentions explicitly the following data (Art. 25, Paragraph 2): name, address(es), date of birth,
335 unique identification numbers. The U.S. Health Insurance Portability and Accountability Act
336 (HIPAA) provides more details. The following is a non-exhaustive list (16):

- 337 • Real names
- 338 • All elements of dates (except year) for dates that are directly related to an individual,
339 including birth date, admission date, discharge date, death date, and all ages over 89
340 and all elements of dates (including year) indicative of such age, except that such ages
341 and elements may be aggregated into a single category of age 90 or older.
- 342 • Addresses and geolocations/-codes past and present (canton/state might be allowed
343 given that the geographic unit contains more than 20,000 persons; MEDSTAT regions
344 might be more appropriate as they were designed to ensure anonymity (17)).
- 345 • Telephone number, email addresses, IP addresses, or any links or
346 aliases/pseudonyms^{Glossary} e.g. Facebook, LinkedIn, WhatsApp, Twitter, or
347 links/URLs to personal websites.

- 348 • Device/implant identifiers^{Glossary} and serial numbers or vehicle identifiers and serial
 349 numbers, including license plate numbers.
- 350 • Any other (non-health) personal identifier (ID), e.g. hospital ID (or PID), social
 351 security numbers (AHV), insurance numbers, passport numbers, account numbers, etc.
- 352 • Full-face photographs and any comparable images or biometric identifiers including
 353 finger and voice prints.

354 Identifying data^{Glossary} can be variables containing information that are by definition unique to
 355 the patient, and therefore the patient can be identified with medium effort, e.g. genetic,
 356 genomics, metabolomics, proteomics, micro-array, biomarkers, or similar high-precision data.

357 Identifying data can be variables containing information which singly or in combination with
 358 other data, can be used to identify the patient with some effort (indirect identifiers), e.g.:

- 359 • Marker of rare disease or subtype of disease
- 360 • Rare medication, treatment, or surgery
- 361 • Rare diagnostic tool or machine used
- 362 • Rare population
- 363 • High-precision variable (while precision depends on the type of data)
- 364 • Any unusual variation or combination of variables as mentioned above

365 **8.3. The process of anonymizing data**

366 Anonymization is a multistep process that requires input by several people, among them the
 367 sponsor and the statistician. The shared data set should in principle contain only the data that
 368 are needed for the intended purpose. For example, to share a dataset underlying a scientific
 369 report only the data needed to reproduce the statistics, graphs, tables etc. in the report should
 370 be in the dataset³.

371 **8.3.1 Assessment of the data**

372 It is necessary to assess the whole dataset with all individual variables. This is best done by a
 373 statistician or data manager and by the sponsor (because content knowledge might be needed).
 374 HIPAA states three criteria relating to a variable or a set of variables that might serve as
 375 guidance to assess the risk of re-identification:

- 376 1. **Replicability:** How consistently is a piece of information related to a specific person?
 377 For example, while laboratory values vary (low replicability), demographics are more
 378 stable (high replicability).

³ Note that (3) refers to the danger that records in a shared data file might be selected because they are “supporting the conclusions of a specific published paper“ (p. 2). This is not an issue of this paper because CTUs conduct analyses according to pre-specified inclusion criteria and are neutral with respect to expected results.

- 379 2. Data source availability: Which external data sources could be used to identify a
 380 specific person? For example, demographics could be obtained from public registries.
 381 3. Distinguishability: How many persons share a specific combination of characteristics?
 382 For example, year of birth and canton is less likely to be unique than complete date of
 383 birth and ZIP code.

384 These criteria are relevant to assess the risk of a linkage attack, the process of re-identification
 385 by linking an external data source with person-identifying data to the original data set. In the
 386 last decades, several cases of successful linkage attacks have been recorded (18). For example,
 387 in 2013 5–7 laboratory values from a known patient were shown to identify the corresponding
 388 records in a de-identified^{Glossary} biomedical research database (19).

389 Each variable should be classified whether it is:

- 390 • (Potentially) Directly identifying (see section 8.2),
- 391 • Indirectly identifying, i.e. identifying in connection with other variable(s). The other
 392 variable(s) should be documented, or
- 393 • Unproblematic, i.e. neither directly nor indirectly identifying.

394 **8.3.2 Detailed specification of required data processing steps**

395 After categorization, the necessary data preparation steps for the directly and indirectly
 396 identifying variables must be defined. This is a non-exhaustive list of potential procedures:

- 397 • Deletion: Variables containing directly identifying information unsuitable for
 398 manipulation must be deleted. The appendix provides some examples.
- 399 • Irreversible pseudonymization^{Glossary}: Irreversible pseudonymization is a
 400 transformation of a variable into a new variable, where the mapping which renders the
 401 process reversible is deleted (database dependent). This usually requires a complex
 402 algorithm and is rarely used.
- 403 • Manipulations to decrease precision: Too much precision bears the risk of making
 404 entries linkable to persons. Possible methods to decrease precision include relative
 405 time in the course of the study instead of precise dates and times, rounding of
 406 continuous data, grouping and aggregation (categorization), introducing random noise
 407 (jittering, perturbation), setting certain values to missing (suppression), data swapping,
 408 resampling or subsampling.

409 The Appendix provides additional details and examples.

410 **8.3.3 Data processing**

411 The steps as defined in 8.3.2 have to be programmed using statistical software and a set of new
 412 data files has to be generated.

413 **8.3.4 Quality control**

414 Two persons should perform a quality control and check the de-identified data:

- 415 1. Sponsor
- 416 In particular to check:
- 417 • Whether the de-identified data set contains free text variables, in which the
- 418 text may potentially lead to identification
- 419 • Whether this data set contains other variables which may alone or in
- 420 combination lead to identification, in particular if infrequent/rare disease or
- 421 population is involved
- 422 • Whether data need to be lumped into categories
- 423 2. Statistician or representative knowledgeable of the data set (e.g. Central Data
- 424 Monitor, Monitor, Data Manager)
- 425 In particular to check:
- 426 • That any combination of indirectly identifying variables results in a
- 427 number >1 (e.g. five) records
- 428 • Whether the de-identified data do not contain personal information variables
- 429 except age without any digit (but not date of birth)
- 430 • Whether the file only contains text variables if specifically requested and that
- 431 those text variables are appropriately redacted
- 432 • Whether digits have been removed/rounded/jittered
- 433 • Whether dates have been replaced
- 434 • Whether the identification numbers have been replaced with a new random
- 435 identifier
- 436 Whether results based on the new dataset are similar to results using the original dataset must
- 437 be checked, and if not, where and to what extent they deviate and any deviations should be
- 438 noted in the same document where the assessment and specifications are described (steps 1 and
- 439 2). Every analysis need not be run. Common sense should be applied to select important ones.
- 440 The statistician corrects the de-identification^{Glossary} coding according to the recommendations
- 441 resulting from quality control.
- 442 If all is in order, the two persons sign a quality control document with a date to document that
- 443 they did the quality control and what was checked. If multiple (repeated) exports need to be
- 444 done based on the same code, then this quality control needs to be done only once, except if the
- 445 sponsor requests a check at each export.

Box 3: Recommendations concerning anonymization

R8. Anonymization should involve at least the sponsor and the statistician.

- R9. Directly identifying variables should be removed, IDs should be replaced by random numbers, string variables should be removed, and rare combinations of values identified and lumped together to achieve larger groups of patients.
- R10. The anonymization process should be quality controlled and appropriately documented.

446 9. Data structure and format

447 Full descriptive information of the data is necessary (see the coding variables section 8) for
448 reproduction of analyses as well as for reuse of the data, which are the two main purposes of
449 data sharing. Details of the de-identification process should be provided for the sake of
450 transparency.

451 Although the European Clinical Research Infrastructure Network (ECRIN) recommends the
452 Clinical Data Interchange Standards Consortium (CDISC) format for sharing data (3), the use
453 of this standard outside of the pharmaceutical industry is relatively rare, particularly in the
454 academic setting where resources to set up CDISC-compliant databases are limited. While we
455 agree that standardization of items and structure aids secondary data processing and reuse, the
456 current reality is that academic databases are rarely (if ever) designed to CDISC standards.
457 Furthermore, the CDISC defines a variety of formats such as the Study Data Tabulation Model
458 (SDTM) and Analysis Data Model (ADaM) on the database side, and seven different extensible
459 markup language (XML) based formats for data exchange. It is therefore a substantial challenge
460 to understand the full CDISC standard structure, let alone work with it. That being said,
461 utilization of certain features of the format is recommended (such as standardized variable
462 naming and encoding). The ECRIN statement highlights that it is difficult to transfer data to a
463 specific standard unless this is done from the project planning stage. Thus, as far as is possible
464 given constraints of cost, time to implement, and technical capabilities CDISC standards should
465 be employed for new trials at the database design stage.

466 9.1. Data structure

467 Clinical research projects typically involve multiple assessments over time (at least two
468 different time points). Data in a study are usually collected on different forms within the case
469 report form. The structure of the database usually reflects this structure, i.e. data are stored in
470 separate tables and keys^{Glossary} serve as the link between these tables (relational database). We
471 recommend that the table structure is preserved when preparing a dataset for sharing, that is,
472 each table remains a separate file within the dataset. Careful description of the keys is needed
473 to ensure that users of the data are able to establish the correct link across the different files (see
474 section 11). The original key-value pairs will usually be replaced with new random unique
475 identifiers (see section 8).

476 9.2. Data format

477 For older projects, where CDISC standards were not considered, data would ideally be shared
478 in a simple format. Text based comma separated value (CSV) or variants thereof (e.g. tab
479 separated value, TSV) are non-proprietary formats which should be future-proof: changes in
480 future versions of software will not render the data unreadable as they are text based formats.
481 Other formats such as XML, while offering the ability to include data, audit trail, coding and

482 database structure, are potentially more difficult to work with. Indeed, some widely used
483 statistical software packages have only very basic XML capabilities. Additionally, the FAIR
484 principles suggest that the data should be usable by most users. Using formats such as XML
485 requires a large degree of specialist knowledge simply to read the data into statistical software.
486 Proprietary formats such as SAS, SPSS, Stata data files, and .xlsx files are also less suitable for
487 sharing as they are generally only accessible using that software (although there are packages
488 available for R to handle many formats), and are typically not suited to long-term storage due
489 to changes between versions. As such, text-based formats such as CSV are preferable. There is,
490 however, some variation in recommendations in this respect. While some institutional
491 repositories recommend plain-text-based formats (Georgia State University, World Wide Web
492 Consortium), many others recommend proprietary formats (Inter-university Consortium for
493 Political and Social Research, ICPSR) or a wide range of formats including text-based and
494 proprietary formats (Oregon State University, Stanford University, UK Data Service). Most
495 also suggest delimited text (e.g. comma separate value format) with setup files (codes to read
496 data in and prepare it). However, setup files are containing software code and the programming
497 language dictate which programs can use the data. In general, data are more ready to use only
498 if in the format of the statistical software used for the original statistical, and data are more
499 accessible in any non-proprietary format. We therefore recommend that data are provided in
500 the format as used during analysis and in comma-separated value format. Metadata and
501 documentation should be uploaded in separate files along with the data (see sections 10 and 11).

502 **9.3. Character encoding**

503 The encoding of files is also an issue, as it determines how special characters (e.g. ä, à, é, è, ö,
504 ü) are interpreted by software. We recommend 8-Bit UCS Transformation Format (UTF-8)
505 (<https://en.wikipedia.org/wiki/UTF-8>) (20) encoding where possible, as this is a widely
506 recognized encoding system and supports the vast majority of characters. The encoding used
507 should be explicitly stated, ideally in the data management plan.

Box 4: Recommendations on data structure and format

- R11. Retain the database structure in the shared data (five case report forms in the database make five tables in the shared data).
- R12. Use text based formats such as CSV to share data, encoded in UTF-8.
- R13. Also provide data in the original format.

508

509 10. Coding of variables

510 The way data are prepared for sharing affects its general usability as well as its interoperability.
511 For data sharing purposes, as few changes as possible should be made to a dataset after
512 exporting the data from the database as it may not be possible to anticipate all the ways in
513 which data might be used further. Thus, in order to avoid wasted effort, we advise not to recode
514 data for data sharing purposes (within the limitations imposed by anonymization, see section
515 8). The use of standardized or controlled vocabularies (e.g. SNOMED, MedDRA, CDASH)
516 increases the interoperability of data. Therefore, we recommend the use of standardized
517 vocabulary. However, this should be considered during database development, rather than
518 coding the data afterwards. Some data manipulation and recoding is inevitable, though, when
519 sharing data.

520 10.1. Variable types

521 Individual variables come usually in four main types: date/time, text, numeric and categorical
522 (binary and ordinal variables can be thought of as special cases of categorical variables). Each
523 type of variable should be handled in a specific manner.

524 **Date variables** should be converted into project days (i.e. days since informed consent or
525 randomization, see section 8). There might be circumstances in which dates/times are necessary
526 such as when seasonal effects are important, as are relationships to historical events. Under
527 such circumstances, we recommend a slightly modified version of the ISO 8601 standard.
528 Date/time variables can be subdivided into three units: date, time, and date-time, each requiring
529 its own handling. The appendix contains further details on formatting standards.

530 **Continuous variables** are relatively simple; they should be provided as they are (e.g., 1.5). The
531 number of decimal places should be the same for all observations (if the most precise
532 observation is 1.5, then all observations should have one decimal place: 1.0 instead of 1). Note
533 that it may be desirable to reduce the precision of some variables (see section 8).

534 **Categorical variables** comprise binary (yes/no), single-choice (male/female), multiple-choice
535 or ordinal type variables (New York Heart Association Functional Classification scores to
536 classify heart failure). They can typically be provided in two ways: a textual description (such
537 as male/female or yes/no) or a numeric representation (e.g., 1 or 2). From a human readability
538 perspective, it would likely be best to save the textual representation, but data saved in such a
539 manner will typically be considerably larger than that saved with the numeric representation
540 instead, and require more work to make it analyzable. It is thus preferable to save the numeric
541 codes with an additional codebook to provide the meaning of the codes (see Appendix for an
542 example: Table 3, p. 48). The codebook can then be used by statistical software to label the data
543 when it is to be reused, albeit with a little programming. Multiple choice questions should be
544 split into as many binary variables as there are options, e.g., if there are options of diabetes,
545 previous myocardial infarction, and previous stroke there would be three binary variables,

546 interpreted as yes/no for each. Other methods are available, but require additional work to make
547 them usable for analyses.

548 We advise that **free text** variables be removed (see also section 8). If the retention of free text
549 is necessary, no special treatment beyond those measures outlined in section 8 is necessary.

550 Some database systems incorporate into the dataset system-level variables such as row numbers
551 in all tables of a data export. Such variables are often of no use and can typically be removed,
552 but this should be confirmed on a case-by-case basis. **Missing values** should be reported as
553 “NA” and clearly distinguished from non-missing categorical answers like “unknown”.

554 **10.2. Variable labels**

555 Variable descriptions are equally important. Without a meaningful name, it is difficult to guess
556 what a particular variable refers to. Short names are preferable for statistical programming and
557 database purposes (some software even imposes limits on the length of names), but this can
558 obscure the meaning of a given variable. Thus, besides the codebook for the meaning of values
559 of (categorical) variables, another file with the labels for each variable is required; for
560 consistency, we call it a labelbook. The labelbook should contain the variable name as it exists
561 in the data (e.g. mi) together with its description (myocardial infarction), any restrictions or
562 dependencies (only if mi == Yes), whether or not the variable is optional, and perhaps some
563 useful notes even if they might also be in other documentation such as the study protocol or
564 data management plan. The level of detail provided in the description depends on context and
565 is likely to evolve over time. We also suggest providing relevant links to the study protocol, for
566 example highlighting endpoints such as "Primary endpoint as described on page XX of the
567 study protocol". A column indicating the data type of each variable is also essential. Different
568 databases use different terms for each type, so a more standardized set of terms is provided in
569 the Appendix (Table 2, p. 47).

570 Of note, we do not list calculated fields here because calculated values returned from electronic
571 data capture systems usually are being re-calculated using statistical software.

572 The appendix provides an example of a labelbook with information on the form/table where the
573 variable is collected/stored, variable name, description/label, data type, unit, applicable value
574 label name, and whether the variable is collected as stored or whether values are
575 calculated/derived (Table 3, p. 48).

576 We would also recommend having a fully annotated version of the (electronic) case report form
577 with example data. Annotations should include variable names, option values, and any logic
578 which defines when a variable should be entered or when a variable/question is shown or hidden
579 in the electronic case report form.

580 As mentioned previously, system variables can typically be removed as they often include
581 potentially identifying information (at least for the study team). The golden rule, though, is that
582 every variable that exists in the data should be described in the labelbook.

583 **10.3. Time structures in the data sampling**

584 If there is a time structure to the data such as multiple follow ups, it is mandatory to include a
585 visit identifier in the data set which allows the discrimination of the visits for a participant. This
586 is particularly important when an individual form is used multiple times. In principle, this can
587 be done by using a key variable containing the visit identifier (long format data) or by a naming
588 convention such as adding a number at the end of the variable name (stub) indicating the order
589 (wide format data). To reduce empty cells, it is advantageous to separate data by form and we
590 recommend providing data in long format although this must be assessed on a case-by-case
591 basis. The appendix provides an example by looking at fictitious eligibility and blood label
592 values forms (Table 4 and

593 Table 5, p. 48).

594 Forms that do not fit into the normal visit structure (sometimes called unscheduled visits or log
595 forms to record medication or events) can be supplied with a "position" variable to indicate the
596 repetition number of the form (starting at either 0 or 1). The visit structure, definition of
597 unscheduled visits and the starting indices should be reported in the data documentation (see
598 section 11). In Table 6 (p. 50) we see that participant 1 reported taking a medication at two time
599 points, while participant 4 reported taking morphine for a period of time, including changing
600 doses. The remaining participants took no medications.

601 Another type of necessary information is information about which variables belong to which
602 form, which can be captured in the labelbook, and which forms are collected during which
603 visits. Following our previous logic, we call this a visitbook (Table 7, p. 50). It requires a
604 column for the visit identifier, and a column for which forms occur in each visit. Each row
605 indicates a visit-form combination (i.e. a visit could have multiple forms, and a form could be
606 in multiple visits). An additional column with the name of the visit is also useful. There
607 should also be a graphical representation of the visit structure as shown in of the Appendix
608 (Table 8, p. 50).

Box 5: Recommendation concerning variables within a shared dataset

R14. Prepare data in a long format, with appropriate keys to link tables together.

R15. Document all variables in all tables, and the tables themselves.

609

610 **11. Metadata and documentation**

611 A data file alone is of limited use, so the concept of data sharing needs comprehensive
612 documentation to go with the data (also see sections 9 and 10). This documentation serves two
613 purposes. It should, on one hand, enable someone not involved with the study to understand
614 and use the data appropriately. On the other hand, it should allow someone with a scientific
615 question in the field of the study to find the data. This section gives a definition of the term
616 metadata and what we think this metadata should contain, at a minimum, to go along with
617 shared clinical research datasets. We also give recommendations with respect to the information
618 and documentation that should be provided. This section therefore deals with findability and
619 re-usability of the data, as claimed in the FAIR principles.

620 **11.1. Metadata schemes**

621 Metadata are data about data, typically structured information such as numbers or classification
622 options, that describe a fixed set of aspects of a data object in a human and, importantly,
623 machine-readable way (21). This definition is in accordance with the concept “metadata scheme”
624 as used in libraries and repositories to denote the fields that describe the stored objects (21,22).

625 The main purpose of metadata is to find and describe a data object such as a data file, a
626 document, or a whole shared package containing different types of artifacts. Because
627 standardized metadata also allows for interoperability between systems, a data object can be
628 made visible from other points of access (<https://www.openarchives.org/pmh/>) as far as the
629 involved metadata schemes cover the same aspects. For example, it might be findable via
630 repositories, databases, or registries.

631 Canham et al. (23) suggest the use of a minimal extension of the DataCite metadata scheme for
632 clinical research data (24) which is a general purpose scheme. Study details can be found
633 basically in one field ("A.3 Study topics"), and the description of the dataset hence remains
634 somehow vague. We think that it is preferable to use a metadata scheme that supports more
635 specific searches. We expect independent reuse of data to evolve into an established scientific
636 research method also in clinical research, so we recommend a metadata scheme that allows
637 researchers to a large extent decide whether or not data are relevant for their research purpose.
638 The World Health Organization (WHO) set out requirements to describe a study (25) while the
639 International Committee of Medical Journal Editors (ICMJE) provided guidelines (26). Section
640 17.3 in the Appendix shows the set of items required by International Standard Randomised
641 Controlled Trials Number (ISRCTN) deemed essential to describe a study which we consider
642 suitable for data sets in most respects. Provided that clinical trials should be registered in WHO
643 compliant registries, this metadata is already publicly available and might be linked to a dataset
644 in a repository via an application protocol interface (API) in the repository. If this is not
645 available, the data should be entered manually. It is important to ensure consistency across the
646 registry ^{Glossary} entry and any data repository entries. Although the scheme gives clear guidance

647 on what information must be provided, it does not mandate how. To improve findability it is
648 recommended to use controlled vocabulary as far as possible. If controlled vocabulary is used,
649 it is important to provide information to the underlying scheme that was used including the
650 version.

651 11.2. Additional documentation

652 In addition to metadata, further documentation is needed to make use of the data. As described
653 above in sections 9 and 10, codebooks, labelbooks, and visitbooks provide necessary
654 information. Someone who wants to understand the data also needs to know how it was
655 collected, which sources were used, what hierarchy there was among data sources, and the
656 definitions applied. The context and purpose of the collection is important, as well as what
657 methods were used to ensure data quality. Information that relates to the conduct of the research
658 project is also needed, such as the reason for missingness of certain data and any adaptations
659 that had to be made. If a new tool or drug is investigated, a comprehensive description/brochure
660 of it is also mandatory. Furthermore, the details of data preparation have to be provided such as
661 derivation of variables, and also the process of rounding or jittering data for de-identification
662 (see section 8) has to be described together with its impact on the result, if applied (this
663 information is typically part of the data management plan).

664 The study protocol and statistical analysis plan^{Glossary} with amendments contain a large part of
665 the information needed, but researchers have to carefully consider whether this information is
666 enough for each individual project.

Box 6: Recommendation for metadata and additional documentation

R16. We recommend selecting a repository with a metadata scheme that allows for meaningfully detailed search on clinical studies (e.g., search options “patients condition”, “intervention”, “study endpoints”, etc.).

R17. We recommend as a minimum to upload with the data:

- a. Readme file describing the data package and containing information to be shared and not contained in the other documents, ideally with a tabular summary of all files (Appendix, 错误!未找到引用源。)
- b. Change log to capture different versions of the data set
- c. Study protocol
- d. Statistical analysis plan
- e. Clinical study reports
- f. Blank consent form
- g. Fully annotated case report form (CRF)
- h. Codebook, labelbook, visitbook
- i. References to any standardized vocabularies or catalogues used
- j. Code for data preparation
- k. Description/brochure of a new tool or drug, if applicable

- l. Documentation of means undertaken for anonymization
- m. Data management plan

667 **11.3. Is statistical analysis code needed for data sharing?**

668 Note that we distinguish between data preparation code and analysis code, and we consider the
669 preparation code to be necessary to go with the data (as it generates an analyzable dataset from
670 the raw data). It is also possible to share the data file after preparation. We see different aspects
671 involved in the question whether sharing analysis code is essential:

- 672 • **Reproducibility:** Undoubtedly, shared code allows for the most precise and quick
673 reproduction of the results because certain analyses might be implemented differently
674 in different software packages, and analyses can be done using different commands
675 within the same software that might even have different implementations. Still, sharing
676 code will often not lead to complete reproducibility because software versions and the
677 underlying operating system might affect usability of the code.
- 678 • **Detection of errors:** Some errors in the analysis can only be detected when scrutinizing
679 the code. Statisticians agree that wrong results are often due to errors in data
680 preparation. From this point of view, sharing of raw data and data preparation code is
681 preferable to sharing data after preparation. Reproducibility of results, even though
682 desirable, does not mean correctness, but is a step in checking it.
- 683 • **Additional information:** Usually, a statistical analysis plan is available for a clinical
684 trial describing in detail all analysis steps. However, statistical code might contain
685 additional details not covered by the statistical analysis plan. Availability of statistical
686 code is therefore essential to fully understand the analyses that were done.

Box 7: Recommendation regarding availability of analysis code

R18. In general, we recommend sharing of code with the dataset and recommend that statisticians keep to programming standards in the scripts, such as:

- Write a master script file that calls all script files of the analysis in the correct sequence.
- Follow a reasonable naming convention.
- Explain each step of the program in (extensive) comments.
- Check logical rigor of the entire code.

687 12. Version control

688 Version control allows one to track changes of objects or files through time. Because it may be
689 difficult to tell whether a dataset has been used, simply replacing an object is likely to be
690 undesirable as it would render the DOI referenced by the data user void (or rather, the DOI
691 would be correct, but the dataset it referred to is no longer available or has changed). Version
692 control may not be relevant for all datasets that will be shared. For example, a dataset that
693 accompanies a publication would be unlikely to require version control as it is a static item—it
694 does not change. Similarly, if a questionnaire performed and shared in 2017 was repeated in
695 2019 but the data were shared separately (2017 data not included), no version control is
696 necessary (although it may be helpful to refer to the other dataset in the metadata). Conversely,
697 extracts from registries might need version control if new data are periodically added to the
698 dataset. Similarly, if the originally shared dataset from a clinical trial is shared but only some
699 variables are cleaned and a second dataset is shared with all variables cleaned, this would ideally
700 be a revision. New data (variables or observation) or changes to data are reasons to make a new
701 version. Replacing only parts in the data object is easier than creating a whole new data object.

702 Where version control is considered necessary, a new DOI should be assigned to the object.
703 Ideally the new objects DOI would indicate that it is a child of the original object. For example,
704 dataset X is assigned a DOI of 1234. A year later, new data are added to X and the dataset is
705 shared. A DOI of 1234.1 would indicate that it is a child of the original dataset (the main part
706 of the DOI has remained the same, but has an extra part appended). If this is not possible and
707 the new dataset is assigned a completely different DOI (e.g. 5678), then the original DOI should
708 be entered into the metadata of the new dataset, and vice versa, to establish a link between the
709 objects.

Box 9: Recommendations regarding version control

- R19. Objects whose content has changed - new data appended to the original dataset (variables or observations) should be versioned.
- R20. A related DOI should be assigned to the new dataset, rather than creating a whole new object. At the minimum, the DOI of the different versions should be stored in the metadata of all objects.

710

711 13. Selection of repository

Infobox 2: Data repository versus (clinical trial) registry

Registries: A clinical trial registry is a collection of records about clinical trials according to an agreed upon set of metadata (27). In registries accepted by the World Health Organization (WHO) and included in their International Clinical Trials Registry Platform (ICTRP), see 9.1, these records contain a minimum amount of information as defined in the WHO Data Set (25). As of 2019, this data set does not define or require attached artifacts or files. Confusingly, the WHO calls the database behind its Search Portal "Central Repository" (27), when it is in fact a registry.

Data repositories: In contrast, a data repository is a (digital) collection of digital datasets. Although not mandatory, the term nowadays implies a function to make these datasets findable, accessible, and reusable (5) and allows for longer term storage. Technically, a repository consists at least of a backend, a database to store metadata and information, and file server to store the datasets and other digital artifacts, and a web-based frontend that allows users to access the backend.

712 13.1. Principles

713 According to the FAIR data principles, research data should be findable, accessible,
714 interoperable, and reusable (4,5), see section 2. Principle F3 mandates that "(meta)data are
715 registered or indexed in a searchable resource" (4). Although the principles do not explicitly
716 mention data repositories, principle F3 implies that research data should be stored in an
717 appropriate repository that follows all principles (5). The European Clinical Research
718 Infrastructure Network (ECRIN) data sharing statement is more explicit and states, that "data
719 and trial documents made available for sharing should be transferred to a suitable data
720 repository" (3) and we support this view. According to the FAIR data principles, research data
721 should be findable, accessible, interoperable, and reusable (4,5), see introduction. Principle F3
722 mandates that "(meta)data are registered or indexed in a searchable resource" (4). Although the
723 principles do not explicitly mention data repositories, principle F3 implies that research data
724 should be stored in an appropriate repository that follows all principles (5).

725 When selecting a repository, clinical researchers therefore should ensure that the repository
726 respects all FAIR data principles as a minimum. Although there are alternative initiatives like
727 CoreTrustSeal (28), the FAIR principles seem to be the most widely accepted. However, other
728 initiatives might evolve over time and become generally agreed standards. Given the lack of
729 generally agreed standards and certification processes, researchers will need to assess the
730 suitability of a repository for their purposes.

731 **13.2. Time point**

732 Ideally, the appropriate repository is identified before writing the Data Management Plan (see
733 section 7) and then described therein. We assume that a sponsor/investigator uses the same
734 repository for all her/his projects so this should be feasible.

735 **13.3. Identifying potential repositories**

736 So far, no repository exists that is specific for clinical research projects. Therefore, clinical
737 researchers need to identify an appropriate repository by themselves. Many institutions
738 involved in clinical research, like universities, currently maintain their own institutional
739 repository. This might be a good starting point in the evaluation process. Alternatively,
740 universities usually have a central contact point that supports researchers with issues related to
741 data sharing and open science in general (29).

742 For projects that were funded by extramural grants, there might be specific requirements for a
743 repository or even a specific repository mandated. For example, the Bill & Melinda Gates
744 Foundation maintains a list of approved repositories for publications published in Gates Open
745 Research (30). It is also expected that the planned European Open Science Cloud (EOSC) will
746 affect how data from projects funded by the European Union will be shared (31). Repository
747 registries maintain a searchable database of repositories. The largest one is probably [r3data](#), a
748 collaborative project of large European academic institutions. R3data can help locating topic
749 specific repositories, which may be a better choice than an institutional repository because data
750 are more likely to be found in a search for that particular topic. Furthermore, Swiss academic
751 research institutions are currently developing a digital repository for long-term preservation and
752 publishing of research data, Olos (32), to support the publication needs of funders and help
753 researchers to manage research data.

754 Another choice might be Zenodo, which is based at CERN (European Organisation for Nuclear
755 Research). There are also for-profit/commercial repositories such as FigShare and Dryad,
756 although we do not explicitly recommend their use.

757 **13.4. Selection criteria**

758 After having identified a set of potential repositories, a researcher will need some explicit
759 criteria to select a repository. We suggest an approach to structure this process which is based
760 on a report by the Digital Curation Centre in Edinburgh (33), shaped as a checklist (Table 10,
761 p. 54). Some items are very specific, others cannot be defined exactly and require adaptations
762 on a project basis and not all aspects might be assessable.

763 Another useful resource are the levels of digital preservation by the National Digital
764 Stewardship Alliance (34).

Box 8: Recommendations selection of repository

- R21. Select a suitable repository, and include this information in the data management plan. Institutional repositories might be a good choice.
- R22. Make data as open as possible, but as closed as necessary (FAIR)

766 14. Requesting and use of data

767 Principle 6 of the ECRIN statement (3) states: “In the context of managed access, any citizen
768 or group that has both a reasonable scientific question and the expertise to answer that question
769 should be able to request access to individual participant data and trial documents.” This begs
770 the question of who decides whether a question is reasonable and an individual/group has the
771 relevant competencies. Decisions made by the original project team could be seen as biased.
772 Accordingly, the ECRIN statement (3) suggests that ideally each repository would have
773 independent boards to assess the "scientific merit, potential impact and appropriateness of the
774 proposed secondary analyses". With slightly different priorities, such a board might also be
775 referred to as Data Access Committee^{Glossary} (DAC). A DAC might evaluate and approve data
776 requests within a reasonable response time. This would of course require separate boards or
777 DACs for different subject areas. From our point of view, it is a good idea to have a board of
778 specialists/DACs supporting new research on existing data, but it might be difficult to find the
779 resources for this work. From a legal point of view, there are few minimal requirements that
780 have to be fulfilled in order to receive data:

- 781 1. The data requester has to confirm that the purpose of the data request is scientific, that
782 the research project will be conducted in accordance with the local legislation (Human
783 Research Act, authorization from ethics committee) and rules of conduct (Good
784 Clinical Practice). Any different purpose would have to be explicitly mentioned in the
785 informed consent (see (35) and see section 6).
- 786 2. The data requester has to confirm that she/he:
 - 787 2.1. Will not try to identify individual persons in the data
 - 788 2.2. Will not give the data to other persons
 - 789 2.3. Will maintain data security
 - 790 2.4. Will report any accidental finding to the data provider

791 We think that publishing the metadata and sharing the data after checking these two
792 requirements will be the usual process in clinical research. The requesting process is obviously
793 determined by the repository, so we only sketch some possible implementation options. With
794 minimal use of resources, requirement 1 might be covered by a checkbox on the request form
795 that a requester has to tick. If she/he does not, a pop-up window might occur saying that the
796 request is going to be rejected. Requirement 2 needs the requester to be a person able to confirm
797 in a legally binding way. There are established ways to check whether an action is done by a
798 human over the Internet, but in the context of data sharing we assume by default that the
799 requester has an academic affiliation, which will be used to verify the requester's identity. A
800 requester without academic affiliation might turn to the data provider directly. The requester
801 might confirm the items of requirement 2 by signing a contract or by ticking a checkbox of a
802 license agreement (35). The agreement might contain an example text of how the original study

803 and its investigators should be acknowledged in any kind of publication to ensure that data
804 generators receive appropriate recognition (36). All requests are stored by the repository to be
805 traced by interested persons such as the principal investigator. If there is a DAC/board of
806 specialists it makes sense that a data request comprises a proposal together with an authorization
807 from the ethics committee (unless the request comes from a country without ethics committees).
808 The proposal briefly describes the aims and objectives of the planned study or reanalysis of the
809 requested data, the planned analysis, the data that are needed and the time frame of the study.
810 The DAC/board of specialists evaluates and approves the request, checks the requesters identity
811 and informs the principal investigator.

812

813 **15. References**

- 814 1. Naudet F, Sakarovitch C, Janiaud P, Cristea I, Fanelli D, Moher D, et al. Data sharing
815 and reanalysis of randomized controlled trials in leading biomedical journals with a full
816 data sharing policy: Survey of studies published in the BMJ and PLOS Medicine. *BMJ*.
817 2018;360.
- 818 2. Taichmann D. Data Sharing Statements for Clinical Trials: A Requirement of the
819 ICMJE. *Ann Intern Med* [Internet]. 2017; Available from: [http://www.icmje.org/news-](http://www.icmje.org/news-and-editorials/data_sharing_june_2017.pdf)
820 [and-editorials/data_sharing_june_2017.pdf](http://www.icmje.org/news-and-editorials/data_sharing_june_2017.pdf)
- 821 3. Ohmann C, Banzi R, Canham S, Battaglia S, Matei M, Ariyo C, et al. Sharing and
822 reuse of individual participant data from clinical trials: Principles and
823 recommendations. *BMJ Open*. 2017;
- 824 4. FORCE11. The FAIR Data Principles - FOR COMMENT [Internet]. 2016 [cited 2018
825 Dec 6]. Available from: <https://www.force11.org/group/fairgroup/fairprinciples>
- 826 5. Wilkinson MD, Dumontier M, Aalbersberg IJ, Appleton G, Axton M, Baak A, et al.
827 Comment: The FAIR Guiding Principles for scientific data management and
828 stewardship. *Sci Data*. 2016;
- 829 6. GO FARE. FAIR Principles - GO FAIR [Internet]. 2018. Available from:
830 <https://www.go-fair.org/fair-principles/>
- 831 7. The Federal Assembly of the Swiss Confederation. CC 235.1 Federal Act of 19 June
832 1992 on Data Protection (FADP). 1992; Available from:
833 <https://www.admin.ch/opc/en/classified-compilation/19920153/index.html>
- 834 8. Confederation TFA of the S. Federal Act on Research involving Human Beings.
835 <Http://WwwAdminCh/Opc/En/Classified-Compilation/20061313/IndexHtml> [Internet].
836 2011; Available from: <http://www.admin.ch/opc/en/classified->
837 [compilation/20061313/index.html](http://www.admin.ch/opc/en/classified-compilation/20061313/index.html)
- 838 9. EDI ED des I. Koordinationsstelle Forschung am Menschen (kofam) [Internet].
839 Koordinationsstelle Forschung am Menschen (kofam) c/o Federal office of public
840 health FOPH CH-3003 Bern. 2017. Available from: [http://kofam.ch/en/applications-](http://kofam.ch/en/applications-and-procedure/projects-that-do-not-require-authorisation/)
841 [and-procedure/projects-that-do-not-require-authorisation/](http://kofam.ch/en/applications-and-procedure/projects-that-do-not-require-authorisation/)
- 842 10. Network SPH. Swiss Personalized Health Network [Internet]. Swiss Personalized
843 Health Network. 2019. p. <https://www.sphn.ch/en.html>. Available from:
844 <https://www.sphn.ch/%0A>
- 845 11. Swiss Personalised Health Network. Ethical Framework for Responsible Data
846 Processing in Personalized Health Research [Internet]. 2018. Available from:
847 [https://www.sphn.ch/dam/jcr:6fb78ffa-95c8-4372-bfb1-](https://www.sphn.ch/dam/jcr:6fb78ffa-95c8-4372-bfb1-5c9b1e2cb53d/Ethical_Framework_20180507_SPHN.pdf)
848 [5c9b1e2cb53d/Ethical_Framework_20180507_SPHN.pdf](https://www.sphn.ch/dam/jcr:6fb78ffa-95c8-4372-bfb1-5c9b1e2cb53d/Ethical_Framework_20180507_SPHN.pdf)
- 849 12. Husedzinovic A, Ose D, Schickhardt C, Fröhling S, Winkler EC. Stakeholders'
850 perspectives on biobank-based genomic research: Systematic review of the literature.
851 *Eur J Hum Genet*. 2015;23(12):1607–14.
- 852 13. Trial Master File Reference Model [Internet]. <https://tmfrefmodel.com/>. 2014.
853 Available from: <http://tmfrefmodel.com/>
- 854 14. Chevrier R, Foufi V, Gaudet-Blavignac C, Robert A, Lovis C. Use and understanding

- 855 of anonymization and de-identification in the biomedical literature: Scoping review.
856 Vol. 21, Journal of Medical Internet Research. Journal of Medical Internet Research;
857 2019.
- 858 15. O'Neill L, Dexter F, Zhang N. The risks to patient privacy from publishing data from
859 clinical anesthesia studies. *Anesth Analg*. 2016 Jun 1;122(6):2017–27.
- 860 16. Services C for M& M. The health insurance portability and accountability act of 1996
861 (HIPAA). Online <http://www.C.hhs.gov/hipaa> [Internet]. 1996; Available from:
862 [https://www.federalregister.gov/documents/2009/10/30/E9-26203/hipaa-](https://www.federalregister.gov/documents/2009/10/30/E9-26203/hipaa-administrative-simplification-enforcement)
863 [administrative-simplification-enforcement](https://www.federalregister.gov/documents/2009/10/30/E9-26203/hipaa-administrative-simplification-enforcement).
- 864 17. Statistik B für. MEDSTAT.
865 [https://www.bfs.admin.ch/bfs/de/home/statistiken/gesundheit/nomenklaturen/medsreg.](https://www.bfs.admin.ch/bfs/de/home/statistiken/gesundheit/nomenklaturen/medsreg.html)
866 [html](https://www.bfs.admin.ch/bfs/de/home/statistiken/gesundheit/nomenklaturen/medsreg.html).
- 867 18. Garfinkel. De-identification of personal information. NISTIR 8035. 2015;
- 868 19. Atreya. Reducing patient re-identification risk for laboratory results within research
869 datasets. *J Am Med Inf Assoc*. 2013;20(95).
- 870 20. Pike R. UTF-8 [Internet]. 2019. Available from: <https://en.wikipedia.org/wiki/UTF-8>
- 871 21. Oulun Yliopisto. Research Data Guide [Internet]. Managing research data: Data
872 documentation and metadata. 2018. Available from:
873 http://libguides.oulu.fi/Researchdata/Data_documentation
- 874 22. University ta trobe. Durable file formats. 2019; Available from:
875 <https://latrobe.libguides.com/dataorganisation/fileformats>
- 876 23. Canham S, Ohmann C. A metadata schema for data objects in clinical research. *Trials*
877 [Internet]. 2016 Dec 24;17(1):557. Available from:
878 <http://trialsjournal.biomedcentral.com/articles/10.1186/s13063-016-1686-5>
- 879 24. DataCite. DataCite Metadata Schema for the Publication and Citation of Research
880 Data. 2011;29.
- 881 25. WHO | WHO Data Set. WHO [Internet]. 2018 [cited 2018 Dec 6]; Available from:
882 <http://www.who.int/ictrp/network/trds/en/>
- 883 26. ICMJE. Up-Dated ICMJE Recommendations [Internet]. 2019. Available from:
884 <http://www.icmje.org/icmje-recommendations.pdf>
- 885 27. WHO | Glossary [Internet]. Who. Geneva, Switzerland: World Health Organization;
886 2018 [cited 2018 Dec 6]. Available from: <https://www.who.int/ictrp/glossary/en/>
- 887 28. CoreTrustSeal – Core Trustworthy Data Repositories [Internet]. [cited 2018 Dec 6].
888 Available from: <https://www.coretrustseal.org/>
- 889 29. Sciences swiss academies of arts and. Swiss Academies Factsheets. 2019. p.
890 <http://www.akademien-schweiz.ch/index/Publikatione>.
- 891 30. Bill & Melinda Gates Foundation. Gates Open Research - Data Guidelines [Internet].
892 Seattle, WA: Bill & Melinda Gates Foundation; [cited 2018 Dec 6]. Available from:
893 <https://gatesopenresearch.org/for-authors/data-guidelines>
- 894 31. European Commission. European Open Science Cloud (EOSC) | Open Science -
895 Research and Innovation - European Commission [Internet]. [cited 2018 Dec 6].
896 Available from: [https://ec.europa.eu/research/openscience/index.cfm?pg=open-science-](https://ec.europa.eu/research/openscience/index.cfm?pg=open-science-cloud)
897 [cloud](https://ec.europa.eu/research/openscience/index.cfm?pg=open-science-cloud)

- 898 32. Swissuniversities. Olos [Internet]. 2018. Available from: <https://www.dlcm.ch/olos>
- 899 33. Whyte A. Where to keep research data: DCC checklist for evaluating data repositories
900 (v.1.1) [Internet]. Edinburgh: Digital Curation Centre; 2015 [cited 2018 Dec 10].
901 Available from: www.dcc.ac.uk/resources/how-guides
- 902 34. Preservation D, Questions FA, Based L, Format O. Levels of Digital Preservation
903 Support [Internet]. 2014 [cited 2018 Dec 10]. p. 1–5. Available from:
904 <https://ndsa.org/activities/levels-of-digital-preservation/>
- 905 35. IWATSUBO T. アルツハイマー病根本治療薬の臨床応用に向けて : Alzheimer’s
906 Disease Neuroimaging Initiative (ADNI)の取り組み(最前線,アルツハイマー病).
907 Vol. 43, Farumashia. 2007. p. 894–8.
- 908 36. Alter G, Gonzalez R. Responsible practices for data sharing. *Am Psychol*. 2018 Feb
909 1;73(2):146–56.
- 910 37. Kirkwood BR, Sterne JAC. *Essential Medical Statistics*. Medical statistics. 2003. 513
911 p.
- 912
913
914

915 **16. Glossary**

| Term | Definition | Reference |
|----------------------------------|---|--|
| Anonymization | <p>Process by which any way of linking data in a data set with a natural person is irreversibly removed/destroyed or only possible with disproportionate effort. <i>De-identification</i> or <i>Pseudonymization</i> with destruction of the <i>Key</i> are needed as a minimum for this process. It must be noted that the required measures for anonymization must be defined on a case-by-case basis because a combination of not directly identifying information might enable identification of a natural person. The Human Research Act acknowledges that absolute irreversible anonymization is impossible. Disproportionate effort is given if linking:</p> <ol style="list-style-type: none"> 1. Is only possible with considerable criminal energy, or 2. Requires extensive technical infrastructure and know-how. | <p>Schweizerischer Bundesrat p. 8096. Eidgenössisches Department des Inneren p. 69-70.</p> |
| Anonymized (health-related) data | (Health-related) Data which cannot (without disproportionate effort) be traced to a specific person. See also <i>Anonymization</i> | Human Research Act Art. 3i. and General Data Protection Regulation (EU) |
| Artifact | <p>The term artifact is used because relevant study information might be recorded in a variety of different ways, including records, documents and data. An artifact is therefore any information that is captured during a clinical trial that meets the purpose or definition described in the protocol. In some cases, the artifact is a single document, data set or piece of information but in other cases it could be represented by multiple document types or data types.</p> | <p>https://tmfrefmodel.com/wp-content/uploads/2018/03/tmf-rm-deliverable-user-guide-v1-2018-03-16.pdf</p> |

| | | |
|------------------|---|--|
| Case Report Form | <p>(1) A printed, optical or electronic document designed to record all of the protocol-required information to be reported to the sponsor for each subject/patient in a clinical trial.</p> <p>(2) A record of clinical study observations and other information that must be completed for each subject in a clinical trial, per study protocol mandate. CRF can refer to either a CRF page (which contains one or more data items linked together for collection and display) or a casebook (which includes all CRF pages on which a set of clinical study observations and other information can be or have been collected, or the information collected by completion of such CRF pages for a subject/patient in a clinical study).</p> | The Free Dictionary |
| Coded data (set) | <p>De-identified data that can be linked to a specific person via a <i>Key</i> (code). This means that the data are anonymized for any person who accesses the data and who has no direct access to the <i>Key</i>. However, the conditions under which the <i>Key</i> is stored and can be accessed are critical for qualifying data as coded:</p> <ol style="list-style-type: none"> 1. Storage of the <i>Key</i> must be separate from the storage of the data. No person directly involved in a research project or who works as a subordinate to someone who wants to use coded data may have access to the <i>Key</i>. This includes but is not limited to investigators, study nurses/coordinators, statisticians, and data managers. Precautions must be taken to ensure that only authorized persons have access to the <i>Key</i> (see 2) and each access must be documented (date and who accessed it for what reason). | Human Research Act Art. 3h. HRO Art. 26-27. |

| | | |
|-----------------------|--|--|
| | <p>2. Decoding i.e. identifying a person is only allowed under the following conditions:</p> <ol style="list-style-type: none"> a. Breaking the code is necessary to avert an immediate risk to the health of the person concerned. b. A legal basis exists for breaking the code. <p>Breaking the code is necessary to guarantee the rights of the person concerned, and in particular the right to revoke consent.</p> | |
| Controlled access | Refers to the way a data set is shared. In a controlled access model, the data are only shared with an entity if they meet certain conditions and on request. | Keerie C et al. 2018. |
| Data | Pieces of information. Within this document, we use a narrow definition of data, denoting the content of structured data files. | |
| Data Access Committee | A Data Access Committee (DAC) is a body of one or more individuals who are responsible for data release to external requestors based on consent and/or National Research Ethics terms. A DAC is typically formed from the same organization that collected the samples and generated any associated analyses. Multiple datasets may be affiliated to a single DAC. | European Genome-phenome Archive |
| Data Management Plan | Document that outlines how data are to be handled both during and after a research project including data preservation. | Wikipedia |
| Data Object | An entity available in electronic format (document, text, program, zip file). In the setting of clinical research data sharing: data and associated documents related to a clinical trial and typically stored in a repository. | Canham and Ohmann. <i>Trials</i> (2016) 17:557 |

| | | |
|---------------------------------|--|---------------------------------|
| Data Sharing/Transfer Agreement | Contract or license that describes the conditions | |
| Data Validation Plan | Document that describes the process of data validation, e.g. which variables have to be checked and what consistency rules have to be met. It might include checks on chronological sequence, completeness, identification of duplicates, checks of range and distribution shape of variables. | |
| De-identified | See <i>De-identification</i> | |
| De-identification | Process by which all directly identifying data is either removed, altered or censored from a data set. It must be noted that the term <i>de-identification</i> as such has no legal basis in Switzerland but rather is a concept originating in the USA based on rules set forth in the Health Insurance Portability and Accountability Act (HIPAA). For the purpose of this document, de-identification relates only to directly identifying data. De-identification is usually insufficient for data sharing. | US Office for Civil Rights 2012 |
| External party | Receiver of de-identified data whose access to the data was not explicitly consented to by the patients (could be a researcher or data repository). Alternative phrase: third party. | |
| Identifier | A number or string that identifies/labels a unique object. <i>Identifiers</i> in a clinical study project usually follow an encoding system; in other words, there are rules behind the generation of the <i>identifier</i> . Such rules might be a pseudonymization algorithm (see <i>pseudonym</i>) or a sequential numbering system. Identifiers are therefore often referred to as <i>ID code</i> , <i>ID number</i> , <i>record ID</i> , or <i>unique identifier</i> (UID) in the clinical research context. | Wikipedia |

| | | |
|---|---|--|
| Identifying data (directly or indirectly) | Any information that solely (directly) or jointly with other data enables identification of a natural person among a data set. | |
| Key | A piece of information that allows decrypting encrypted data. In the clinical research context this is usually a participant/patient log/list that allows linking a (unique) <i>identifier</i> (record) with the <i>identifying data</i> usually the full name, birth date, and hospital/practice identification number. The <i>key</i> is usually stored on site under restricted access (e.g. in the study-binder). | Wikipedia |
| Limited data (set) | A data set that has been de-identified and which contains only the absolute minimum number of <i>variables</i> required to conduct an analysis by an <i>External party</i> . It includes <i>variables</i> needed to derive variables which are needed to conduct the analyses by the <i>External party</i> unless these <i>variables</i> increase the risk for identification. | |
| Metadata | Data about data; a vector of structured information, typically numbers or classification options that describes a fixed set of aspects of a data object in a human and machine-readable way. | |
| Open access | Refers to the way a data set is shared. In an open access model, the data is shared publicly and can be accessed without restriction or request. | Keerie C et al. 2018. |
| Personal Data (health-related) | Any information relating to an identified/specific or identifiable natural person (<i>data subject</i>); an identifiable person is one who can be identified, directly or indirectly, in particular by reference to an identification number or to one or more factors specific to his physical, physiological, mental, economic, cultural or social identity. | Cantonal Data Protection Act (Kantonales Datenschutzgesetz, KDSG) Art.2 Par. 1 (Federal Act on Data Protection Art. 3a). |

EU Directive
95/46/EC 4.

Pseudonym A pseudonym or alias is a unique name (or more generally, a string consisting of alphabetic and potentially numeric characters) used to conceal *identifying data*. The pseudonym is generated using a set of rules (pseudonymization algorithm). A pseudonym can be generated with or without the possibility of restoring the underlying *identifying data* (reversible or irreversible pseudonymization). If the same algorithm is used across systems, pseudonymization allows for data to be linked to the same person across multiple data records or information systems without revealing the identity of the person. It must be noted that the term does not appear in any of the following laws: HRA, ClinO, HRO. Derivation of a new variable from other variable(s) using simple rules like calculating age from date of birth and enrolment date is not considered pseudonymization as this does not generate a unique attribute.

Pseudonymization (reversible or irreversible) See *Pseudonym*

Registry A clinical trial registry is an entity that houses clinical trial registers i.e. a record containing information about a clinical trial (27). In registries accepted by the World Health Organization (WHO) and included in their International Clinical Trials Registry Platform (ICTRP), these records contain a minimum amount of information as defined in the WHO Data Set (25). As of 2019, this data set does not define or require attached artifacts or files. Confusingly, the WHO calls the database behind its Search Portal "Central

| | |
|---------------------------|--|
| | Repository" (27), when it is in fact a registry. |
| Repository | Collection of digital datasets. Technically, it consists at least of a backend, a database to store metadata and information; a file server to store the datasets and other digital artifacts; and a web-based frontend that allows users to access the backend. Although not mandatory, the term implies that there is a function to make these datasets findable, accessible, and reusable (5) and allows for longer term storage. |
| Statistical Analysis Plan | A statistical analysis plan is a document that contains a more technical and detailed elaboration of the principal features of the analysis described in the protocol, and includes detailed procedures for executing the statistical analysis of the primary and secondary variables and other data. |
| Third party | See <i>External party</i> |
| Variable | A measured or recorded attribute that (37) characterizes an object, e.g. a participant. A variable is the operationalized way in which the attribute is represented for data processing i.e. a variable contains attributes. There are different types of variables (data types). The most common ones are: nominal/categorical with the special case of binary (only two categories), ordinal, numeric/continuous, date & time, string. |

917 **17. Appendix**

918 **17.1. Further detailed specification of required data** 919 **processing steps**

920 **17.1.1 Example data to be considered for deletion**

- 921 • Names, address, etc. have to be deleted, see section 8.2.
- 922 • All freetext variables should be deleted unless the content is checked and redacted
923 where necessary to ensure privacy.
- 924 • Any internal record identifier of the clinical database.
- 925 • Any identification numbers that are not needed for analysis purposes such as
926 biosample/kit numbers etc.
- 927 • Any variables that contain data that is particular or has low prevalence e.g. multiples
928 (twins, ...), special comorbidities.

929 **17.1.2 Examples and details on manipulations to decrease precision**

- 930 • Dates (time): The enrolment date (time) should be set to zero. All other date variables
931 (including date of birth) should be replaced by variables containing time relative to the
932 enrolment date using the appropriate scale e.g. years for age, or days for study visits
933 (relative study day). Consider to deliver age bands (e.g. 5 year bands) instead if the
934 disease or population is infrequent or rare. To protect persons in rare age groups, those
935 above 89 should be grouped together in a “90 or older” category. Use accordingly for
936 young ages as appropriate.
937 Alternatively, a random offset can be added to all dates in the data for a specific person.
938 It is recommended to use different offsets for each person, as long as relative
939 differences between persons are not relevant. For some dates, e.g. birthdays, or when
940 seasonal effects are of interest, other methods such as the generalization into certain
941 categories like month or years, may be required.
- 942 • Geographic information: Consider whether aggregation to MEDSTAT(17) or other
943 higher level unit is appropriate.
- 944 • Unusual data: If a variable contains data that allows identification of individuals
945 because it is special or has low prevalence consider grouping or aggregation into
946 categories.
- 947 • Height and body weight: Consider whether Body Mass Index (BMI) is sufficient and
948 derive BMI and delete height and body weight.
- 949 • Renal function: Consider whether Serum Creatinine can be replaced by estimated
950 Glomerular Filtration Rate (GFR).
- 951 • (High precision) continuous/numerical data: Round data to the next higher digit or
952 introduce random jitter on the last digit (Perturbation).

-
- 953 • Identification numbers that are needed for analysis, participant ID, study site ID
954 (cluster ID, country ID, etc.):
- 955 • All identification numbers must be replaced by a unique random number. It is
956 important to ensure that records with the same identification number, e.g. participant
957 or study site identifier, are assigned the same new random number. The general process
958 is:
- 959 1. Check all data files for the variable (identification number) of interest.
 - 960 2. Collect the maximum amount of data i.e. make sure that you get all identification
961 numbers of interest across all data files and save in a separate data file.
 - 962 3. Randomly shuffle the IDs (1. generate a new variable with random numbers (no
963 seed⁴), 2. sort data accordingly, and 3. replace the new variable with integers in
964 ascending order (new ID). Make sure that the new variable contains only unique
965 numbers).
 - 966 4. Merge the new ID into all relevant data files.
 - 967 5. Delete the original ID from all relevant data files.
 - 968 6. Repeat for other identification numbers.
- 969 If the number of records is unique for a particular identification number e.g.
970 study site ID, consider to aggregate.

971 General approaches:

- 972 • Aggregation (generalization) might be a strategy to achieve de-identification and
973 should be considered if other manipulations remain unsatisfactory. For example,
974 numerical data can be transformed into categorical variables and categorical variables
975 may be combined into new (less informative) categories. As outliers have a larger risk
976 of re-identification, one could aggregate outliers only and leave non-outlier values
977 unchanged.
- 978 • Replacing the observed value of specific record with "missing", thereby increasing the
979 frequency of certain rare combination (suppression).
- 980 • Data swapping: For a fraction of records, values of quasi-identifiers might be
981 exchanged, with the possibility of adding constraints on which pairs of records can be
982 swapped. For example, given two "similar" records, one may swap the values of one
983 quasi-identifier, e.g. age.
- 984 • Resampling: One identifies the probability distribution of the quasi-identifying data
985 and replaces its values with a random sample from its distribution. Care must be taken
986 if correlations with other variables need to be preserved.
- 987 • Subsampling: Only a subsample of the data might be shared, thereby reducing the risk
988 of re-identification.

⁴ Alternatively, a random seed might be used, but removed from any documentation after the final dataset was created and underwent the anonymization process.

989 **17.2. Further details on coding of variables**

990 **17.2.1 Formatting of date and time variables**

- 991 • Date variables should be provided in the ISO 8601 standard of year-month-day (e.g.
992 12th October 2018 would be 2018-10-12).
- 993 • Time with seconds should be coded as hours:minutes:seconds (e.g. 07:59:45 or
994 15:32:01). Where seconds are unavailable, leaving away seconds is considered
995 acceptable (e.g. 15:32), so long as all observations are coded consistently (same
996 applies to minutes).
 - 997 – Where data come from multiple time zones, the offset from Coordinated
998 Universal Time (UTC) should be added (e.g. 15:32+01:00 for Central European
999 Time). Conversion to UTC is encouraged.
- 1000 • Date-time variables should follow the rules for both date and time, and have the date
1001 part followed by the time part, separated by a space (e.g. 2018-10-12 07:59:45 or 2018-
1002 10-12 07:59; the strict ISO 8601 standard separates dates and times by T, but the space
1003 is readily recognized as a date-time variable by statistical software).
 - 1004 – As with times, the offset from UTC is vital for datasets including multiple time
1005 zones.

1006

1007 **17.2.2 Examples for further documentation of the dataset**1008 **Table 1:** Codebook example

| Labelname | Code | Value label |
|-----------|------|---------------|
| yn | 0 | No |
| yn | 1 | Yes |
| sex | 1 | Male |
| sex | 2 | Female |
| route | 1 | Oral |
| route | 2 | IV |
| route | 3 | Anal |
| unit | 1 | mg/dL |
| unit | 2 | mg |
| unit | 3 | ug/dL |
| unit | 4 | ug |
| unit | 5 | g |
| freq | 0 | less frequent |
| freq | 1 | daily |
| freq | 2 | twice daily |
| freq | 3 | every 8 hours |
| freq | 4 | every 6 hours |
| freq | 5 | more frequent |

1009 **Table 2:** Recommended data type names. These types would be referenced in the labelbook

| Data type | Description |
|-----------|---|
| Str | Free text (short for string). See above for notes |
| Int | Integer |
| Num | Numbers without specific accuracy |
| Num_Xdp | Number with X decimal places (e.g. num_1dp for values with 1 decimal place) |
| Date | Date variables (formatted to ISO 8601 standards)* |
| Time | Time variables (formatted to ISO 8601 standards)* |
| Datetime | Date and time variable (formatted to ISO 8601 standards)* |
| Cat | Categorical variable (e.g. male/female/undifferentiated/unknown) |
| Bin | Binary variables (e.g. yes/no) |

1010 * would ideally be converted to study time (e.g. days since randomization/informed consent/some other
1011 reference point); see section 6.

1012

1013

1014 **Table 3:** Labelbook example

| Form | Variable | Label | Type | Unit | Label name | Note |
|------|------------|---|---------|-------|------------|-----------------------------|
| | visit | Visit ID | Int | | | |
| | pid | Participant ID | Int | | | |
| | position | Position in repeating form sequence | Int | | | |
| elig | sex | Sex | Cat | | sex | |
| elig | age | Age | Int | Years | | |
| elig | ic | Informed Consent given | Cat | | yn | |
| elig | ic1 | Age 18 years or older | Cat | | yn | |
| elig | ic2 | Recurrent kidney stone disease | Cat | | yn | |
| elig | ex1 | More than 5 instances of kidney stone disease | Cat | | yn | |
| elig | ic_date | Date of Informed Consent | Date | | | |
| lab | lab_bl_yn | Blood sample taken | Cat | | yn | |
| lab | lab_bl_rbc | Red blood cell count | num_1dp | mcl | | |
| lab | lab_bl_ldl | Blood LDL cholesterol | Int | mg/dl | | |
| drug | uvisit | Unscheduled visit ID | Int | | | |
| drug | position | Drug name | Str | | | |
| drug | route | Administration route | Cat | | route | |
| drug | Dose | Dose | Num | | | see unit for relevant units |
| drug | Unit | Unit | Cat | | unit | |
| drug | Freq | Frequency | Cat | | freq | |
| drug | freq_det | Frequency details | Str | | | if freq = 0 or 5 |
| drug | Start | Start | Date | | | |
| drug | ongoing | Ongoing? | Cat | | yn | |
| drug | End | End | Date | | | |

1015 **Table 4:** Structure of dataset with one row per participant (part of eligibility form)

| Visit* | pid | sex | age | ic1 | ic2 | ex1 | ic | ic_date |
|--------|-----|-----|-----|-----|-----|-----|----|------------|
| 1 | 1 | 1 | 58 | 1 | 1 | 0 | 1 | 2016-01-09 |
| 1 | 2 | 2 | 54 | 1 | 1 | 0 | 1 | 2016-01-15 |
| 1 | 3 | 1 | 54 | 1 | 1 | 0 | 1 | 2016-07-11 |
| 1 | 4 | 1 | 41 | 1 | 1 | 0 | 1 | 2016-09-01 |
| 1 | 5 | 1 | 32 | 1 | 1 | 0 | 1 | 2017-09-11 |
| 1 | 6 | 2 | 36 | 1 | 1 | 0 | 1 | 2017-09-28 |
| 1 | 7 | 2 | 30 | 1 | 1 | 0 | 1 | 2017-10-24 |
| 1 | 8 | 2 | 51 | 1 | 1 | 0 | 1 | 2018-10-27 |

1016 * The visit variable in this case is optional as the eligibility form is only used once.

1017

1018 **Table 5:** Structure of dataset with multiple rows per participant (part of blood laboratory
1019 values form)

| visit | pid | lab_bl_yn | lab_date | lab_bl_rbc | lab_bl_chol |
|-------|-----|-----------|------------|------------|-------------|
| 1 | 1 | 1 | 2016-01-09 | 5.1 | 123 |
| 1 | 2 | 1 | 2016-01-15 | 5.6 | 144 |
| 1 | 3 | 1 | 2016-07-11 | 4.7 | 103 |
| 1 | 4 | 0 | | | |
| 1 | 5 | 0 | | | |
| 1 | 6 | 0 | | | |
| 1 | 7 | 1 | 2017-10-24 | 5.2 | 110 |
| 1 | 8 | 1 | 2018-10-27 | 4.2 | 90 |
| 2 | 1 | 0 | | | |
| 2 | 2 | 0 | | | |
| 2 | 3 | 1 | 2016-08-05 | 4.8 | 66 |
| 2 | 4 | 1 | 2016-10-02 | 4.5 | 142 |
| 2 | 5 | 1 | 2017-10-12 | 4.7 | 103 |
| 2 | 7 | 0 | | | |
| 2 | 8 | 1 | 2018-11-25 | 6.1 | 125 |
| 3 | 1 | 1 | 2016-03-10 | 5.5 | 140 |
| 3 | 2 | 1 | 2016-03-20 | 5.4 | 130 |
| 3 | 3 | 0 | | | |
| 3 | 4 | 1 | 2016-11-06 | 6 | 129 |
| 3 | 5 | 0 | | | |
| 3 | 7 | 1 | 2017-12-20 | 5.2 | 111 |
| 3 | 8 | 1 | 2018-12-28 | 4.5 | 121 |

1020

1021

1022 **Table 6:** Example for a log form data table

| pid | position | drug | route | dose | unit | freq | freq_det | start | ongoing | end |
|-----|----------|-------------|-------|------|------|------|---------------|------------|---------|------------|
| 1 | 0 | amoxicillin | 1 | 500 | 2 | 2 | | 2016-02-25 | 0 | 2016-03-05 |
| 1 | 1 | amoxicillin | 1 | 500 | 2 | 2 | | 2018-10-20 | 0 | 2018-11-01 |
| 2 | 0 | | | | | | | | | |
| 3 | 0 | | | | | | | | | |
| 4 | 0 | morphine | 1 | 15 | 2 | 5 | every 4 hours | 2017-01-20 | 0 | 2017-01-25 |
| 4 | 1 | morphine | 1 | 30 | 2 | 5 | every 4 hours | 2017-01-26 | 0 | 2017-01-30 |
| 4 | 2 | morphine | 1 | 5 | 2 | 1 | | 2017-01-31 | 0 | 2017-02-05 |
| 5 | 0 | | | | | | | | | |
| 6 | 0 | | | | | | | | | |
| 7 | 0 | | | | | | | | | |
| 8 | 0 | | | | | | | | | |

1023 We see that participant 1 reported taking a medication at two time points, while participant 4 reported
1024 taking morphine for a period of time, including changing doses. The remaining participants took no
1025 medications.

1026 **Table 7:** Visitbook example (first three visits only)

| visit | visitlevel | form | formname |
|-------|----------------|-------|-------------------|
| 1 | Baseline visit | elig | Eligibility |
| 1 | Baseline visit | lab | Laboratory values |
| 2 | 1 month | visit | Visit info |
| 2 | 1 month | lab | Laboratory values |
| 3 | 2 month | visit | Visit info |
| 3 | 2 month | lab | Laboratory values |

1027 **Table 8:** Visit structure

| Form | Baseline | 1 month | 2 month |
|-------------------|----------|---------|---------|
| Eligibility | X | | |
| Laboratory values | X | X | x |
| Visit info | | X | x |

1028

1029

1030 **17.3. Meta data scheme from ISRCTN**

1031 Options are added in curled brackets if provided, an empty filed on the right hand side indicates
 1032 free text, “M” denotes mandatory fields.

1033 *General data*

| | | |
|-----------------------|---|---|
| Public title | M | |
| Overall trial status | | |
| Recruitment status | | |
| Plain English Summary | M | Who can participate? What does the study involve? Where is the study run from? When is the study starting and how long is it expected to run for? Who is funding the study? Who is the main contact? Trial website |

1034

1035 *Contact information*

| | | |
|--------------------|---|----------------------|
| Type | M | {Public, Scientific} |
| Primary contact | M | |
| ORCID ID | | |
| Contact details | M | |
| Additional contact | | |
| Type | | {Public, Scientific} |
| ORCID ID | | |
| Contact details | | |

1036

1037 *Additional identifiers*

| | | |
|---------------------------|---|--|
| EudraCT number | | |
| ClinicalTrials.gov number | | |
| Protocol/serial number | M | |

1038

1039 *Study information*

| | | |
|---------------------------|---|---|
| Scientific title | M | |
| Acronym | | |
| Study hypothesis | M | |
| Ethics approval | M | |
| Study design | M | Free text |
| Primary study design | M | {Not Specified, Interventional, Observational, Other} |
| Trial setting | | {Not Specified, Hospitals, GP practices, Other, Home, Internet, Community, Schools} |
| Trial type | M | {Not specified, Diagnostic, Other, Prevention, Quality of life, Screening, Treatment} |
| Patient information sheet | | |

| | | |
|-------------------------------------|---|---|
| Condition | M | Free text |
| Intervention | M | Free text |
| Intervention type | M | {Not specified, Drug, Supplement, Device, Biological/Vaccine, Procedure/Surgery, Behavioral, Genetic, Other, Mixed} |
| Phase | | |
| Drug names | | |
| Primary outcome measure | M | |
| Secondary outcome measures | M | |
| Overall trial start date | M | |
| Overall trial end date | M | |
| Reason abandoned (if study stopped) | | |

1040

1041 *Eligibility*

| | | |
|--------------------------------|---|--|
| Participant inclusion criteria | | |
| Participant type | M | {Not Specified, Healthy volunteer, Patient, Health professional, Carer, All, Mixed, Other} |
| Age group | M | {Not Specified, Adult, Senior, Neonate, Child, All, Mixed, Other} |
| Gender | M | {female, male, both} |
| Target number of participants | M | |
| Participant exclusion criteria | M | |
| Recruitment start date | M | |
| Recruitment end date | M | |

1042

1043 *Locations*

| | | |
|----------------------------|---|--|
| Countries of recruitment | | |
| Trial participating centre | M | |

1044

1045 *Sponsor information*

| | | |
|-----------------|---|---|
| Organisation | M | |
| Sponsor details | M | |
| Sponsor type | M | {Not defined, Charity, Government, Hospital/treatment centre, Industry, Other, Research council, Research organisation, University/education} |
| Website | | |
| Privacy | M | {Show all contact details, Hide telephone and email details} |

1046

1047 *Funders*

| | |
|----------------------|---|
| Funder type | M |
| Funder name | M |
| Alternative name(s) | |
| Funding Body Type | |
| Funding Body Subtype | |
| Location | |

1048

1049 *Results and Publications*

| | | |
|------------------------------------|---|---|
| Publication and dissemination plan | | |
| Intention to publish date | | |
| Participant level data | M | {Available on request, Not expected to be available, Stored in repository, Other, Not provided at time of registration, To be made available at a later date} |
| Basic results (scientific) | | |
| Publication list | | |
| Publication citations | | |

1050 **17.4. Information required for additional documentation**1051 **Table 9:** Information on supplied documentation

| Title | Size | Type | Format |
|---------------------------|--------|--------------|--------|
| Study_Protocol_final | 420 KB | Text | Pdf |
| Data_preparation | 67 KB | Stata script | Do |
| Statistical_analysis_plan | 180 KB | Text | Pdf |
| Consent_form | 56 KB | Text | Word |
| analysis_final | 120 KB | Stata script | Do |

1052

1053

1054 **17.5 Checklist for selecting a data repository**1055 **Table 10:** Selection criteria

| Item | Yes | No | Unsure | Potential indicators | Explanation |
|--|--------------------------|--------------------------|--------------------------|--|---|
| Is the repository trustworthy? | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <ul style="list-style-type: none"> • Certifications or public institution behind the repository? | |
| Will my data, information, and documentation be hosted? | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <ul style="list-style-type: none"> • Any restrictions on file type? • Any restrictions on file size? | |
| Will any legal requirements be met? | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <ul style="list-style-type: none"> • Licensing • Storage of sensitive data | |
| Does the repository support the sharing process? | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <ul style="list-style-type: none"> • On request ... | See chapter 11 |
| <i>FAIR data principles</i> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | Does the repository make the data findable, accessible, interoperable, and as reusable as possible for as long as required? | In order to sustain the value of the data, the repository has to comply with the FAIR principles. |
| Basic functionality | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <ul style="list-style-type: none"> • Single landing page per dataset • Unique identification number • Digital Object Identifier | |
| Does the repository allow for enough and the right meta-information? Is the metadata scheme specific for medical research? | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <ul style="list-style-type: none"> • Specific metadata fields on disease, intervention, outcome etc. | See chapter 8 |
| Long term preservation, sustainability | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <ul style="list-style-type: none"> • (might not be possible to assess) | Is there any plan in how long term preservation is ensured? For how long is storage guaranteed (for example, the repository of the Open Science Framework has a preservation fund that ensures hosting for 50+ years (based on present costs)). |
| Does the repository track usage and provide sufficient statistics? | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <ul style="list-style-type: none"> • Page views for each object/dataset • Number of downloads per object/dataset | |

1056