

1 1 September 2022  
2 EMA/787647/2022  
3 European Medicines Agency

4 **Good Practice Guide for the use of the Metadata**  
5 **Catalogue of Real-World Data Sources**  
6 **V 1.0**

Start of public consultation	27 September 2022
End of consultation	16 November 2022

7

Comments should be provided using this [template](#). The completed comments form should be sent to [metadata@ema.europa.eu](mailto:metadata@ema.europa.eu)

8

Keywords	Data sources, studies, metadata, study protocol, study report, data flows, data management, vocabulary, glossary, use cases, population
----------	---

9



10 Contents

11 **Abbreviations** ..... 3

12 **Glossary** ..... 3

13 **1. Introduction** ..... 5

14 **2. Purpose of this document**..... 5

15 **3. Format of the catalogue** ..... 6

16 **4. Use of the catalogue to assess the suitability of data sources** ..... 6

17 4.1. Reliability and relevance of data sources .....6

18 4.2. Assessing suitability of data sources with the catalogue .....7

19 4.3. Use cases .....9

20 4.3.1. *Planning of a study* .....9

21 4.3.2. *Assessment of a study protocol* ..... 11

22 4.3.3. *Assessment of a study report* ..... 11

23 4.3.4. *Writing of a study protocol or study report*..... 11

24 4.3.5. *Benchmarking of several data sources*..... 12

25 4.3.6. *Analysis of a data source used in a study*..... 12

26 **User guides** ..... **13**

27 **5. Description of the metadata list and definitions** ..... **13**

28 5.1. Metadata characterising the 'data source' ..... 13

29 5.1.1. Data source – Administrative details ..... 13

30 5.1.2. Data source – Data elements collected ..... 15

31 5.1.3. Data source - Quantitative descriptors ..... 19

32 5.1.4. Data source – Data flows and management ..... 20

33 5.1.5. Data source – Vocabularies and standardised dictionaries ..... 22

34 **6. Registering a data source in the Data source catalogue** ..... **26**

35 **7. Maintenance of information in the Data source catalogue** ..... **26**

36 **References** ..... **27**

37

38

## 39 Abbreviations

CDM	Common Data Model
EMA	European Medicines Agency
ENCePP	European Network of Centres for Pharmacoepidemiology and Pharmacovigilance
ETL	Extract, Transform, Load
EU	European Union
EUDPR	Regulation (EU) 2018/1725 on the protection of natural persons with regard to the processing of personal data by the Union institutions, bodies, offices and agencies and on the free movement of such data
EU PAS Register	European Union electronic register of post-authorisation studies
FAIR	Findable, Accessible, Interoperable, and Reusable
FDA	Food and Drug Administration
GDPR	Regulation (EU) 2016/679 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation)
HARPER	HARmonized Protocol template to Enhance Reproducibility
HMA	Heads of Medicines Agencies
ID	identification
IMI	Innovative Medicines Initiative
MINERVA	Metadata for data dIScoverability aNd study rEPLICability in obseRVAtional studies
OMOP	Observational Medical Outcomes Partnership
RWE	Real-world evidence
SIFPD	Structured Process to Identify Fit-for-Purpose Data
TEHDAS	Towards the European Health Data Space

40

## 41 Glossary

- 42 • Catalogue: A collection of dataset descriptions, which is arranged in a systematic manner and  
43 consists of a user-oriented public part, where information concerning individual dataset parameters  
44 is accessible by electronic means through an online portal.
- 45 • Common data model (CDM): Common structure and format for data that allows for interoperability,  
46 e.g., the efficient execution of the same analysis code against different local database for an efficient  
47 execution of programs against local data.
- 48 • Contributor: An institution that contributes content to the metadata catalogue.
- 49 • Data quality: Set of attributes of a data source that define its fitness for purpose for users' needs in  
50 relation to health research, policy making and regulation.
- 51 • Data source: Data set sustained by a specified organisation, which is the data holder. The data  
52 source is characterised by the underlying population that can potentially contribute records to it, the  
53 trigger that leads to the creation of a record in the data source, and the data model used in the data  
54 source.
- 55 • Dataset: a structured collection of electronic health data.
- 56 • Data characterisation: The summarisation of features of a data source, including quantitative  
57 measures.

- 58 • Data holder: any natural or legal person, which is an entity or a body in the health or care sector, or  
59 performing research in relation to these sectors, as well as Union institutions, bodies, offices and  
60 agencies who has the right or obligation, in accordance with this Regulation, applicable Union law or  
61 national legislation implementing Union law, or in the case of non-personal data, through control of  
62 the technical design of a product and related services, the ability to make available, including to  
63 register, provide, restrict access or exchange certain data.
- 64 • Extract, transform, load (ETL): A repeatable process for converting data from one format to another,  
65 such as from a source native format to a common data model format. In this process, mappings to  
66 the standardised dictionary are added. It is typically implemented as a set of automated scripts.
- 67 • FAIR (findable, accessible, interoperable, and reusable) principles:
- 68 ○ Findability: Any (healthcare) database that is used for analysis should, from a scientific  
69 perspective, persist for future reference and reproducibility. A comprehensive record of the  
70 database in terms of purpose, sources, vocabularies and terms, access-control mechanisms,  
71 licence, consents, etc., should be available.
- 72 ○ Accessibility: Data should be accessible through a standardised and well-documented  
73 method.
- 74 ○ Interoperability: The ability of organisations as well as software applications or devices from  
75 the same manufacturer or different manufacturers to interact towards mutually beneficial  
76 goals, involving the exchange of information and knowledge without changing the content of  
77 the data between these organisations, software applications or devices, through the  
78 processes they support.
- 79 ○ Reusability: For data to be reusable, the data licences should explicitly allow the data to be  
80 used by others, and the data provenance (understanding how the data came into existence)  
81 needs to be specified and updated as needed.
- 82 • Institution: An organisation connected to one or more data sources—such as a Data Holder, or a  
83 research organisation running a study.
- 84 • Metadata: A set of data that describes and gives information about a dataset. More specifically,  
85 information describing the generation, location, and ownership of the data set; key variables; and  
86 the format (coding, structured versus not) in which the data are collected is needed to enable  
87 accurate identification and qualification of the exposure and outcome information available. Metadata  
88 also include the provenance and time span of the data, clearly documenting the input, systems, and  
89 processes that define data of interest. Finally, metadata include details on the storage, handling  
90 processes, access, and governance of data.
- 91 • Underlying population: The population of individuals in a geographical location who can *potentially*  
92 contribute information to a data source. This is a population defined by an administrative  
93 characteristic, a disease, a medical condition or any other relevant characteristic.
- 94 • Vocabulary: Standardised medical terminologies; may be an international standard  
95 (e.g., International Classification of Diseases, Anatomical Therapeutic Chemical) or a country/region-  
96 specific system or modification.

## 97 **1. Introduction**

98 Identification of appropriate data sources is becoming an increasing need for regulatory decision making.  
99 While data needs are becoming more complex, standardised information and statistics on real-world data  
100 sources is lacking. Metadata are descriptive data that characterise other data to create a clearer  
101 understanding of their meaning and to achieve greater reliability and quality when using the data for a  
102 specific purpose. Access to a standard and electronic set of complete and accurate metadata information  
103 can contribute to identifying the data sources suitable for a specific study, facilitate description of the  
104 data sources planned to be used in a study protocol or research proposal, and contribute to assessing  
105 the evidentiary value of the results of studies.

106 The Heads of Medicines Agencies–European Medicines Agency (HMA-EMA) joint Big Data Task Force  
107 recommended “to promote data discoverability through the identification of metadata” as part of its  
108 Recommendation III: “*Enable data discoverability. Identify key meta-data for regulatory decision making*  
109 *on the choice of data source, strengthen the current European Network of Centres for*  
110 *Pharmacoepidemiology and Pharmacovigilance (ENCePP) resources database to signpost to the most*  
111 *appropriate data, and promote the use of the FAIR principles (Findable, Accessible, Interoperable and*  
112 *Reusable)*” (HMA-EMA, 2020). This goal is therefore included in the 2020-2021 Work Plan of the HMA-  
113 EMA joint Big Data Steering Group (HMA-EMA Big Data Steering Group, 2022).

114 To fulfil this mandate, EMA in November 2020 the study “Strengthening Use of Real-World Data in  
115 Medicines Development: Metadata for Data Discoverability and Study Replicability” (MINERVA; EU PAS  
116 Register number EUPAS39322). The main focus of the study was the definition of a set of metadata on  
117 real-world data sources, including engagement with stakeholders to reach broad agreement and the  
118 development of a good practice guide describing the metadata and recommendations based on a pilot.

119 Based on the results of the MINERVA study and the consultation of the ENCePP community and other  
120 stakeholders, the EMA is developing an electronic catalogue that will provide metadata for real-world  
121 data sources. This catalogue has two objectives: 1) to facilitate the *discoverability* of data sources to  
122 generate adequate evidence for regulatory purpose, i.e., the initial identification of data sources suitable  
123 to investigate a specific research question, and 2) to support the assessment of study protocols and  
124 study results by providing quick access to information on the suitability of data source(s) proposed to be  
125 used in the study protocol or referred to in the study report.

126 The *Good Practice Guide for the use of the Metadata Catalogue of Real-World Data Sources* has been  
127 developed to provide regulators, researchers and other interested stakeholders with recommendations  
128 on the use of the EU metadata catalogue of real-world data sources.

## 129 **2. Purpose of this document**

130 The Good Practice Guide aims to provide recommendations for the use of the EU metadata catalogue to  
131 identify real-world data sources suitable for specific research questions and to assess the suitability of  
132 data sources proposed to be used in a study protocol or referred to in a study report.

133 It also provides a detailed description of all the metadata elements as envisaged to be used in the EMA  
134 catalogue, which have been published by HMA/EMA in the List of metadata for Real World Data  
135 catalogues<sup>1</sup>, and it guides the user for the insertion and maintenance of data in the catalogue.

136 The catalogue is targeted for release in late 2023.

---

<sup>1</sup> HMA/EMA. [List of metadata for Real World Data catalogues](#) (2022).

### 137 **3. Format of the catalogue**

138 The structure of the catalogue is based on the MINERVA catalogue pilot project.<sup>2</sup> A **data source** is a  
139 data collection (or a set of linked data collections) sustained by a specified organisation, which is the  
140 data holder. It is characterised by the underlying population that can potentially contribute records, the  
141 event triggering the creation of a record in the data source and the data model. The mechanisms that  
142 put data into existence are heterogeneous across data sources. The catalogue is therefore divided into  
143 the following sections allowing to capture the variety of existing data sources and facilitate data  
144 discoverability: *Characteristics, Population, Data elements, Data flows and management* and  
145 *Vocabularies*. It is composed of qualitative information and quantitative metadata, e.g. counts and  
146 demographic distributions of the underlying population.

147 The catalogue follows good practices for data management:

- 148 • FAIR principles are complied with: the data are Findable, Accessible, Interoperable and Reusable,<sup>3</sup>  
149 and there is interoperability with the EU PAS register for studies conducted with the data sources  
150 and with other catalogues to be developed in the future.
- 151 • A controlled data entry process is run for the initial collection of metadata by the data holder, regular  
152 updates of metadata are foreseen with trusted relationship between the data holder and the EMA.
- 153 • Change management and reproducibility are supported by enabling data holders of a data source to  
154 edit the corresponding metadata while ensuring that the attribution of each data entry is traceable  
155 via appropriate version control, and by enabling the creation of a copy of the metadata and their  
156 update by the data holders.
- 157 • Quantitative metadata for data sources are provided at the level of the total and active populations.
- 158 • Personal data will be processed in compliance with European data protection legislation and, in  
159 particular, Regulation (EU) 2018/1725 (EUDPR) and Regulation (EU) 2016/679 (GDPR) as applicable.  
160 In this regard, EMA will publish a record of processing activity and a data protection notice as  
161 required. A quality management process is in place, including an incident management system, a  
162 disaster recovery plan and a quality assurance office.

163

## 164 **4. Use of the catalogue to assess the suitability of data** 165 **sources**

### 166 **4.1. Reliability and relevance of data sources**

167 The assessment of the suitability of data sources for studies needs to consider the differences between  
168 studies with primary data collection and studies based on secondary use of data already collected for  
169 another purpose, such as patient monitoring, healthcare reimbursement, quality management or another  
170 administrative purpose. In primary data collection, the study itself applies and controls all the quality  
171 management steps related to the data collected. In secondary data collection, use of already collected  
172 data relies on existing processes for data quality, i.e., which data have been collected for the initial  
173 purpose and how they were generated, and many aspects of the data processes, i.e., how the data were  
174 coded, curated, validated and stored.

---

<sup>2</sup> MINERVA: Strengthening Use of Real-World Data in Medicines Development: Metadata for Data Discoverability and Study Replicability (2022). [EUPAS39322](https://eupas39322.europa.eu/)

<sup>3</sup> FAIR Principles. <https://www.go-fair.org/fair-principles/>

175 The assessment of the suitability of data sources should therefore differentiate between two broad  
176 aspects of data quality<sup>4,5</sup>:

177 - quality in relation to the *reliability* of the primary data, based on e.g. the detection and correction  
178 of errors, missing data and implausible values, the verification and validation of formats, codes,  
179 values, time components and underlying calculations, the presence of unique identification numbers  
180 for each person and the documentation of standardised processes leading to entry and exit of person;  
181 this aspect of quality is a characteristic of the data source independent from its use for a specific  
182 study.

183 - quality in relation to the *relevance* of the data source to provide adequate and valid evidence  
184 informing a specific research question following the application of appropriate epidemiological and  
185 statistical techniques; this aspect requires adequate information on the format and content of the  
186 data source, such as the presence of the data needed for the study, the numbers of individuals  
187 included, population characteristics, coding terminologies, the availability and completeness of data  
188 elements and the time span of the data; this aspect of quality is partly dependent on the research  
189 question as some data characteristics (such as some data elements or age range of the population)  
190 may be required for some studies and not for others.

191 Several data quality frameworks have been proposed to help understand the strengths and limitations  
192 of a data source to answer a research question and the impact they may have on the suitability of data  
193 sources for a specific study<sup>6,7,8</sup>. These data quality frameworks differ as to the specific dimensions  
194 included (with varying levels of details and names used to describe these dimensions) and the methods  
195 used to assess them, and some frameworks address both the data reliability and relevance or only one  
196 of these. In Europe, the Towards European Health Data Space (TEHDAS) project has set out and defined  
197 six dimensions deemed the most important ones at data source level: reliability, relevance, timeliness,  
198 coherence, coverage and completeness.<sup>4</sup>

## 199 **4.2. Assessing suitability of data sources with the catalogue**

### 200 **Reliability**

201 The metadata catalogue provides information allowing an initial evaluation of the suitability of data  
202 sources. Information on the following aspects of *reliability* is provided:

- 203 • Data management, including the possibility of data validation (elements C2.7, C2.9, C8.5 and  
204 C8.5.1), the mapping to a CDM (D1.2.1.1, D1.2, D1.2.1, D1.4 and D1.7)
- 205 • The data source ETL process and status (B7.1 to B7.5)
- 206 • Any qualification received (C3.1, C3.1.1)
- 207 • Governance details as regards data capture and management, data quality checks and validation of  
208 results (C2.3)
- 209 • The process of collecting and recording the data (C4.3), linkage information (B5.2, B.5.2.1, B5.3,  
210 B4.1)

---

<sup>4</sup> ENCePP Guide on Methodological Standards in Pharmacoepidemiology, 10<sup>th</sup> Rev. (2022). [Chapter 12.1 General principles of quality management](#)

<sup>5</sup> Wang S., Schneeweiss S. [Assessing and Interpreting Real-World Evidence Studies: Introductory Points for New Reviewers](#). Clin Pharmacol Ther. 2022;111(1):145-149.

<sup>6</sup> ENCePP Guide on Methodological Standards in Pharmacoepidemiology, 10<sup>th</sup> Rev. (2022). [Chapter 12.2. Data Quality Frameworks](#).

<sup>7</sup> TEHDAS. [European Health Data Space Data Quality Framework](#) (2022).

<sup>8</sup> HMA/EMA. Data Quality Framework for EU medicines regulation (2022).

- 211 • All vocabularies used in the data source
- 212 • A link to the publications describing the data sources (e.g. validation, data elements,  
213 representativity).

214 Access to raw data and computational resources would be required for a more in-depth assessment of  
215 reliability, for example a verification of the records and values, data validation against reference or  
216 plausible values and other computations. Such assessment should be performed by the data holders and  
217 periodically updated. The data holders should make the methods and the results of the assessment  
218 publicly available for consultation to support the assessment and replication of studies.

## 219 **Relevance**

220 The metadata catalogue is also suitable for an initial evaluation of the *relevance* of the data sources to  
221 generate valid evidence informing a specific research question based on the study design, e.g. to  
222 implement step 3 of the Structured Process to Identify Fit-for-Purpose Data (SIFPD)<sup>9</sup> or the Population,  
223 Intervention, Comparison, Outcome and Time horizon (PICOT) format.<sup>10</sup> The catalogue also provides the  
224 data elements to be included in the table of data sources recommended by the HARmonized Protocol  
225 template to Enhance Reproducibility (HARPER).<sup>11</sup> The assessment of relevance is supported by the  
226 availability of the following variables:

- 227 • Setting: county(-ies) (C1.5), region(s) (C1.5.1), type of data source (C5.1 and C5.1.1), care setting  
228 (C1.14).
- 229 • Population: total and active population size (C7.1), percentage of the population covered by the data  
230 source in the catchment areas (C1.11.2) and description of the population for which data are not  
231 collected (C1.11.1), age groups (C1.8), sociodemographic information (C6.7), lifestyle factors  
232 (C6.8), family linkage (C6.6, C6.6.1), availability of data on pregnancy and neonates (C1.9), trigger  
233 for registration (C1.6, C1.6.1) and de-registration (C17.1, C1.7.1), median time between first and  
234 last records for all individuals (B6.3) and active individuals (B6.3.1).
- 235 • Exposure: availability of data on prescriptions and/or dispensing (C6.13), ATMPs (C6.16),  
236 contraception (C6.17), vaccines (C6.19), other injectables (C6.19), medical devices (C6.20),  
237 procedures (C6.21), medicinal products (C6.15.1) and indication (C6.18), biomarker data (C6.26).
- 238 • Outcomes: availability of data on hospital admission or discharge (C6.10), ICU admission (C6.10.1),  
239 death and cause of death (C6.11), clinical measurements (C6.23), genetic data (C6.25), patient-  
240 generated data (C6.27), health care utilisation (C6.29), diagnostic codes (C6.9), specific diseases  
241 (C1.10), with disease information collected (C1.10.1).
- 242 • Time elements: date when the data source was established (C4.5), first collection date (C1.12) and  
243 last collection date (C1.13), median time between the first and the last available records for unique  
244 individuals captured in the data source (B6.3) and for unique active individuals (B6.3.1).

245 Links to the EU PAS Register also allow to identify studies that have been performed with the same data  
246 source, allowing an evaluation of the analyses that can be performed.

---

<sup>9</sup> Gatto, N. M., Campbell, U. B., Rubinstein, E., Jaksa, A., Mattox, P., Mo, J., & Reynolds, R. F. (2022). The Structured Process to Identify Fit-For-Purpose Data: A Data Feasibility Assessment Framework. *Clin Pharmacol Ther.* 2022;111(1), 122–134. <https://doi.org/10.1002/cpt.2466>

<sup>10</sup> Brown, P., Brunnhuber, K., Chalkidou, K., Chalmers, I., Clarke, M., Fenton, M., Forbes, C., Glanville, J., Hicks, N. J., Moody, J., Twaddle, S., Timimi, H., & Young, P. How to formulate research recommendations. *BMJ.* 2006;333(7572), 804–806. <https://doi.org/10.1136/bmj.38987.492014.94>

<sup>11</sup> Wang S, Pottgard A, Crown W et al. HARmonized Protocol Template to Enhance Reproducibility (HARPER) of Hypothesis Evaluating Real-World Evidence Studies on Treatment Effects: A Good Practices Report of a Joint ISPE/ISPOR Task Force. *Pharmacol Drug Saf.* 2022;

247 In order to provide adequate evidence, appropriate epidemiological and statistical methods must be  
248 applied to the study design and the analysis and interpretation of data generated from a real-world data  
249 source. These methods are not addressed by the metadata catalogue but are described in other  
250 guidance, e.g. the ENCePP Guide on Methodological Standards in Pharmacoepidemiology, 10<sup>th</sup> Rev.  
251 (2022).

### 252 **4.3. Use cases**

#### 253 **4.3.1. Planning of a study**

254 *Use case: An investigator wants to identify suitable data sources for a planned study.*

255 The process for identification of suitable data sources may follow six successive steps (Figure 1):

- 256 1. In a first step, the investigator searches the catalogue to identify relevant data sources fulfilling the  
257 specifications of the research question or, if there is a prior interest in using a specific data source,  
258 to access the record for this data source and consult the available information. The search may  
259 initially use the data elements useful to assess pre-defined PICOT criteria (see section 4.1) in order  
260 to identify possibly suitable data sources.
- 261 2. In a second step, the investigator accesses the record of each potential data source and screens  
262 more detailed information on the availability of data (incl. quantitative metadata) on the population,  
263 exposures, outcomes and confounding variables to confirm that the data source may be relevant to  
264 answer the research question.
- 265 3. In a third step, the investigator consults information on the governance, accessibility and availability  
266 of the data sources (C2.3) to determine whether they are accessible, as well as the conditions related  
267 to this use, and whether the investigator would be eligible to receive aggregated information or get  
268 access to raw data.
- 269 4. In a fourth step, the investigator screens the metadata allowing to perform a preliminary assessment  
270 of the reliability of each potential data source based on important quality aspects of the data source  
271 that are relevant for the specific study (see section 3.1). Publications describing the data source and  
272 its validation can be extracted and consulted. Missing information for some of these variables may  
273 raise doubts about the presence of an adequate quality management process or may question  
274 whether the data holder gives sufficient attention to quality management.

275 At this stage, the investigator should establish a first list of candidate data sources (if there is no *a*  
276 *priori* choice of a specific data source).

- 277 5. In a fifth step, the investigator uses the link providing access to the EU PAS Register of studies that  
278 have been performed with the same data source and addressed research questions similar (as to the  
279 topic or study design) to the current one. After selecting studies with a similar topic or design as for  
280 the planned study, the investigator accesses the study information to:

- 281 • confirm the suitability of the data source as regards to the PICOT criteria; if the study protocol  
282 and/or the study report have been uploaded, more granular information can be extracted on  
283 the time frame for the use of the database, the number of active study participants originating  
284 from the data source (providing useful information for the sample size calculation of the current  
285 study), the data elements used for the study (e.g. exposure and outcome variables,  
286 confounding factors), variable definitions and vocabularies (and any need for mapping of  
287 terms), the transformation of data into categories and the analyses that could be performed  
288 with the data;

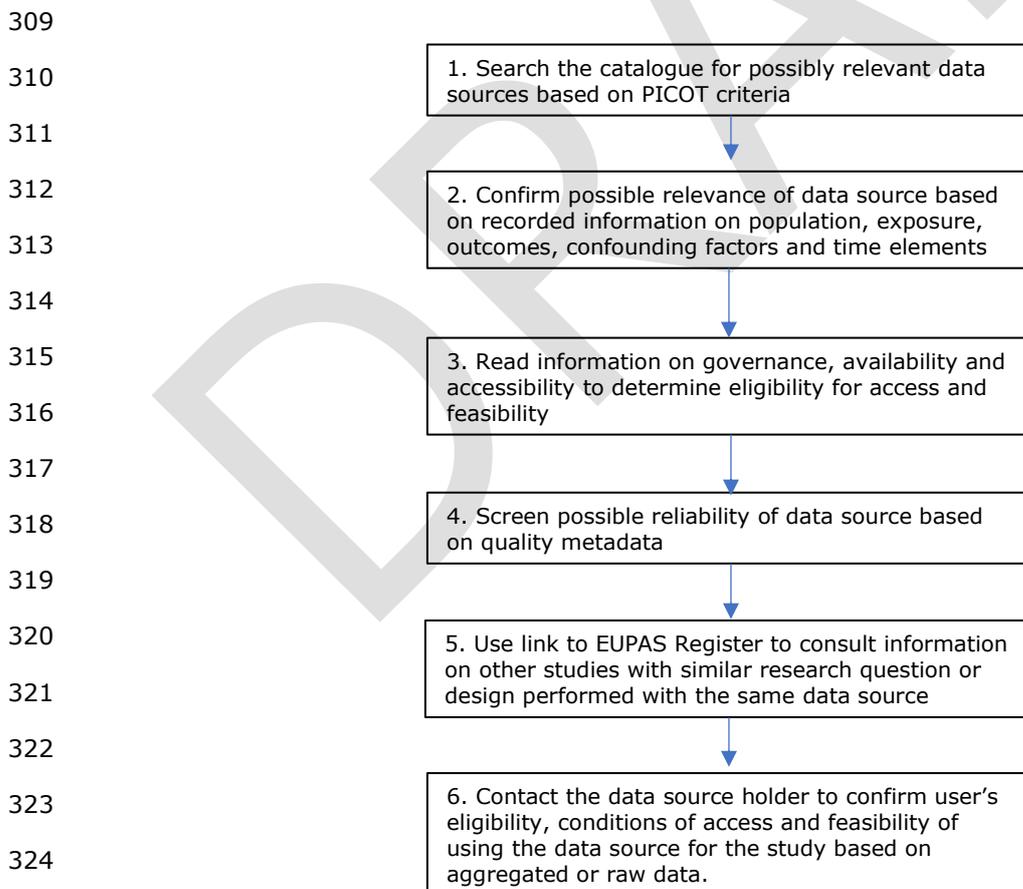
- 289 • check in the study protocol or study report (if available) the algorithms that have been used to  
290 identify diseases or outcomes of interest and their severity (for example persons with a rare  
291 disease if applicable) and which prompts, and contents were used in such algorithm(s);
- 292 • learn about the strengths and weaknesses of the data source encountered in the study conduct;  
293 in case a limitation is acknowledged that the data source is not optimal to identify all the  
294 variables of interest (e.g. diagnosis of the disease, levels of severity, treatments, confounding  
295 variables), use of the data source should be reconsidered or a strategy could be devised to  
296 complement the information obtained from the data source with that from another, possibly by  
297 data linkage;
- 298 • search for use of the data source in studies published in peer-review journals and comments  
299 made on study limitations.

300 If there are remaining uncertainties as regards the reliability and relevance of the data source for  
301 the specific study, the investigators of similar studies in terms of topic or study design can be  
302 contacted to gather additional information.

303 If past studies using the same data source cannot be found, it may be preferable to investigate the  
304 information available for another relevant data source.

- 305 6. If the previous steps have been successful, the data holders of the data sources of interest can be  
306 contacted to discuss the feasibility of using the data sources for the specific study and the conditions  
307 of this use.

308 Figure 1. Steps for using the metadata catalogue when planning a real-world study



### 325 **4.3.2. Assessment of a study protocol**

326 *Use case: A data source is mentioned in the study protocol submitted for a study and the assessor needs*  
327 *to understand in detail the suitability of the data source proposed to be used.*

328 The user may verify if the data source has been registered in the catalogue.<sup>12</sup> Depending on the  
329 information that is already available in the protocol, that is missing or that needs verification, the user  
330 accesses different sections of the catalogue. In order to verify the representativeness of the study  
331 population described in the report, the user may verify qualitative information, such as the geographical  
332 coverage, the type of data source, the care setting and the trigger for registering a person in the data  
333 source, as well as quantitative metadata on the percentage of the population covered by the data source  
334 in the catchment area and the estimated sample size of active patients per age category.

335 In the Data elements section, the assessor may find information on exposure, outcomes and covariates  
336 collected in the data source and identify those that have not been proposed to be extracted but could be  
337 useful to include for the study.

338 The assessor can also explore technical information supporting the evaluation of the protocol such as  
339 the vocabularies used to define variables, the process of data collection, the CDM, the ETL specifications  
340 and any linkage strategy.

341 The extent of the validation of the data source and the possibility to contact patients provides regulatory  
342 assessors of studies required to pharmaceutical companies information about the need and the possibility  
343 to request additional data validation. The link to studies using the same data source and registered in  
344 the EU PAS register will allow to further document use cases where the data source was used with its  
345 strengths and limitations.

### 346 **4.3.3. Assessment of a study report**

347 *Use case: A data source is mentioned in the study report or publication and the reader needs to*  
348 *understand the suitability of the data source used in the study to interpret its results.*

349 The process is similar to the process described above for the assessment of a study protocol. The main  
350 difference resides in the fact that the study report contains results and generally quantitative information  
351 on the characteristics of the study population originating from the data source. The assessor may  
352 therefore identify, and investigate if needed, differences between the information provided in the study  
353 report and in the metadata catalogue.

354 Some verification may be applied to the description of the study population, the sample size originating  
355 from the data source included in the report, the nature and categories of variables included in the analysis  
356 and the coding system provided. Insight into the characteristics of the data source also helps interpret  
357 the study results and understand the strengths and limitations of the study independently from the  
358 investigator's own interpretation.

### 359 **4.3.4. Writing of a study protocol or study report**

360 *Use case: An investigator writes a study protocol or a study report for which he needs to describe the*  
361 *data source(s) proposed to be used or used in the study. The information on the data source he finds in*

---

<sup>12</sup> Except of specific circumstances, there is no legal obligation to register a data source into the metadata catalogue. It is however expected that data source holders will register their data source, and update the record, whenever it will be used for public health or regulatory purpose, as absence of public information on the data source may affect the scientific credibility and public confidence on study results. In case where a data source user has got access to a data source based on a contractual agreement, the contract may include a provision that the data source is registered, or the record updated, in the metadata catalogue as part of the agreement.

362 *other publications or other documentation is heterogeneous, and a comparison between the*  
363 *characteristics of several databases used in to be used or used in the study is difficult to perform.*

364 The investigator can extract from the metadata catalogue standardised information on each data source  
365 and provide a reference to public information for the registered data sources. He can provide in the  
366 Methods section of the protocol or report the identification number and the link of the data source in the  
367 catalogue.

368 If a data source is not registered in the metadata catalogue, this registration can be made simultaneously  
369 to the writing of the protocol or report. If access to the data source has been obtained through a  
370 contractual agreement, this agreement could provide for the registration of the data source, or updating  
371 of its record, before the study commences.

#### 372 **4.3.5. Benchmarking of several data sources**

373 *Use case: A data holder or data user may wish to compare the characteristics of a specific data source*  
374 *with other ones covering fully or partially the same population.*

375 The different data sources may have different primary purposes, contain different data elements and  
376 cover different population groups. It is nevertheless important to be able perform comparisons to help  
377 understand the heterogeneity of results obtained in some analyses conducted in the same country or  
378 region or to perform a validation of a data source in comparison to another one considered a gold  
379 standard. For this purpose, the metadata catalogue provides:

- 380 • a harmonised description of the characteristics of each data source that allow to compare differences,  
381 e.g., in age groups covered
- 382 • information on common variables and variable categories by which analyses can be stratified to map  
383 sources of heterogeneity
- 384 • information on possible linkages with other data sources, including availability of linkages to the same  
385 data sources (or cross-linkage between data sources) allowing to harmonise data on the same  
386 individuals and provide additional information, e.g. on confounding factors.

#### 387 **4.3.6. Analysis of a data source used in a study**

388 *Use case: An investigator, statistician or analyst wants to benefit from the experience of others for the*  
389 *programming of the data transformation and statistical analysis.*

390 If the study is implemented in a CDM, the analyst may find in the catalogue the specifications of the ETL  
391 procedure from the data source to the CDM. Irrespective of whether the data holder has converted to  
392 the CDM the entire data source, or only an extraction thereof, this information supports the programmer  
393 in developing the study script. Using the link to the EU PAS Register, the analyst can also access detailed  
394 information on the studies performed with the same data source and registered in the EU PAS Register,  
395 and select the studies that investigated the same topic and/or study design. The study protocol or  
396 statistical analysis plan of these studies may contain information on how to operationalise the variables  
397 of the study in their respective data sources. The detailed programming script may also be available in  
398 a public repository, e.g., a GitHub repository.

399 At the end of the analysis, the analyst should also record the script of the analysis in a public repository  
400 and provide the link in the EU PAS Register, thus enabling transparency and quality control and  
401 facilitating reproducibility.

## 402 **User guides**

### 403 **5. Description of the metadata list and definitions**

404 Following several prioritisation exercises and consultations with stakeholders, the below metadata  
405 elements, which have been published by HMA/EMA in the List of metadata for Real World Data  
406 catalogues<sup>13</sup>, have been selected for a first iteration of this process. The data elements aim to describe  
407 the data sources, with a view of facilitating the choice of data source for the specific use cases listed in  
408 chapter 4.

#### 409 **5.1. Metadata characterising the 'data source'**

410 A data source is described by the data holder that sustains the collection of records in the data source,  
411 the underlying population that can potentially contribute records to a data source, and the prompt that  
412 leads to creation of a record in the data source.

##### 413 **5.1.1. Data source – Administrative details**

###### 414 **5.1.1.1. Name of the data source (C1.2)**

415 The name of the data source, as used in European projects, must be provided. If the database  
416 is widely known by several names, these can be provided in this field, separated by a '/' sign.  
417 Where the name of the data source is in a local language, the English translation should also be  
418 provided, using parentheses.

###### 419 **5.1.1.2. Data source acronym (C1.3)**

420 Where the data source is generally known under a specific acronym, this should be provided.

###### 421 **5.1.1.3. Data holder (C4.1)**

422 The data holder must be provided, selecting one of the existing entries from the 'institutions'  
423 available look-up. For the purpose of this catalogue, a data holder is defined as an organisation  
424 that sustains the collection of records in a data source.

###### 425 **5.1.1.4. Data source contact name (M1.3)**

426 A contact name should be provided for queries related to the data source. The contact details  
427 would be visible in the publicly available catalogue.

###### 428 **5.1.1.5. Data source contact email (M1.6)**

429 An e-mail contact should be made available for queries related to the data source. This  
430 information will be visible in the publicly available catalogue.

###### 431 **5.1.1.6. Data source countries (C1.5)**

432 The country where the data originate should be selected from the list of country codes (ISO  
433 3166-1).

434 Where needed, multiple countries can be selected.

---

<sup>13</sup> HMA/EMA. [List of metadata for Real World Data catalogues](#) (2022).

435 **5.1.1.7. Data source language(s) (C6.2)**

436 The data source language should be specified using the appropriate ISO 639 code.

437 **5.1.1.8. Data source regions (C1.5.1)**

438 The geographical regions that the data source covers should be provided using regions codes  
439 (ISO 3166-2). Multiple regions can be selected where required.

440 **5.1.1.9. Date when the data source was first established (C4.5)**

441 The date when the data source was first set-up. This date can be different from the 'first collection  
442 date' (C1.12).

443 **5.1.1.10. First collection date (C1.12)**

444 The date when data started to be collected or extracted.  
445 It is expected that this information is populated only once, when the data source is first described  
446 (with the exception of error corrections from the initial submission).

447 **5.1.1.11. Last collection date (C1.13)**

448 Where applicable, the date when the data collection ended. This information should only be  
449 provided for data sources where the data collection has stopped permanently.

450 **5.1.1.12. Data source website (C11.1)**

451 Where such an information is available, a link to the dedicated webpage describing the data  
452 source should be provided. The information listed would capture information such as data  
453 content, release notes etc.

454 **5.1.1.13. Data source publications (C11.2)**

455 A list of peer-reviewed papers or documents describing the data source (validation, data  
456 elements, representativity) or its use for pharmacoepidemiologic research

457 **5.1.1.14. Data source qualification (C3.1, C3.1.1)**

458 If the data source has successfully undergone a formal qualification process (e.g., from the EMA,  
459 or ISO or other certifications), this should be described.

460 **5.1.1.15. Main financial support (C4.6)**

461 The source of finance for the data source in the last three years should be specified using the  
462 below categories:

- 463 - Funding by own institution
- 464 - National, regional, or municipal public funding
- 465 - European public funding
- 466 - Funding from industry or contract research organisation
- 467 - Funding from public-private partnership
- 468 - Funds from patients organisations, charity or foundation

469 **5.1.1.16. Data source type (C5.1, C5.1.1)**

470 Data source may fit in one of more of the following categories:

- 471 Administrative
  - 472 - population registry
  - 473 - death registry

- 474 - registration with healthcare system
- 475 - exemptions from co-payment
- 476 - diagnostic tests or procedures reimbursement
- 477 - administrative healthcare claims
  
- 478 Primary care
- 479 - primary care medical records
- 480 - pharmacy dispensation records
  
- 481 Secondary care
- 482 - hospital discharge records
- 483 - hospital inpatient records
- 484 - hospital outpatient visit records
- 485 - emergency care discharge records
- 486 - specialist ambulatory care records
- 487
- 488 Registries
- 489 - birth registry
- 490 - induced terminations registry
- 491 - congenital anomaly registry
- 492 - cancer registry
- 493 - disease registry
- 494 - vaccination registry
- 495 - drug registry
- 496
- 497 Other
- 498 - biobank
- 499 - spontaneous reporting of adverse drug reactions
- 500 If none of the listed categories apply to the data source, it's type should be described in the
- 501 available free text field (C5.1.1).

502 **5.1.1.17. Care setting for data source (C1.14)**

- 503 Where the data source describes a care setting, this can be further characterised as:
- 504
- 505 - primary care – GP, community pharmacist level
  - 506 - primary care – specialist level (e.g. paediatricians)
  - 507 - secondary care – specialist level (ambulatory)
  - 508 - hospital inpatient care
  - 509 - hospital outpatient care

510 **5.1.2. Data source – Data elements collected**

511 **5.1.2.1. Data source characteristics**

512 To characterise the content of the data source the specific data elements should be selected as applicable.

513

514

Value (yes/no)	Description
Specific diseases (C1.10)	Data source collects information with a focus on specific diseases. This might be a patient registry or other similar initiatives. Where this is applicable
Hospital admission discharge (C6.10)	Information on hospital admission and/or hospital discharge is available in the data source.
ICU admission (C6.10.1)	Information on intensive care admission available.

Value (yes/no)	Description
Cause of death (C6.11)	The cause of death is captured, either as structured or unstructured information
Rare diseases (C6.12)	The data source captures rare diseases, where the prevalence of the condition in the EU is less than 5 in 10,000
Prescriptions and/or dispensing (C6.13)	The data source contains information on prescriptions or dispensing of medicines
ATMP (C6.16)	A medicinal product for human use that is either a gene therapy medicinal product, a somatic cell therapy product or a tissue engineered products as defined in Regulation (EC) No 1394/2007 [Reg (EC) No 1394/2007 Art 1(1)].
Contraception (C6.17)	Any information on use of any type of contraception (oral, injectable, devices etc.)
Indication for use (C6.18)	Therapeutic indication for the use of medicinal product
Administration of vaccines (C6.19)	Information on any vaccines administered
Administration of other injectables (C6.19.1)	Information on medicinal products administered via an injectable route (e.g.: solutions for perfusion, solutions for injection)
Medical devices (C6.20)	Where data source captures information on medicinal devices (e.g.: pens, syringes, inhalers)
Procedures (C6.21)	Medical procedures (e.g. surgical interventions, tests)
Clinical measurements (C6.23)	Information on clinical measurements (e.g.: BMI, blood pressure, height)
Healthcare provider (C6.24.1)	Data on individual health professionals or a health facility organization licensed to provide health care diagnosis and treatment services including medication, surgery and medical devices
Genetic data (C6.25)	Data related to genotyping, genome sequencing
Biomarker data (C6.26)	The term "biomarker" refers to a broad subcategory of medical signs – that is, objective indications of medical state observed from outside the patient – which can be measured accurately and reproducibly. For example, haematological assays, infectious disease markers or metabolomic biomarkers.
Patient-generated data (C6.27)	Health-related data created, recorded, or gathered by or from patients (or family members or other caregivers) to help address a health concern
Units of healthcare utilisation (C6.29)	Quantification of the use of services for the purpose of preventing or curing health problems (e.g.: number of visits to GP per year, number of hospital days)
Unique identifiers for persons (C6.4)	Where applicable, if patients are uniquely identified
Diagnostic codes (C6.9)	If diagnostic codes are captured; further information will be captured in section 5.1.5.11

Value (yes/no)	Description
Pregnancy and neonates (C1.9)	Where data on pregnant women and neonates (under 28 days of age), infant, and child development

515

DRAFT

516 **5.1.2.2. Disease information collected (C1.10.1)**

517 The disease or diseases for which information is collected should be specified in this field, using  
518 MedDRA terminology.

519 **5.1.2.3. Population age groups (C1.8)**

520 The information on the following age groups are being captured separately:

- 521 - newborn infants (0 to 27 days),
- 522 - infants and toddlers (28 days to 23 months),
- 523 - children (2 to 11 years),
- 524 - adolescents (12 to 17 years),
- 525 - adults (18 to 45 years),
- 526 - adults (46 to 64 years),
- 527 - adults (65 to 74 years),
- 528 - adults (75 to 84 years),
- 529 - adults (85 years and over)

530 **5.1.2.4. Family linkage (C6.6, C6.6.1)**

531 Where family linkage is made available in the data source this should be characterised using one  
532 or more of the following values: household (where the information on individuals sharing a  
533 household can be identified), mother-child, father-child, sibling.

534 If family linkage is not available, it should be specified if familial linkage can be created on an  
535 ad-hoc basis (C6.6.1).

536 **5.1.2.5. Sociodemographic information collected (C6.7)**

537 Where one or more of the following specific sociodemographic information are captured by the  
538 data source, these should be selected:

- 539 - age
- 540 - gender
- 541 - ethnicity
- 542 - country of origin
- 543 - indicator of socioeconomic status
- 544 - marital status
- 545 - education level
- 546 - type of residency
- 547 - living in rural area
- 548 - health area
- 549 - deprivation index

550 **5.1.2.6. Lifestyle factors (C6.8)**

551 Where the data source captures this information, one or more of the following lifestyle factors  
552 can be selected:

- 553
- 554 - tobacco use
- 555 - alcohol use
- 556 - amount of physical exercise
- 557 - diet



558 **5.1.2.7. Population covered by the data source (C1.11.2)**

559 The percentage of the population covered by the data source in the catchment area should be  
560 specified.

561 **5.1.2.8. Population not covered by the data source (C1.11.1)**

562 The description of the population covered by the data source in the catchment area whose data  
563 are not collected, where applicable (e.g.: people who are registered only for private care).

564 **5.1.3. Data source - Quantitative descriptors**

565 This section aims to collect a limited amount of data elements that look at the quantitative details  
566 of the data source. In future iterations of the data source catalogue this section can be further  
567 expanded as found useful.

568 **5.1.3.1. Population size (C7.1)**

569 The total number of unique individuals with records captured in the data source.

570 **5.1.3.2. Population size by age (C7.3)**

571 Where this information can be extracted, the number of unique individuals split by age groups  
572 should be captured.

573 **5.1.3.3. Active population size (C7.1.1)**

574 An active population for administrative healthcare data refers to the collection of patients for  
575 which there is an active record in the practice, i.e. the record was created and not closed  
576 (because patient moved or died).

577 **5.1.3.4. Active population size by age (C7.3.1)**

578 Where this information can be extracted, the number of unique active individuals split by age  
579 groups should be captured.

580 An active population for administrative healthcare data refers to the collection of patients for  
581 which there is an active record in the practice, i.e. the record was created and not closed  
582 (because patient moved or died).

583 **5.1.3.5. Median time (B6.3)**

584 The median time, in years, between first and last available records for unique individuals  
585 captured in the data source.

586 **5.1.3.6. Median time active (B6.3.1)**

587 The median time, in years, between first and last available records for unique **active** individuals  
588 (alive and currently registered) captured in the data source.

589 An active population for administrative healthcare data refers to the collection of patients for  
590 which there is an active record in the practice, i.e. the record was created and not closed  
591 (because patient moved or died).  
592

593 **5.1.4. Data source – Data flows and management**

594 **5.1.4.1. Governance details (C2.3)**

595 Description of the documents or links to webpages that describe the overall governance,  
596 processes and procedures for data capture and management, data access, data quality check  
597 and validation results, utilisation for research purposes.

598 **5.1.4.2. Follow-up (C2.13, C2.13.1, C2.7)**

599 If further follow-up would be needed, the availability of below access options should be specified:  
600 Accessing biospecimens: if this is possible (C.2.13) then also the biospecimen access conditions  
601 should be described (or a reference source can be added) (C2.13.1)  
602 Contacting patients or practitioners (C2.7)

603 **5.1.4.3. The process of collection and recording (C4.3)**

604 The process or manner in which recording of data in the data source occurs should be described;  
605 this could include the tools used, such as surveys, or a description of the system that the data  
606 holder uses to gather data and store it the data source.

607 **5.1.4.4. Record creation (C5.2)**

608 The event triggering the creation of a record in the data source should be described (e.g.:  
609 hospital discharge, specialist encounter, medicinal product dispensing).  
610 This refers in general to the creation of a record in the data source (and not to the registration  
611 of a person, see below).

612 **5.1.4.5. Registration of a person (C1.6, C1.6.1)**

613 The event triggering registration of a person in the data source should be selected from the  
614 following available values:  
615 - Birth  
616 - Immigration  
617 - Residency obtained  
618 - Start of insurance coverage  
619 - Disease diagnosis  
620 - Start of treatment  
621 - Practice registration  
622 Where none of the above values apply, the triggering event for a person to be registered in the  
623 data source should be described separately (C1.6.1).

624 **5.1.4.6. De-registration of a person (C1.7, C1.7.1)**

625 The event triggering de-registration of a person in the data source: The event triggering de-  
626 registration of a person in the data source should be selected from the following available values:  
627 - Death  
628 - Emigration  
629 - End of insurance coverage  
630 - Practice deregistration  
631 - Loss to follow up  
632 - End of treatment  
633 Where none of the above values apply, the triggering event for a person to be de-registered in  
634 the data source should be described separately (C1.7.1).

635 **5.1.4.7. Linkage (B5.2, B5.2.1, B5.3, B4.1)**

636 Where the data source is created by the linkage of other data sources, the elements of the  
637 linkage should be briefly captured as follows:

- 638  
639 - The linkage strategy (B5.2): whether the linkage is deterministic, probabilistic or a combination  
640 of the two.  
641 - The linkage variable used (B5.2.1) (e.g.: patient ID, date of birth etc.)  
642 - The completeness of the linkage (B5.3), described as a percentage along with the reference  
643 used  
644 - Names of the linked data sources (B4.1). Where these data sources are available in the data  
645 source catalogue, these should be cross-referenced.

646 **5.1.4.8. Data management specifications (C2.7, C8.5, C8.5.1, C2.9):**

647 The following information related to data management specifications should be selected, as  
648 applicable to the data source:

- 649 - Whether or not the data source allows data validation (e.g.: access to original medical charts)  
650 - If the records are preserved indefinitely (C8.5)  
651 - Where the records are not indefinitely preserved, the number of years for which the records  
652 are kept should be specified (C8.5.1)  
653 - Whether approval is needed for publishing results of a study using its data (C2.9)

654 **5.1.4.9. Informed consent for use of data for research (C2.5, C2.5.1)**

655 The need for informed consent in the context of research should be captured here. The type of  
656 informed consent could be categorised as:

- 657 - Not required  
658 - Required for general use of the data source  
659 - Required for all studies run on the data source  
660 - Required for intervention studies only  
661 - Waiver

662 Where the informed consent does not fit in the above categories, the value 'Other' can be used  
663 and further details should be provided (C2.5.1).

664 **5.1.4.10. Data source refresh (C8.2)**

665 Where the data source is refreshed on fixed dates around the year, this should be provided by  
666 selecting the month as applicable (e.g.: every June). The field can be repeated where the refresh  
667 happens more often than once a year (e.g.: every May and November).

668 **5.1.4.11. Data source last refresh (C8.3)**

669 Where the data source is refreshed at particular times throughout the year, the date when the  
670 last refresh of the data source occurred should be provided.

671 **5.1.4.12. CDM (Common Data Model) specifications (D1.2.1.1, D1.2, D1.2.1, D1.4, D1.7)**

672 The following data elements should be captured for data sources being transformed using a  
673 Common Data Model (CDM), (D1.2.1.1) as follows:

- 674 - The CDM name should be selected from the existing predefined list as follows: OMOP,  
675 ConcepTION, Nordic, Sentinel, PCORnet, VSD, i2b2, CDISC SDTM, PEDSnet (D1.2).

676 Where the common data model used is not listed in the values offered, further details should be  
677 provided (D1.2.1)

- 678 - The CDM website reference should be provided where available (D1.4)  
679 - The CDM release frequency, in number of months, should be provided (D1.7)

680 **5.1.4.13. Data source ETL to a CDM (B7.1, B7.5, B7.3, B7.4)**

681 Where applicable, further information on the data transformation (ETL) to a common data model  
682 (CDM) should be provided as follows:  
683 - The status of the transformation (ETL) of the data source should be described as either:  
684 planned, in progress or completed.  
685 - The frequency in months of the ETL frequency  
686 - The version(s) of CDM(s) to which the data source has been ETL-d  
687 - Data source ETL specifications: documents describing the mapping of the data source to the  
688 CDM (including codes and scripts to transform original data to CDM)

689 **5.1.5. Data source – Vocabularies and standardised dictionaries**

690 **5.1.5.1. Medicinal product information available**

691 The type of information captured with regards to the medicinal product should be selected from the  
692 values described in the table below.

Vocabulary	Description
Brand name	Specific name or trademark under which a medicine is sold
Batch number	The designation printed on the medicine label that allows the history of its production to be traced
Formulation	Pharmaceutical form of the medicinal product (e.g.: tablets, capsules etc.)
Strength	The amount of active ingredient contained in the medicinal product.
Package size	Number of individual formulations contained in a package (e.g.: 30 tablets per package)
Dose	The medicinal product dose prescribed or administered to the patient
Dosage regime	The schedule of doses of a medicinal product per unit of time (e.g.: every 6 hours)
Route of administration	The manner in which a medicinal product enters the body (e.g.: oral, intravenous)

693

694 **5.1.5.2. Medicinal product vocabulary used (C6.15.1)**

Vocabulary	Description
Art 57	Authorised medicines information in EU and EEA. Further reference <a href="#">here</a> .
IFA GmbH	Informationsstelle für Arzneispezialitäten. Further reference <a href="#">here</a> .

Vocabulary	Description
EDQM	European Directorate for the Quality of Medicines. Further reference <a href="#">here</a> .
SPN	Standard Product Nomenclature. Further reference <a href="#">here</a> .
MTHSPL	FDA Structured Product labelling. Further reference <a href="#">here</a> .

695

696 Where the medicinal product information is not coded (i.e.: provided as free text) this should be  
697 marked accordingly.

698 If other dictionaries than the listed ones are used, the value 'Other' should be used.

699 **5.1.5.3. Cause of death vocabulary**

Vocabulary	Description
ICPC	International Classification of Primary Care. Further reference <a href="#">here</a> .
ICD9	International Classification of Diseases, 9 <sup>th</sup> revision. External reference <a href="#">here</a> .
ICD10	International Classification of Diseases, 10 <sup>th</sup> revision. External reference <a href="#">here</a> .
ICD1	International Classification of Diseases, 1 <sup>st</sup> version. External reference <a href="#">here</a> .
Read	External reference <a href="#">here</a> .
SNOMED	Systematized Nomenclature of Medicine. Further reference <a href="#">here</a> .
SNOMED CT	Systemized Nomenclature of Medicine – Clinical Terms. Further reference <a href="#">here</a> .
MedDRA	Medical Dictionary for Regulatory Activities. Further reference <a href="#">here</a> .
OPCS	Classification of Interventions and Procedures. Further reference <a href="#">here</a> .

700

701 Where the cause of death is not coded (i.e.: provided as free text) this should be marked accordingly.

702 If other dictionaries than the listed ones are used, the value 'Other' should be used.

703 **5.1.5.4. Quality of life measurements (C6.28, C6.28.1)**

Vocabulary	Description
AQoL-8D	Assessment of Quality of Life 8-Dimension. Further reference <a href="#">here</a> .
QOLS	Quality of Life Scale. Further reference <a href="#">here</a>
MQOL	McGill Quality of Life Questionnaire. Further reference <a href="#">here</a> .
MQOL-E	The McGill Quality of Life Questionnaire – Expanded. Further reference <a href="#">here</a> .
HRQOL	Health-related quality of life. Further reference <a href="#">here</a> .
WHOQOL	World Health Organization External Measuring Quality of Life. Further reference <a href="#">here</a> .
EQ5D	Standardised measure of health-related quality of life developed by the EuroQol Group. EQ-5D assesses health status in terms of five dimensions of health. Further reference <a href="#">here</a> .
15D	The 15D is a generic 15-dimensional self-administered instrument for measuring HRQoL (Health-related quality of life). Further reference <a href="#">here</a> .
SF-36	The Short Form (36) Health Survey is a 36-item, patient-reported survey of patient health. Further reference <a href="#">here</a> .
SF-6D	An abbreviated variant of SF-36 commonly used in health economics as a variable in the quality-adjusted life year calculation to determine the cost-effectiveness of a health treatment. External reference <a href="#">here</a> .
HUI	Health Utilities Index. Further reference <a href="#">here</a> .

704  
705 Where the quality of life is captured but not coded (i.e.: provided as free text) this should be marked  
706 accordingly.

707 If other dictionaries than the listed ones are used, the value 'Other' should be used. In this case, the  
708 name of the 'quality of life' scale used should be provided in the free text field accordingly (C6.28.1).  
709

710 **5.1.5.5. Prescription vocabulary (C6.13.1)**

Vocabulary	Description
ATC	Anatomical Therapeutic Chemical code. Further reference <a href="#">here</a> .
RxNorm	A normalized naming system for generic and branded drugs. Further reference <a href="#">here</a> .
EphMRA	Anatomical Classification of Pharmaceutical Products maintained by EphMRA. Further reference <a href="#">here</a> .

Vocabulary	Description
DrugBank	DrugBank Online is a comprehensive, free-to-access, online database containing information on drugs and drug targets. Further reference <a href="#">here</a> .

711 Where the prescription is captured but not coded (i.e.: provided as free text) this should be marked  
712 accordingly.

713 If other dictionaries than the listed ones are used, the value 'Other' should be used.

714 **5.1.5.6. Dispensing vocabulary (C6.14.1)**

715 The dictionary used to code the dispensing information captured in the data source should be selected.

716 For further details on the available values, see section 5.1.5.5.

717 Where the dispensing information is captured but not coded (i.e.: provided as free text) this should be  
718 marked accordingly.

719 If other dictionaries than the listed ones are used, the value 'Other' should be used

720 **5.1.5.7. Indication vocabulary (C6.18.1, C6.18.2)**

721 The dictionary used to code the therapeutic indication captured in the data source should be selected.

722 For further details on the available values, see section 5.1.5.3.

723 Where the therapeutic indication is captured but not coded (i.e.: provided as free text) this should be  
724 marked accordingly.

725 If other dictionaries than the listed ones are used, the value 'Other' should be used. In this case, the  
726 name of the 'quality of life' scale used should be provided in the free text field accordingly (C6.18.2)

727 **5.1.5.8. Procedures vocabulary (C6.22)**

728 The dictionary used to code the procedures captured in the data source should be selected.

729 For further details on the available values, see section 5.1.5.3.

730 Where the procedure is captured but not coded (i.e.: provided as free text) this should be marked  
731 accordingly.

732 **5.1.5.9. Genetic data vocabulary (C6.25.1)**

Vocabulary	Description
OGG	A biological ontology in the area of genes and genomes. Further reference <a href="#">here</a> .
GO	Gene Ontology. Further reference <a href="#">here</a> .
EGO	Eukaryotic Gene Orthologues. Further reference <a href="#">here</a> .
SOPHARM	Suggested Ontology for Pharmacogenomics. Integrates OBO ontologies and formalizes specific gene variants. Further reference <a href="#">here</a> .

Vocabulary	Description
PHARE	PHarmacogenomic RElationships Ontology. Further reference <a href="#">here</a> .

733

734 Where the genetic data is captured but not coded (i.e.: provided as free text) this should be marked  
735 accordingly.

736 If other dictionaries than the listed ones are used, the value 'Other' should be used

737 **5.1.5.10. Biomarker data vocabulary (C6.26.1)**

Vocabulary	Description
SMASH	Semantic Mining of Activity, Social, and Health data. Further reference <a href="#">here</a> .
FOBI	Food-Biomarker Ontology. Further reference <a href="#">here</a> .

738 Where the biomarker data is captured but not coded (i.e.: provided as free text) this should be marked  
739 accordingly.

740 If other dictionaries than the listed ones are used, the value 'Other' should be used.

741 **5.1.5.11. Diagnosis/ medical event vocabulary (C6.9.1)**

742 The dictionary used to code the diagnosis, or any other medical event captured in the data source should  
743 be selected.

744 For further details on the available values, see section 5.1.5.4.

745 Where the diagnosis or medical event is captured but not coded (i.e.: provided as free text) this should  
746 be marked accordingly.

747 If other dictionaries than the listed ones are used, the value 'Other' should be used.

748 **6. Registering a data source in the Data source catalogue**

749 A Data holder would be able to request to register on a voluntary basis a data source in the Data source  
750 catalogue via a dedicated webform to be made available in the second half of 2023. An e-mail address  
751 supporting this process is available: [metadata@ema.europa.eu](mailto:metadata@ema.europa.eu).

752 Additionally, EMA is proactively contacting data holders requesting the addition of the metadata  
753 information in the catalogue, looking to current data sources registered in ENCePP Resources Database.

754 **7. Maintenance of information in the Data source catalogue**

755 It is important that the metadata information is kept up-to-date; this refresh of information is expected  
756 to be run on a yearly basis or more often for particular data sources if found necessary.

757 The data holder will be provided with the technical means to update the information provided directly,  
758 via a dedicated webform. This will be made available in the second half of 2023.

## 759 **References**

- 760 Brown, P., Brunnhuber, K., Chalkidou, K., Chalmers, I., Clarke, M., Fenton, M., Forbes, C., Glanville, J.,  
761 Hicks, N. J., Moody, J., Twaddle, S., Timimi, H., & Young, P. How to formulate research  
762 recommendations. *BMJ*. 2006;333(7572), 804–806. <https://doi.org/10.1136/bmj.38987.492014.94>
- 763 ENCePP Guide on Methodological Standards in Pharmacoepidemiology, 10<sup>th</sup> Rev. (2022). [Chapter 12.1](#)  
764 [General principles of quality management](#)
- 765 ENCePP Guide on Methodological Standards in Pharmacoepidemiology, 10<sup>th</sup> Rev. (2022). [Chapter 12.2.](#)  
766 [Data Quality Frameworks](#).
- 767 FAIR Principles. <https://www.go-fair.org/fair-principles/>
- 768 HMA/EMA. [List of metadata for Real World Data catalogues](#) (2022).
- 769 Gatto, N. M., Campbell, U. B., Rubinstein, E., Jaksa, A., Mattox, P., Mo, J., & Reynolds, R. F. (2022).  
770 The Structured Process to Identify Fit-For-Purpose Data: A Data Feasibility Assessment Framework.  
771 *Clin Pharmacol Ther*. 2022;111(1), 122–134. <https://doi.org/10.1002/cpt.2466>
- 772 MINERVA: Strengthening Use of Real-World Data in Medicines Development: Metadata for Data  
773 Discoverability and Study Replicability (2022). [EUPAS39322](#)
- 774 TEHDAS. [European Health Data Space Data Quality Framework](#) (2022).
- 775 Wang S, Pottegard A, Crown W et al. HARmonized Protocol Template to Enhance Reproducibility  
776 (HARPER) of Hypothesis Evaluating Real-World Evidence Studies on Treatment Effects: A Good  
777 Practices Report of a Joint ISPE/ISPOR Task Force. *Pharmacol Drug Saf*. 2022;
- 778 Wang S., Schneeweiss S. [Assessing and Interpreting Real-World Evidence Studies: Introductory Points](#)  
779 [for New Reviewers](#). *Clin Pharmacol Ther*. 2022;111(1):145-149.