



1 4 November 2024
2 EMA/503781/2024
3 Committee for Medicinal Products for Human Use (CHMP)

4 **Data Quality Framework for EU medicines regulation:**
5 **application to Real-World Data**
6 **Draft**

7

Draft agreed by Methodology Working Party (MWP)	September 2024
Adopted by CHMP for release for consultation	4 November 2024
Start of public consultation	29 November 2024
End of consultation (deadline for comments)	31 January 2025

8

Comments should be provided using this [EUSurvey form](#). For any technical issues, please contact the [EUSurvey Support](#).

9

Keywords	<i>Data quality, framework, real-world data, real-world evidence, use of data, primary, metadata, reliability, extensiveness, coherence, timeliness, relevance, maturity models, validation</i>
-----------------	--

10



11 Data Quality Framework for EU medicines regulation:
12 application to Real-World Data

13 **Table of contents**

14 **Executive summary 3**

15 **1. Background: Real-World Data and Data Quality 4**

16 1.1. Definition of Real-World Data4

17 1.2. Distinctive traits of RWD4

18 1.3. RWD use-based quality control.....5

19 1.4. Impact of secondary use of RWD on data quality5

20 1.4.1. Impact on Reliability6

21 1.4.2. Impact on Extensiveness and Representativeness.....6

22 1.4.3. Impact on Coherence7

23 1.5. Responsibility for DQ in RWD7

24 **2. Application of the EMRN DQF to RWD 7**

25 2.1. Purpose of this document7

26 2.2. Scope of the RW-DQF8

27 2.3. Structure of the RW-DQF.....8

28 2.3.1. Understanding relevance9

29 **3. Guidelines for the characterisation of systems and processes**

30 **underpinning data 10**

31 3.1. Systems and process characterisation checklist and maturity model 10

32 3.2. General considerations on the characterisation of systems and processes..... 17

33 **4. Data quality metrics for RWD 18**

34 4.1. Framework for the categorisation and identification of metrics 18

35 4.2. Metrics for DQ assessments..... 20

36 4.3. Considerations for the implementation of RWD DQ metrics 23

37 4.3.1. Different roles of metrics 23

38 4.3.2. Additional considerations on level of application and maturity for metric assessments

39 24

40 **5. Guidelines to assess quality in relation to a specific research question. 24**

41 5.1. General principles for assessment of data quality in relation to a research question 24

42 5.2. Framework for detailed fitness-for-use assessment..... 27

43 5.3. Illustrative example for detailed fitness-for-use assessment 28

44 5.4. Toward a generalisation of question-specific aspects..... 31

45 5.5. Providing supporting information for RWD in regulatory submissions 32

46 **6. Concluding remarks..... 32**

47 **7. References 32**

48 **Definitions..... 34**

49 **Glossary 34**

50

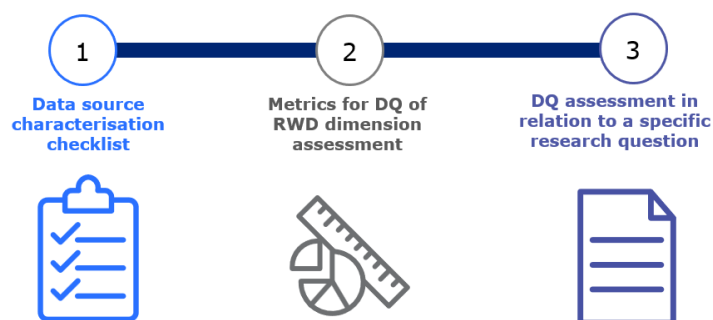
51 **Executive summary**

52 This document describes the Real-World Data (RWD) specific recommendations as derived from the
53 Data Quality Framework (DQF) for EU Medicines regulation endorsed by the Committee for Medicinal
54 Products for Human Use (CHMP) [1] (hereafter referred to as EMRN DQF). The EMRN DQF sets out the
55 principles, concepts, and definitions as intended to be applied widely across datasets used in medicine
56 regulatory use cases. It also provides examples and in-depth clarifications on the developed framework
57 elements for characterising, assessing, and assuring DQ in the regulatory context. It is therefore
58 recommended to use the EMRN DQF as a companion document when reading this chapter.

59 The application of the EMRN DQF to RWD (hereafter referred to as RW-DQF) sets out the specificities of
60 RWD and enable regulators to evaluate the quality of data underpinning Real-World Evidence (RWE) as
61 used in the regulatory assessment. It also provides guidance on the relevance assessment of such data
62 to a research question based on DQ metrics and evidence of systems and processes underpinning data.
63 These parts provide actionable and focused recommendations for assessing DQ of RWD, with the goal
64 of improving the usefulness of RWE for regulatory purposes. The RW-DQF is intended for the use of
65 stakeholders involved in regulatory processes, primarily aimed at members of the European Medicines
66 Regulatory Network (EMRN), but also other actors involved in this process, such as the Data Analysis
67 and Real-World Interrogation Network (DARWIN EU®), pharmaceutical industry, academia, contract
68 research organisations, and data holders.

69 With a view of maintaining the consistency with parallel activities ongoing in the context of European
70 Health Data Space (EHDS), the document was developed in close collaboration with the Towards
71 European Health Data Space (TEHDAS) and QUANTUM (The Health Data Quality label) projects. These
72 initiatives aim to address the wider use of health data, whereas the RW-DQF specifically focuses on the
73 challenges faced when using this data within the medicine regulation assessment.

74 The topics addressed in this document are: an introduction to RWD key considerations on quality
75 (Chapters 3 & 4), practical recommendations on characterisation of the systems and processes that
76 underpin data (Chapter 5), a set of metrics to assess data quality (DQ) dimensions (Chapter 6), and a
77 guideline on how to assess DQ in relation to a research question via the use of a framework and an
78 illustrative example (Chapter 7) (Figure 1).



79

80 **Figure 1- Representation of the key points of the DQF for EU medicines regulation:**
81 **application to RWD.**

82

83 **1. Background: Real-World Data and Data Quality**

84 **1.1. Definition of Real-World Data**

85 In the context of RWE studies, Real-World Data (RWD) are data that describe patient characteristics
86 (including treatment utilisation and outcomes) in routine clinical practice [2, 3, 4]. In broader terms,
87 RWD represent data captured in routine care which are not collected in a clinical trial and are relevant
88 to the subject (e.g., age, sex, ethnicity etc.), the disease, the treatment, interactions with the
89 healthcare system, as well as social and environmental factors influencing health status. RWD may
90 originate from primary data collection (primary use of data), i.e., data collected specifically for the
91 study in question, or secondary use of data initially recorded in the context of different primary
92 purposes (such as the clinical management of patients or for administrative reasons). The secondary
93 use of data, driven by specific research objectives, can involve a single or multiple RWD sources (i.e.,
94 multi-database studies). For example, the assessment of treatment pathway disparities in a given
95 region and for a given indication could entail combining and analysing RWD from multiple sources.

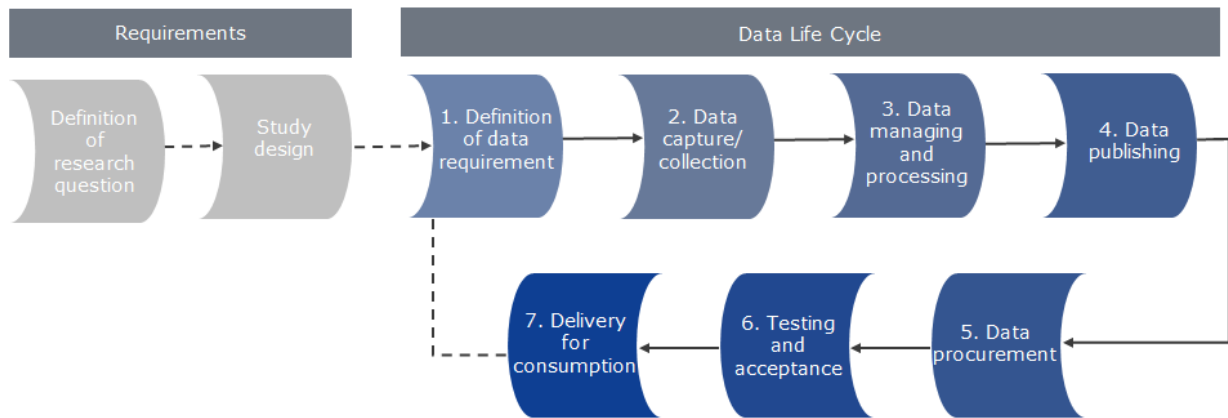
96 Through the analysis of RWD, Real-World Evidence (RWE) is generated to answer research questions.
97 The main areas where RWD analyses can aid medicines regulatory use cases across a medicinal
98 product's lifecycle are to [2]: 1) support planning and validity of applicant studies (e.g., comparison of
99 a study population with patients from the real-world setting to ensure representativeness of the clinical
100 study, patient recruitment), 2) understand the clinical context (e.g., disease epidemiology, disease
101 prevalence/incidence, description of drug utilisation patterns including switching and off-label use), 3)
102 investigate associations and impact (e.g., medicines post-marketing surveillance, assessment of
103 effectiveness of risk minimisation measures) [5].

104 **1.2. Distinctive traits of RWD**

105 RWD have several distinctive traits, including the variety of sources capturing them, their structure,
106 format, variables collected, terminologies used, processes related to data collection or data recording.
107 RWD can be leveraged for their primary intended purpose or for secondary use [6, 7]. Types of RWD
108 sources include, for example, electronic healthcare databases (e.g., from primary care, specialist care,
109 and hospital care settings, claims databases), longitudinal drug prescription, dispensing or other drug
110 utilisation data or patient registries. The latter is defined as a system which records uniform data
111 (clinical and other) to identify specified outcomes for a population defined by a particular disease,
112 condition or exposure [13]. Data captured in a registry can involve data collection with the primary
113 purpose of generating RWE or secondary use of data. In addition, RWD can be collected via
114 compassionate use programmes on products in development for patients who cannot enter clinical
115 trials, as they can facilitate the understanding of the best conditions of medication use and which
116 patients can be benefited the most [22].

117 The RWD lifecycle may involve several manipulations with or without changing hands between
118 organisations, e.g., aggregation, transformation, cleaning, metadata creation, publishing (see Figure
119 1). In cases of secondary use of RWD, the data may not be entirely fit for the study at hand, as their
120 primary purpose might differ.

121 The metadata of RWD for RWE generation is typically included in data catalogues for dissemination.
122 These "research-ready" data are usually optimised for quality and usability as much as possible without
123 referring to a specific research question.



124

125 **Figure 1 Data life cycle for secondary use of RWD. The definition of research**
 126 **question and the study design are required steps preceding the definition of data**
 127 **requirement (Step 1).**

128 **1.3. RWD use-based quality control**

129 The primary purpose of recording RWD is the provision of health services to assess, maintain or
 130 restore the state of health of the person that the data belong to, and in some cases, RWE generation.
 131 The quality of RWD capture should be under the control of a Quality Management System (QMS) as
 132 explained in more detail in the EMRN Data Quality Framework (DQF) [1].

133 When RWD are leveraged for RWE (either when that is their primary or secondary use) and for other
 134 secondary purposes, they must be considered as uncontrolled, i.e., not produced through a defined
 135 process with feedback loops to detect and correct errors and prevent their future occurrence. A QMS
 136 such as the ISO 9000 family, Good Clinical Practices, Good Laboratory Practices or Good Manufacturing
 137 Practice can only be expected to assure quality with respect to the primary use of RWD.

138 **1.4. Impact of secondary use of RWD on data quality**

139 Secondary use of RWD has a significant impact on data quality (DQ) for the reasons outlined here:

140 **Decoupling of data collection and purpose:** Data for secondary use are, by definition, originally
 141 recorded for other purposes than generation of RWE, and the secondary purpose consists of many
 142 different individual research questions not always known at the time of capturing and publishing
 143 research-ready data. Fitness-for-use depends on a defined research question, and therefore cannot be
 144 controlled or imposed to a source at the time of data recording. In other words, DQ can be measured
 145 at source, but cannot be fully assessed for adequacy at the time of data collection or data recording.

146 **Anonymisation and pseudonymisation:** Sensitive patient-level healthcare data need to be
 147 anonymised or pseudonymised to facilitate data protection, publication, and re-use by different
 148 researchers. This process sometimes involves masking, replacing or removing information.
 149 Consequently, quality and operational aspects of quality control may be negatively affected given that
 150 people with access to identifiable information are usually tightly restricted and may not include staff
 151 involved in quality determination.

152 **Data linkage:** RWD captured from a single source may not systematically provide a comprehensive
 153 view over the whole lifetime of the patient. Data linkage can be used to address this problem by
 154 combining RWD from several individual data sources. Linkage methodologies such as *person matching*
 155 (after pseudonymisation) and *de-duplication of records* are often used but may create quality issues.
 156 For example, probability-based matching of patients based on non-identifiable or incomplete

157 information can combine two different patient records, resulting in inaccurate information. De-
158 duplication may have the effect of introducing incoherence, since both duplicate records are valid in
159 each of the sources. On the other hand, identifier-based techniques such as deterministic and
160 probabilistic linkage are often adequate [7].

161 Additionally, RWD originating from multiple RWD sources are often highly disparate, with structure and
162 terminology that are not standardised. Given this variety, the use of Common Data Model (CDMs)
163 plays an important role in RWD, facilitating the systematic implementation of some aspects of quality
164 control and the development of consistent methodologies (e.g.: data cleaning, profiling, reporting,
165 analysis) applied to different sources. The process to transform an original source to a CDM, called ETL
166 (extract, transform, load) can improve coherence, but at the same time have a negative impact on
167 other dimensions of DQ. For instance, reliability can be affected, both because any transformation
168 increases the risks of error (accuracy) and because the CDM will define some level of precision that
169 may be lower to that of the original source. Extensiveness can also be affected if for instance, some
170 data are removed as non-conformant to the target model, as well as timeliness, which is also affected
171 as the transformation process introduces delays.

172 **1.4.1. Impact on Reliability**

173 In a secondary use of data scenario, there is limited possibility to control most DQ factors of reliability
174 (e.g., accuracy and precision) at the source (point of data recording). Therefore, the primary focus of
175 DQF implementation is error detection that could lead to record removal or amendment with
176 approximated values, while only in some limited cases it can lead to the correction of the data capture
177 processes.

178 For RWD datasets that are in the category of “Big Data” as defined by the HMA/EMA joint Big Data
179 Steering Group [9], error detection cannot be practically achieved through manual checking of all data
180 records. Therefore, automation plays a key role in error detection and can sometimes help identify
181 inconsistencies or outliers in the data. The use of common data models (CDM)¹ and standardised
182 analytics can prevent coding errors and differences in data curation [10]. However, if these errors are
183 missed during the conversion process, they will remain in the converted data and may be hidden if
184 conversion back to the source data is not possible. Plausibility metrics also play a key role in assessing
185 the reliability of RWD, as it is possible to provide an alternative to validation against primary institution
186 source records.

187 **1.4.2. Impact on Extensiveness and Representativeness**

188 In a secondary use of RWD scenario, it is possible to measure the amount of information in a dataset
189 (e.g.: measuring completeness), but it can be challenging to characterise how this relates to the data
190 recording process. For instance, DQ control on secondary use of data often cannot adequately detect
191 missing information, especially when this relates to an outcome or event that is not necessarily
192 expected to be present². This is made even more challenging as missing data can be the result of
193 flawed data transfer rather than incomplete RWD capture.

194 In addition, RWD tends to be interpreted following a closed-world assumption that means that an
195 outcome or event is assumed as non-present because it has not been recorded³ (an example is a
196 patient with cardiovascular disease and type II diabetes, who is considered however non-diabetic in
197 their medical records because they were never tested for the disease).

¹ The harmonisation to a common data model itself has an impact on precision and potentially accuracy.

² In some cases (e.g.: age), where the information is known to necessarily exist, the lack of such information can be clearly detected as missing information.

³ This is unlike in clinical trial data, where absence of events is explicit, and there is no concept of missing values.

198 These presumptions, which are rarely made explicit, have implications for the analytical methods of
199 generating RWE and are fundamental for the assessment of DQ [11].

200 When it comes to secondary use of RWD, lack of representativeness of the target population for the
201 study objective may lead to biased outcomes in some types of studies. For example, a study on
202 disease prevalence using RWD from a primary care database may produce skewed results if individuals
203 in the data source are not representative of the entire population.

204 **1.4.3. Impact on Coherence**

205 RWD are often recorded from different healthcare actors and is varied both due to different data
206 representation, i.e., format, structure, content, etc., and due to differing processes for data collection
207 or data recording and DQ control. Coherence, which refers to the homogeneity/uniformity and
208 consistency of data within a single source or across multiple RWD sources, is a critical aspect that
209 needs to be assessed at the time of data publication or consumption. Additionally, coherence is
210 essential to be re-assessed whenever new RWD sources or elements are introduced, especially in the
211 case of data linkage, to ensure data integrity.

212 Though assessment of several aspects of coherence can be facilitated by measuring conformance vs a
213 specific (common) target model (e.g.: format coherence or structural coherence), some aspects are
214 more subtle:

- 215 • Semantic coherence may vary as diverse sources adopt different approaches to map between
216 terminologies. For instance, the term “anuria” can describe a condition of total cessation of
217 urine production in one source, while the same term in another source can be used to
218 specifically note instances where the measurement of urine output is below a specific
219 threshold. The mapping strategy of each source to a target model, coupled with the limitations
220 of terminologies to fully capture the semantic meaning of a mapped term, can lead to
221 coherence issues across diverse sources.
- 222 • Temporal coherence can be an issue for long-term datasets, as medical practices (and
223 therefore the meaning of data) may change along the data recording timeline.

224 **1.5. Responsibility for DQ in RWD**

225 Given the variety and complexity of the processes related to RWD recording and utilisation, the variety
226 of actors and data processing involved, DQ in RWD is a distributed responsibility. The responsibility to
227 ensure that DQ is properly characterised is divided among all actors in the RWD life cycle. This includes
228 the measures by which any processes involved in the various steps of the data life cycle (e.g., data
229 capture, aggregation, processing) can impact DQ. To allow an efficient DQ assessment, each party is
230 responsible for making evidence on DQ available when suitable or required, as well as to maintain DQ
231 within declared or acceptable standards, while documenting the processes followed and the data tools
232 used. More detailed information regarding DQ responsibility in RWD can be found under subsection
233 5.1. (Systems and process characterisation checklist and maturity model).

234 **2. Application of the EMRN DQF to RWD**

235 **2.1. Purpose of this document**

236 This document, hereafter referred to as the RW-DQF, extends the EMRN DQF to provide more
237 actionable and focused recommendations for assessing the DQ of RWD with the goal of improving the
238 quality of RWD for regulatory use. It also serves as a guide for enhancing the assessment and

239 documentation of RWD quality by providing guidelines for characterising the systems and processes
240 underpinning data and their impact, key metrics to evaluate different aspects of DQ within a dataset,
241 and guidelines for using these metrics to assess the suitability of a dataset through a fitness-for-use
242 assessment in relation to a specific research question.

243 **2.2. Scope of the RW-DQF**

244 This document focuses on the subset of RWD recorded within routine clinical practice, i.e.,
245 administrative or claims data, EHRs, pharmacy/prescription data, patient registries etc., when used in
246 the context of a specific research question and in line with the European Network of Centres for
247 Pharmacoepidemiology and Pharmacovigilance (ENCePP) guide on Methodological Standards in
248 Pharmacoepidemiology [12]. Existing guidance specific to particular data types (e.g., guidance on
249 registry-based studies [13] for the use of patient registry data) still applies and the RW-DQF should be
250 read in conjunction with any other relevant guidance documents.

251 The following should be considered out of scope:

- 252 • DQ concerning RWD arising from repurposing of previously published analyses, e.g., meta-
253 analyses etc.
- 254 • DQ of direct-from-patient data, i.e., PROs, patient engagement data, patient preferences,
255 mobile health data, social media data, etc., as these may have peculiarities in terms of
256 measuring DQ and would be subject to further guidance.

257 **2.3. Structure of the RW-DQF**

258 The RWD chapter is composed of three parts, presented in the three following sections:

- 259 • Guidelines to characterise systems and processes (and their impact)
- 260 • Metrics to appraise different aspects of DQ within a given dataset
- 261 • Guidelines to assess the suitability of a dataset for answering a specific research question.

262 The RW-DQF inherits the concepts and design of the EMRN DQF. This includes the categorisation of
263 quality aspects into three types of determinants (foundational, intrinsic and question-specific), the
264 maturity model and the definitions of the DQ dimensions. Such concepts are presented in this
265 document using a varied terminology, that is more commonly understood among RWD stakeholders.
266 They are also further specialised and altered to the RWD context.

267 In particular:

- 268 • **Foundational determinants** are defined as “everything that impacts DQ, but it is not related
269 directly to the dataset and does not depend on any specific research question”. In this chapter,
270 foundational DQ aspects are referred to as the characterisation of the **systems and processes**
271 that have an impact on DQ. In most establishments where RWD are captured, processed and/or
272 consumed, information on foundational determinants is often limited to onboarding documents [2].
273 In addition, the impact of systems and processes is usually considered with respect only to data
274 accrual [7] and then approximated [14]. This chapter considers the impact of systems and
275 processes in some detail, as well as the impact along the whole evidence generation process.
- 276 • **Intrinsic determinants** are defined as “DQ aspects that can be observed only on the basis of a
277 given dataset, without requirement for information about how the data were captured, or about its
278 primary/intended use”. In this chapter, intrinsic determinants are considered as **Metrics** that can
279 be used to characterise DQ, and the chapter provides guidelines on how to use such metrics.

280 • **Question-specific determinants** are defined as “aspects of DQ that cannot be assessed
 281 independently of a research question”. In this chapter, they are considered in terms of how to
 282 assess the **suitability of a dataset for a specific research question**.

283 As for dimensions and sub-dimensions, the terminology introduced in the EMRN DQF sometimes differs
 284 from what is found in other DQFs focused on RWD [15-17]. In general, there is a lack of consensus on
 285 terminology, among RWD practitioners (researchers, data analysts, data custodians) and more broadly
 286 among people involved in the RWD recording process. For instance, the term “validation” refers to
 287 checking whether the data correspond to the source [15], but is commonly understood as “checking
 288 that the data conform to a schema” by database administrators.

289 This RW-DQF will use the terminology introduced in the EMRN DQF [1]. These are reported in the
 290 glossary of this document (Chapter 9) for convenience.

291 **2.3.1. Understanding relevance**

292 In the EMRN DQF, the notion of relevance (i.e., the extent to which a dataset presents the data
 293 elements useful to answer a given research question) is considered as something cross-cutting through
 294 all DQ dimensions and applying to each of them. All aspects of DQ can be measured by metrics
 295 independently of a research question, while no thresholds or acceptance criteria can be established
 296 independently from it⁴.

297 **Error! Reference source not found.** summarises the interplay between systems and processes
 298 characterisation, metrics, and suitability to a research question across five dimensions for DQ. It can
 299 be seen in this table that “relevance” is the only dimension purely determined by the research
 300 question.

301 **Table 1- Interplay between definitions and dimensions for DQ**

	Are data correct? <i>Reliability</i>	Are data enough? <i>Extensiveness</i>	Are data homogenous? <i>Coherence</i>	Are data timely? <i>Timeliness</i>	Is this the right type of data? <i>Relevance</i>
Systems and processes	Determine reliability	Determine extensiveness	Enable coherence	Determines timeliness	
Metrics	Can assess reliability	Can measure extensiveness	Can measure and improve coherence		
Suitability to a research question	Defines “acceptable” reliability	Defines if data are sufficient	Defines if the level of coherence is adequate	Defines acceptable timeliness	Defines if the content of the data is what is needed

302 A major difference between the EMRN DQF and other proposed DQFs [15-17] for RWD is the role of
 303 relevance. In other frameworks, relevance is often considered as its own dimension, and includes
 304 completeness and reliability as sub-dimensions [15]. The notion of “Relevance DQ dimension”
 305 introduced in the EMRN DQF is much more restricted, and only meant to capture some DQ aspects not
 306 covered by other dimensions (corresponding to the question: is the type of data fit-for-use?).

⁴ As discussed in the EMRN DQF, there may be “general questions” that a dataset may be expected to be used for, and from which some quality threshold could be derived. However, establishing such threshold is easily discretionary without a clear definition of such target uses. Even in this case, an “unqualifying” dataset may still be useful in a different use case, e.g.: if data are very scarce and critical.

307 In the context of RWD, reference to the Relevance DQ Dimension is omitted in the discussion around
308 metrics but is considered in the context of assessing the suitability of datasets used to answer a
309 specific research question.

310 **3. Guidelines for the characterisation of systems and** 311 **processes underpinning data**

312 The quality of data cannot be assured unless the systems and processes responsible for their collection
313 or recording and transformation are reliable and offer the necessary guarantees. If RWD are
314 considered for regulatory use, no matter the content of an RWD source, its use would be unfeasible
315 unless there is some reasonable evidence that the information provided is true and not accidentally or
316 intentionally altered. This section provides guidelines on how to characterise systems and processes,
317 so that their effect on DQ can be assessed.

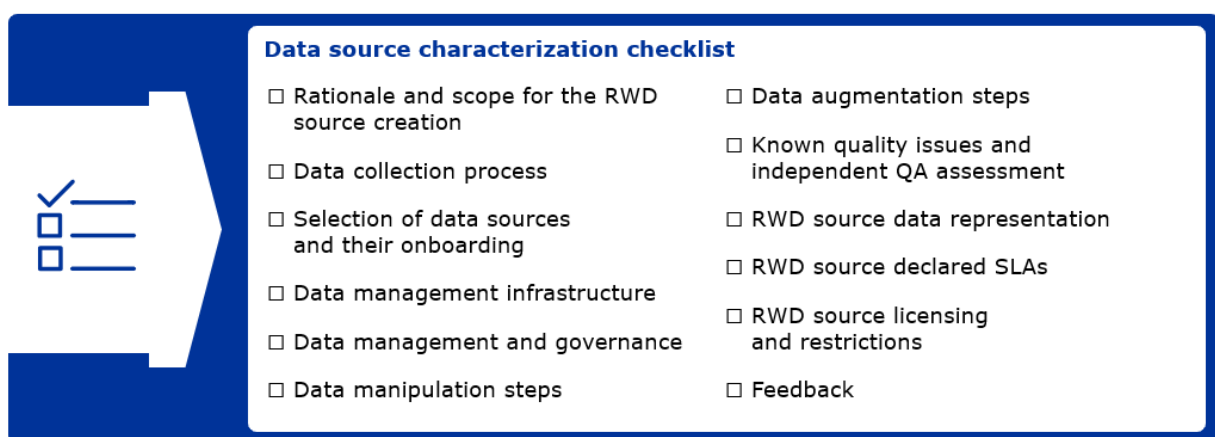
318 This section builds on the “Foundational Determinants” definition of the EMRN DQF [1], and the related
319 maturity model. In the RWD context, the maturity model is adapted and takes the form of a practical
320 checklist (Figure 3).

321 At its core, maturity depends on the ability to produce documented evidence of good DQ practices and
322 quality-related actions. Maturity advances if this documented evidence is standardised and systematic
323 by nature, to allow it to be coherently interpreted across different RWD sources.

324 Maturity depends on the level of automation, as automated DQ processes can be both more extensive
325 and less subject to accidental or human errors. Therefore, the maturity model presented here focuses
326 on what information should be provided, how it could be standardised, and how it could be automated.

327 **3.1. Systems and process characterisation checklist and maturity model**

328 This section suggests how to characterise aspects of systems and processes that have an impact on
329 DQ. These aspects are grouped by areas that reflect different steps found in the data life cycle (see
330 Figure 2) and are presented in Table 2 that can be used as checklist to verify if foundational
331 information necessary to assess DQ is provided.



332 **Figure 2 - RWD source characterisation checklist overview. SLA – service level**
333 **agreement.**
334

335 For each area of interest, the table lists what information should be provided together with how this
336 information should generally be generated, following the maturity model derived from the EMRN DQF
337 [1]⁵:

338 **Level 1: Documented.** Some information should be provided at least as simple documentation (e.g.:
339 short text) and/or supporting links. More extensive documentation can include standard operating
340 procedures (SOPs), as well as key performance indicators (KPIs).

341 **Level 2: Formalised.** The information should be provided in a way that follows established or
342 emergent standards (e.g.: following recommendations such as the ones in the EMA Guideline on
343 registry-based studies [13] or the REQuEST tool [18], or more general frameworks and standards
344 e.g.,: [19]. SOPs and KPIs, when reported, follow established standards and guidelines.

345 **Level 3: Automated.** The information should be derived in a way that guarantees higher DQ by
346 design. This means the data are generated by systems and platforms by a computed process rather
347 than being entered ad-hoc or a-posteriori. For example, in the case of data lineage, provenance
348 information is generated by an ETL engine or derived from some executable electronic specification
349 instead of being provided independently from the system by manual documentation.

350 These three maturity levels are reflected in table 2, where the “Documented” column lists what
351 information should be provided. The “Formalised” column provides suggestion of how data can be
352 presented in a standardised manner. The “Automated” column provides suggestions of ways to
353 improve and verify data captured through automated processes.

354 Not all required information is relevant for all workflows (e.g.: a registry implies different processes
355 than a source aggregating and repurposing claims or prescription data). For pragmatic reasons,
356 alongside information describing systems and processes in place, descriptions of the intended
357 characteristics of RWD are also included in Table 2⁶.

358 Achieving higher levels of DQ is a distributed responsibility across the data lifecycle from data capture
359 to data use. Essential DQ checks cannot be fully performed solely by a single stakeholder as for
360 instance the responsible person for a study submission. In addition, data characterisation such as data
361 pertaining to the data’s fitness-for-use, cannot be anticipated by the stakeholder capturing the data or
362 from the data holder alone, however, the data holder has usually knowledge whether a research
363 question can be answered. Therefore, Table 2 proposes to refer to “upstream” quality assessments for
364 all RWD sources under the responsibility of RWD holders. To allow an adequate DQ assessment, the
365 RWD submitters or the RWD end users, i.e., the stakeholders that get access to RWD and use the RWD
366 for secondary purpose, should be responsible for providing the available evidence on DQ when suitable
367 or required.

368 Finally, this table proposes to capture information on what feedback mechanisms are provided. The
369 concept of feedback mechanisms is derived from the EMRN DQF overall maturity model. Given the
370 complexity and heterogeneity of RWD capture, detailed feedback loops on all aspects of DQ may be
371 unfeasible. Therefore, feedback is here proposed as an explicit aspect to report for the overall RWD
372 source, rather than as an extra maturity level.

373

⁵ The maturity model in the RWD Chapter differ from the EMRN DQF in that “feedback” is not considered as a maturity level, but as an element of systems and processes.

⁶ Sections IX, X, XI describe RWD source characteristics pragmatically useful in a “data checklist”: even if not strictly affecting DQ, they are relevant for actions supporting DQ assessment.

374 **Table 2: RWD Source characterisation checklist with 12 areas of interest**

		Level 1: Documented	Level 2: Formalised	Level 3: Automated
Item	Rationale	Documentation to be provided: free text and/or online link(s)	Suggestions for standardised documentation to enhance interpretability	Suggestions for robust documentation generation “by design” and machine readable
I. Rationale and scope for the RWD source creation	Relevant for all DQ dimensions as it provides a general understanding of the strengths and limitations of an RWD source.	<p>A) The primary purpose(s) for which data are collected</p> <p>B) The justification or criteria used for the selection of the data being collected (or integrated)</p> <p>C) Publications describing this RWD source</p>	Provide information using standardised / widely used templates, to make the relevant information easy to digest and interpret ⁷	Provide information as Metadata , in a standard format and with clear definitions ⁸ , to allow Metadata to be automatically processed and their quality (e.g.: completeness) be adequately verified.
II. The data collection or recording process	<p>Essential to understand coverage and to assess reliability (that can be affected by errors or biases in the collection process).</p> <p>Also, essential to evaluate SOP for data collection or recording practices that may impact coherence (e.g., where “curation at source” is involved and provide hard constraints for timeliness).</p>	<p>A) Description of the data provider, including:</p> <ul style="list-style-type: none"> its nature (patients self-reported, carers or third parties, healthcare professionals with specified speciality) its geographical and organisational setting. <p>B) Description of data collection or recording SOPs, including the rationale for the SOP design.</p> <p>C) Information of how SOPs are implemented, and their execution monitored.</p> <p>D) Characteristics of key data elements captured, e.g.:</p> <ul style="list-style-type: none"> core, optional elements planned size planned coverage over time 	<p>In addition to I: make use of standard vocabularies where available.</p> <p>SOPs are made available and are based on common shared standards.</p>	<p>In addition to I:</p> <p>SOPs specify KPIs so that adherence to these KPIs can be monitored and reported.</p>

⁷ An example of a such a template can be found at [REQuest]. We encourage the development of such shared template, when not available.

⁸ In essence, information provided should follow FAIR standards, when feasible.

<p>III. The selection of RWD sources and their onboarding</p> <p>(Applies to RWD sources that integrate or repurpose other RWD sources)</p>	<p>When data are provided by a data aggregator, ensure that all the available evidence related to systems and processes potentially affecting DQ can be followed. Provide information of impact on both reliability and evidence (as well as other dimensions if relative constraints are formulated in inclusion/exclusion (I/E) criteria)</p>	<p>A) The data providers' selection processes and criteria, e.g.:</p> <ul style="list-style-type: none"> Inclusion and exclusion criteria for the acceptance of a RWD source <p>B) A comprehensive DQ assessment of RWD sources being consumed (as a reference, or as evidence of the frameworks being followed)</p>	<p>The DQF assessment includes this checklist for each RWD source.</p> <p>Reference to datasets with unambiguous identifiers that can distinguish between datasets versions, when relevant.</p>	<p>(Aspirational)</p> <p>The DQ assessment provided can be processed so that upstream DQ can be incorporated, to reconstruct a full "chain of evidence" of an RWD resource.</p>
<p>IV. The data management infrastructure</p>	<p>Essential for reliability regarding data alterations resulting from system accidents, software errors or malicious intervention.</p>	<p>A) The list of systems used to manage the RWD source, from data collection or recording to processing to making it available (version, features used).</p> <p>B) The software⁹ testing and QA processes in place.</p> <p>C) Measures to prevent accidental physical data alterations (e.g.: backups, redundant systems, checksums).</p>	<p>The hardware or software implementation complies with recognised quality standards that can be reported.</p>	
<p>V. Data management and governance</p>	<p>Data management and governance impact reliability, as well as all quality</p>	<p>A) A description of the overall data management principles adopted (e.g.: ALCOA+, FAIR)</p>	<p>SOPs and data management processes adhere to standards that can be referred to: e.g., GCP, ENCePP, ISO 25012, ISO 25101, ISO 8000-</p>	<p>Data management and governance is implemented in the data platforms 'Digital Quality Measures' (DQMs) so that reports of performance and</p>

⁹ We assume the HW testing is not an issue as necessarily performed a-priori.

	dimensions for metadata.	<p>B) A description of data management processes in place:</p> <ul style="list-style-type: none"> * SOPs in place * Responsibilities and roles * DQ controls * KPIs <p>C) Measures to prevent unauthorised data alterations (e.g.: cybersecurity approach)</p> <p>D) Monitoring, auditing, and quality improvement procedures in place.</p> <p>E) Metadata management practices and SOPs</p>	<p>6x, ISO 25024:2015.</p> <p>The representation of metadata follows FAIR standards.</p> <p>The use of best practices with a direct impact on DQ is explicitly reported (e.g.: variables for which explicit negation is used, variables for which absolute values are reported).</p>	<p>deviations are automated.</p> <ul style="list-style-type: none"> • Submitted metadata are generated "by design"
VI. Data manipulation steps¹⁰	<p>Impacts reliability both in terms of accuracy (possible errors) and precision (i.e., the degree of approximation by which data represents reality). Essential to ensure traceability of information. Also impacts coherence and potentially timeliness.</p>	<p>A) A description of data onboarding procedures, e.g.:</p> <ul style="list-style-type: none"> * Frequency and modality of updates * "acceptance tests" performed on RWD sources. e.g.: sources are monitored over time for sudden variation of content, as a proxy to detect process errors <p>B) A description of data manipulation steps, including:</p> <ul style="list-style-type: none"> • Data transformations performed (e.g.: unit of measure conversions, formatting, pivoting, deriving new values, such as BMI from weight and height). 	<p>Tests performed follow some standard or shared set of tests, that can be re-used across RWD sources.</p> <p>Key performance indicators (KPIs) for data cleaning (e.g., data duplications, mislabelling, etc.) are provided.</p> <p>Data mapping tables and algorithms are described with a standard characterisation of their performance.</p> <p>Lists of standard test batteries</p>	<p>Information about data onboarding is directly provided by the platform, e.g.:</p> <ul style="list-style-type: none"> * Transaction logs are available including deviations and actions that required manual intervention <p>Actual data transformation code is accessible and verifiable.</p> <p>Quality checks and KPIs reported are automatically generated by the data platform (e.g.: unit testing)</p> <p>Lineage information is</p>

¹⁰ By "data manipulation" we consider transformations that, in the absence of error, don't affect data reliability: e.g.: unit of measures conversion.

		<ul style="list-style-type: none"> • Data cleaning steps (e.g.: duplicate detection) • Data mapping steps (e.g.: terminology mapping). • Include information about loss of precision expected (e.g., loss of time detail if time of data capture is rounded up to nearest minute; or loss of precision resulting from terminology mapping). <p>C) A description of testing procedures</p> <ul style="list-style-type: none"> • SOP for testing (e.g.: test of pipelines vs test of executions) <p>D) Lineage information</p> <ul style="list-style-type: none"> • Provide justification for the level of data manipulation. • Provide lineage information to specified level sought. 	<p>used to detect loss of accuracy or precision are provided.</p> <p>All lineage information is provided as metadata associated to the dataset</p>	<p>automatically generated by the processing platform</p>
VII. Data augmentation steps¹¹	Data augmentation steps impact accuracy.	<p>A) Information on data augmentation steps (e.g.: imputation or linkage)</p> <ul style="list-style-type: none"> • Justification, methods (algorithms), assumptions, expected error rate • Detail on where such methods are applied. • algorithm such as name, source description and justification for use. 	<ul style="list-style-type: none"> • Algorithms are published, shared and their performance documented. Reference to algorithms is to a specific version. • Information on which values result from imputation is provided as part of the dataset (e.g.: in metadata, 	

¹¹ We consider here data transformations that produce new information subject to reliability issues: e.g.: imputation of missing values, or extraction of codes via natural language processing.

			or data dictionary).	
VIII. Known quality issues and independent QA assessment of the RWD source	Explicit description of known DQ issues, as well as external validation performed (all dimensions affected)	<p>A) Self-reported known DQ issues with an explanation of factors leading to issues (e.g., poor overall completeness in Q3 2020 due to COVID-19)</p> <ul style="list-style-type: none"> • Include a description of known approximations of loss of precisions in mappings. <p>B) Known independent data validations:</p> <p>* Validation studies</p> <ul style="list-style-type: none"> • Publications resulting from this RWD source 		
IX. The RWD source representation	Descriptive of the intended coherence of a dataset and its metadata.	<p>A) Description of the data model used</p> <p>Description of non-standard data model</p> <p>B) Data dictionary and ontologies (vocabularies) in use.</p>	<p>The description refers to a model such as FHIR, OMOP, I2B2, a subset or eventually an extension of such.</p> <p>Information on data dictionaries is:</p> <ul style="list-style-type: none"> • Based on standard vocabularies (such as ICD-9 to ICD-10 diagnosis) • Refers to a specific version of vocabularies used. • When non-standard dictionaries are used, a rationale is provided, and full dictionary is made available 	Data dictionaries are provided using standard formats that facilitate the mappings across different vocabularies and across languages.
X. The RWD source declared Service Level	Descriptive of guaranteed timeliness and possible variations of	A) Data resource declared SLAs	Provide details of established data processes	SLA compliance is automatically assessed and reported

Agreements (SLAs)	extensiveness/reliability provided.	<ul style="list-style-type: none"> Guaranteed frequency of updates Guaranteed incident response time (e.g.: corrections in case of errors) <p>B) Processes and resources accompanying data (e.g.: documentation, training material, help desk).</p> <p>C) Extended capabilities related to DQ: e.g.: possibility to collect additional data if needed</p>	followed by the SLA provider.	
XI. The RWD source licensing and restrictions	Descriptive of aspects that can limit extensiveness and coherence in downstream data aggregations.	<p>A) Details on conditions and processes under which data are made available, such as:</p> <ul style="list-style-type: none"> Features of data use agreements that may limit data use or access (consent, limitations of use). Licensing constraints <p>B) Dataset retention and accessibility policies.</p>	Policies and licensing reported are standardised and applied to a broad range of RWD sources	
XII. Feedback	Descriptive of feedback mechanisms in place to improve all aspects of DQ	A) Provide a contact for QA and follow-up on DQ issues detected.	The contact provided allows tracking of issues and follow-up after standard service support patterns	The feedback mechanism provided includes notification of automatically detected DQ issues.

375 When filling the above form, all fields should be used, eventually clarifying when something doesn't
376 apply and why (e.g.: no processing of the dataset was done).

377 **3.2. General considerations on the characterisation of systems and**
378 **processes**

379 Since the recording or processing of data may have an impact on DQ, every actor involved in any such
380 process should therefore ensure that their actions adhere to the checklist above. Generally, an end
381 user preparing a dataset to support regulatory activities would therefore provide the above checklist
382 for any eventual processing they did on the datasets, plus a checklist for each source of data used,
383 while the RWD source provider is the entity better positioned to provide a checklist covering its specific
384 data assets.

385 Overall, independent of who takes responsibility for the information provided, how DQ information is
386 aggregated, or whether this information is provided in full, summarised, and accessible on demand

387 (e.g.: in case of audit), all the available evidence related to systems and processes potentially affecting
 388 DQ should be clear at the time of submission.

389 4. Data quality metrics for RWD

390 Metrics are the most obvious case of what is introduced in the DQF as “intrinsic determinants”. This is
 391 defined as all DQ aspects that could be assessed based on a dataset itself, without information on how
 392 data have been produced, nor its intended usage.

393 This section introduces metrics that can be used to measure different aspects of DQ. An overall
 394 framework that groups DQ metrics in terms of their requirements and dimensions is introduced in a
 395 way that can help assemble and systematise existing quality metrics into balanced sets, as well as
 396 identify gaps in existing metric sets.

397 This section also provides a list of example metrics. The presented list is not meant to be exhaustive:
 398 there are many DQ metrics outlined in the literature [15] and many more that could be created based
 399 on the individual characteristics of a data type. In the EMRN DQF [1], metrics were presented at an
 400 abstract level and were covering a very wide range of scenarios, including examples beyond clinical
 401 RWD, mostly with the goal of illustrating each quality dimension. What is presented here, is a sample
 402 of concrete metrics that are highly relevant and broadly applicable for characterisation of RWD,
 403 together with RWD-specific examples to illustrate a potential output these metrics can be applied.

404 4.1. Framework for the categorisation and identification of metrics

405 To categorise and identify metrics, a simple framework can be used to test the completeness of test
 406 sets in use, as well as to identify gaps, redundancy, or complementary metrics (See Figure 4). This
 407 framework can be visualised as a simple table, with dimensions in the columns and metric groups in
 408 the rows. Each dimension shown in Figure 4 consists of multiple sub-dimensions, detailed in Tables 3
 409 to 6. Note that not all metric groups apply to every dimension in this table.

	Reliability	Extensiveness	Coherence	Timeliness
Independent data checks				
Plausibility checks				
Conformance checks				
Comparison to other data sources				
Checks on dataset metadata				

410
 411 **Figure 4 – Proposed 2-dimensional framework for metrics identification.**

412 These dimensions are classes of the DQ features that the metrics are meant to measure. Requirement
 413 for metrics assessment is on the other axis. These metric assessment requirements identify what
 414 resources are needed and what additional information or know-how a metric embeds. In terms of
 415 requirements, this section outlines the following metric groups:

416 **Independent data checks.** These are metric groups for which no additional knowledge or information
 417 on the content of the dataset is required. Examples may include the number of empty or corrupted
 418 fields or the number of potential duplicates. Independent data checks can be designed and applied to a
 419 broad range of data.

420 Other metrics can instead embed knowledge about general know-how on the data being measured.
421 Key examples are plausibility metrics and conformance metrics introduced below.

422 **Plausibility checks.** These are metrics that capture DQ aspects based on general knowledge about
423 the world represented in data. For instance, the number of (un)reliable values could be assessed by
424 detecting patterns that are impossible to be present in the Real-World: e.g.: female patients that have
425 observations only occurring in males, measured quantities that exceed a certain magnitude (e.g.: a
426 blood pressure of 1000/500 mmHg), or patterns that are impossible (e.g.: the timing of a causal effect
427 occurring after its effect), etc.

428 **Conformance checks.** Metrics assessing conformity to standards dictating data structures,
429 dictionaries, or format, e.g., all values to represent a condition come from a prescribed terminology
430 source.

431 **Checks on dataset metadata.** This class is based on the know-how on a specific dataset, such as
432 what is provided in metadata or supporting documentation. In some cases, it is useful to consider
433 metrics that are based on the 'descriptors' that come with a dataset (e.g., metadata) that reflect the
434 processes or the standards behind a dataset. For instance, a dataset could be provided with additional
435 dataset detailing what values are recorded, and what is imputed. Metrics summarising the percentage
436 of imputed data could then be used to assess the reliability¹² of a dataset. It is also useful to verify
437 how data values match expectations with respect to metadata constraints. In principle such metrics
438 could measure, by direct verification, the effect of a full data process.

439 **Comparison to other data sources.** Metrics resulting from the comparison against reference RWD
440 sources can support extensiveness and reliability assessments, particularly against broadly recognised
441 RWD sources with demonstrated quality assurance. For example, it is useful to compare the proportion
442 of missing data to a reference dataset to gain an understanding of the possibility of bias in data
443 collection or recording. This kind of metrics can be used to determine the true accuracy or validity of
444 data only in rare cases, for instance when the same type of data has been collected for the same
445 patient in real world and in a randomised clinical trial (RCT) and the latter can be leveraged as gold
446 standard to assess the validity of the RWD.

447 These metrics don't include results of validation of accuracy against original data, as that is expected
448 to be covered in foundational documentation (see section 3.1. . Certain data can also be valid when
449 observed individually, but the collective trend of all data of a kind should follow expected distributions
450 or trends, based on clinical expectations. For example, the prevalence of a disease is unlikely to grow
451 drastically (e.g., from 2% to 80%) in a population from one year to another. In that case, metrics are
452 difficult to determine. Instead, a visual representation of data may be needed to detect abnormal
453 trends and data with low plausibility. This process is also called clinical validity.

454 In the following section, some examples of metrics are provided in relation to this framework. Tables
455 3-6 refer first to the EMRN DQF and more broadly to commonly used DQ checks [20] . This chapter
456 does not refer to the verification (within data) versus validation distinction (compared to other RWD
457 sources), as this is made more detailed and operational by the above implementation categories.

¹² We note that the term "reliability" is used here with the definition presented in the EMRN DQF ("that data correspond to reality"), that differs from its interpretation in statistics ("consistency of repeated measures")

458 **4.2. Metrics for DQ assessments**

459 **Reliability dimension**

460 These metrics are meant to measure the degree to which data correspond to what they intend to
 461 represent.

462 **Table 3 - Overview of reliability metrics by sub-dimensions with examples**

Sub-dimension	Metric group	Metrics	Example
Accuracy	Plausibility checks¹³	<ul style="list-style-type: none"> Number and percent of records where data values don't agree with standards or knowledge or feasible ranges Number and percent of records where values of repeated measurement of the same fact don't show expected variability Number and percent of records with logical inconsistencies Number and percent of records where observed or derived values don't conform to expected temporal properties Number and percent of records where duplicate records aren't flagged. Number and percent of records where data values don't agree with common expectations (e.g., human with 4 arms) 	<ul style="list-style-type: none"> For X% of records/rows, systolic blood pressure values are higher than 250 mmHg X% of records showed >2kg difference when weight was measured by separate nurses within the same facility using the same equipment X% of records of pregnancy were attributed to males For X% of records, discharge date happens before admission date
	Checks on dataset metadata	<ul style="list-style-type: none"> Number and percent of variables/datasets that are based on imputation or derivation 	<ul style="list-style-type: none"> End of treatment date is derived for X% of patients from treatment start date and treatment cycle length
	Comparison to other data sources	<ul style="list-style-type: none"> Number and percent of records where corresponding variables yield identical results across independent or dependent databases 	<ul style="list-style-type: none"> X% of EHR records had a date of birth value that matched the value in a birth

¹³ Accuracy metrics based on general knowledge are typically plausibility metrics, where a dataset is assessed regarding its likelihood to be correct, based on common expectations regarding data distribution or general constraints between different values.

			registry for the same patient
Precision	Independent data checks	<ul style="list-style-type: none"> The number of decimal points used in data values, and their distribution 	<ul style="list-style-type: none"> "Height" in meters recorded with two decimal digits, but the last being always 0.
Traceability	Checks on dataset metadata	<ul style="list-style-type: none"> Number and percent of datasets/variables for which traceability information is available in metadata. 	<ul style="list-style-type: none"> Metadata regarding traceability are available for only 2 out of the 3 datasets feeding into an overall disease registry (treatment data and death records contain traceability-related metadata, but not medical history data)

463 **Extensiveness dimension**

464 **Table 4 - Overview of extensiveness metrics by sub-dimensions with examples**

Sub-dimension	Metric group	Metrics	Example
Completeness	Independent data checks	<ul style="list-style-type: none"> Percentage of records for which data are missing for a given variable 	<ul style="list-style-type: none"> X% of patients have a value recorded for their Date of birth
		<ul style="list-style-type: none"> Percentage of patients who have a certain number of measurements for a given variable 	<ul style="list-style-type: none"> X% of patients have 2 or more ECOG assessments
	Comparison to other data sources	<ul style="list-style-type: none"> Relative percentage of records for which a variable is missing with respect to a trusted source of knowledge 	<ul style="list-style-type: none"> X% of patients have date of diagnosis missing for a diabetes database, compared to a National and institutionally validated diabetes registry

Coverage	Comparison to other data sources	<ul style="list-style-type: none"> Percentage of a target population present in a database 	<ul style="list-style-type: none"> X% of diabetic patients covered in a regional diabetes registry as compared to the National Patient Registry
-----------------	---	---	--

465 **Coherence dimension**

466 **Table 5 - Overview of coherence metrics by sub-dimensions with examples**

Sub-dimension	Metric group	Metrics	Example
Format coherence (conformance)	Conformance checks	<ul style="list-style-type: none"> For relevant variables, % of records where that conform to formatting constraints 	<ul style="list-style-type: none"> X% of records have Sex as only one ASCII character (0 or 1)
		<ul style="list-style-type: none"> For relevant variables, % of records where data values conform to allowable values or ranges 	<ul style="list-style-type: none"> X% of records have sex with one of the 3 allowable values "M", "F". or "U".
Relational coherence (conformance)	Independent data checks	<ul style="list-style-type: none"> Data values conform to relational constraints based on external standards 	<ul style="list-style-type: none"> X% of records having the Provider ID in the Drug exposure data correspond to the record in the Provider table
	Conformance checks	<ul style="list-style-type: none"> For computed values, % of records where computed values conform to programming specifications 	<ul style="list-style-type: none"> For X% of patients, database calculated, and hand calculated BMI (body mass index) values are identical at a 0.2 margin of error
Semantic coherence (conformance)	Conformance checks	<ul style="list-style-type: none"> For relevant variables which employ code lists according to external standards, % of patient records which use a given code list 	<ul style="list-style-type: none"> X% of diagnoses are coded with ICD-10 (as required by CDM)
Uniqueness	Independent data checks	<ul style="list-style-type: none"> Number of records flagged as potential duplicates 	<ul style="list-style-type: none"> Out of X records, 2 are flagged as potential duplicates: William Smith's DOB and ID matches with Bill Smith's DOB and ID.

467 **Timeliness dimension**

468 **Table 6 - Overview of timeliness metrics by sub-dimensions with examples**

Sub-dimensions	Metric group	Metrics	Example
Currency (Is your data acceptably up to date?)	Independent data checks	<ul style="list-style-type: none">Average time of updates in a database	<ul style="list-style-type: none">Timestamps between 2 consecutive forms indicate patient records are updated on average every 3 months

469 These metrics have been provided at a general level where one would apply the metric to all records,
470 however, there can be some hypothesis-driven stratification to look at the data with more granularity
471 (in context of a particular question/context, see below section 7). E.g., for completeness and coverage,
472 one may want to look at the metrics in a stratified way, where there may be a sub-population of
473 particular interest/criticality or where there is an expectation for lower quality.

474 **4.3. Considerations for the implementation of RWD DQ metrics**

475 **4.3.1. Different roles of metrics**

476 This section distinguishes different primary roles of DQ metrics: such roles correspond to different
477 optimal sets of metrics.

478 **4.3.1.1. Quality assurance**

479 When metrics are used for DQ assurance, the intention is to identify issues with the aim of correcting
480 these issues when possible. Such metrics are naturally automated and tend to be as extensive as
481 possible. Test sets comprising hundreds of metrics are possible: anomalies and unexpected values
482 detected can then be screened and lead to follow-up actions, including inspection of sources of data
483 pipelines to identify errors. Often such issues can be prioritised with respect to frequency and severity,
484 hence there are little downsides on test set being extensive, especially when automation is in place.

485 **4.3.1.2. Quality reporting**

486 When metrics are used for DQ reporting, they are meant to provide some high-level assessment of
487 quality that can be used for an assessment of DQ. In this case, metrics should be more high-level and
488 limited in number, and such that some relative assessment of DQ among datasets is possible (e.g.:
489 typical metrics would be average completeness, or average conformance). The value of such high-level
490 characterisation is limited but useful when datasets are presented in a catalogue for a first
491 characterisation of DQ.

492 **4.3.1.3. DQ assessment**

493 More precise than reporting, DQ metrics can be used to reflect the quality of specific data elements,
494 with the view of assessing whether a dataset is or is not suitable to answer a specific research
495 question, e.g.: whether the precision of age is suitable for a paediatric study. In this case metrics need
496 to be presented at a level of detail that makes them unsuitable for high-level reporting. Furthermore,
497 such metrics should be assessed on the population of interest, rather than an overall generic dataset.

498 **4.3.2. Additional considerations on level of application and maturity for**
499 **metric assessments**

500 As per the above description of the different roles of metrics, metrics may be used at different points in
501 the chain of RWD creation and aggregation. For example, for a registry collecting data from multiple
502 hospitals, metrics can be generated at different levels: within the individual hospitals, or at the whole
503 registry level, varying in relevance of metrics (e.g., coherence between sites).

504 As described in the EMRN DQF, metrics may be assessed and reported on with varying levels of
505 maturity. For example, at the lowest level, metrics may have to be estimated and self-reported by the
506 data owner with approximate knowledge of general data trends ('qualitative assessment') and may be
507 generated "ad-hoc". While at higher levels of maturity, 'quantitative assessments' (based directly on
508 the data) should be fully automated. Fully automated checks should take place in capture systems
509 during data collection or recording and then throughout the generation system.

510 **5. Guidelines to assess quality in relation to a specific**
511 **research question**

512 This framework distinguishes the measurement of DQ (in metrics, but also via supporting evidence)
513 from the evaluation of its fitness-for-use, as this relates, by definition, to the relevance of the dataset
514 for a specific question addressed.

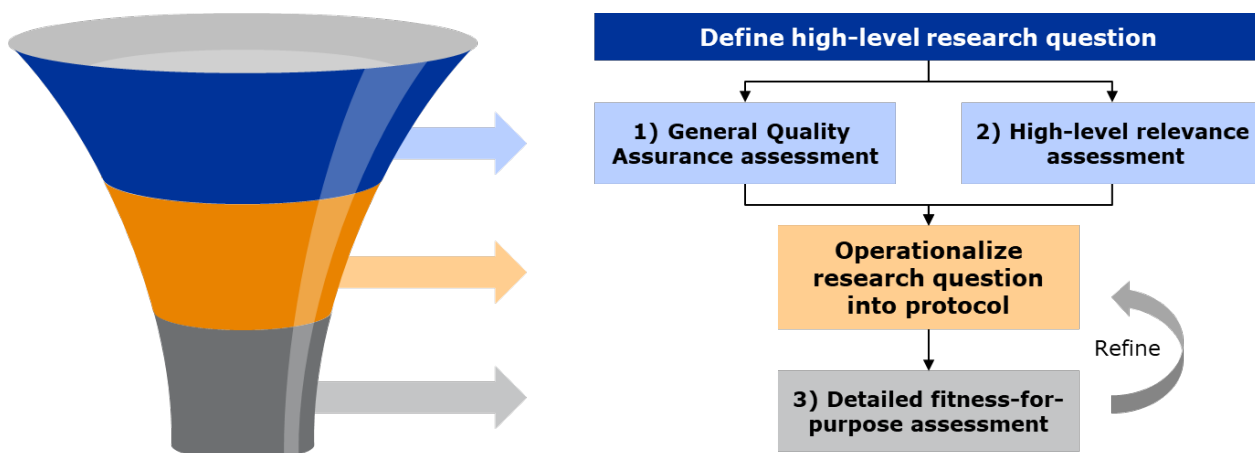
515 It is in fact difficult to pre-specify thresholds or minimum criteria for the fitness-for-use assessment:
516 generally, a lot depends on the type of study, and on disease-dependent and analysis-dependent
517 factors. In addition, there may be some other considerations, such as lack of other RWD sources in a
518 therapeutic area or disease frequency that have an impact on setting acceptability thresholds (e.g. for
519 rare diseases, it might be challenging to identify enough cases via secondary use of data).

520 At a more detailed level, data relevance to a specific question is demonstrated if the data capture key
521 data elements of the research to address such question (e.g., diagnosis, exposure, outcome, and
522 covariates) in a reliable, coherent and timely way, or if the number of patients and follow-up time are
523 sufficient to demonstrate the impact of the intervention/determinant under investigation [14]. To
524 assess the relevance of a RWD source, an in-depth and systematic evaluation of the data source in
525 relation to its design elements is required.

526 In the next section, guidelines are provided to assess the relevance of data to a research question
527 based on DQ metrics and the evidence of systems and processes provided. Note that relevance is not
528 limited to accepting thresholds or evaluating the "usability" of a study. The quality characterisation of a
529 source is also useful input to define applicable methods as well as additional RWD sources that may be
530 required to answer the research question.

531 **5.1. General principles for assessment of data quality in relation to a**
532 **research question**

533 Once a research question has been defined, a set of steps that guide the assessment of the suitability
534 of a RWD source with respect to quality can be considered (See Figure 5).



535

536 **Figure 5 - General principles for assessment of DQ in relation to a research**
 537 **question**

538 **Step 1 – General Quality Assurance assessment**

539 To determine if a RWD source may be relevant for regulatory purposes, the first step is to ensure that
 540 the available information and documentation meet the overall quality requirements in terms of
 541 reliability and, if applicable, timeliness.

542 Since some research questions (e.g., related to pharmacovigilance) are time-sensitive, the overall
 543 characterisation of a source with respect to timeliness (e.g., overall time lag) can be a key criterion for
 544 its acceptability.

545 It is also possible to assess how much a dataset is represented in a coherent way that facilitates
 546 analysis. Documentation on terminologies and standards used, can help in assessing the fitness of a
 547 dataset to a specific analysis goal and process. Coherence assessments do not generally result in
 548 yes/no decisions on the suitability of a dataset, as it can be usually improved with extra efforts¹⁴, but it
 549 can be a criterion when multiple RWD sources are available.

550 This first step of the assessment is typically done by inspecting documentation (e.g.: the systems and
 551 processes checklist), for instance to assess the reliability of data one would look at description on the
 552 presence of QA processes, documentation of any data curation, data transformation/enrichment steps,
 553 etc.

554 **Step 2 – High-level relevance assessment**

555 The next step is to assess the “relevance” of a dataset to a question. In the EMRN DQF, the definition
 556 of “relevance” is narrowed to data having the right kind of variables for the question at hand. To
 557 assess the “fitness-for-use” aspect of a dataset to a specific question, an essential step is indeed to
 558 identify if the content of the information fits the requirements posed by the question. Whether a
 559 dataset presents the right kind of variables can be assessed at high-level based on the overall data
 560 documentation (e.g.: how data are collected, the purpose, the data dictionary).

561 A preliminary assessment of relevance can be conducted by inspecting metadata, without directly
 562 accessing the data or relying on detailed metrics. Increasing knowledge of available data sources (e.g.,
 563 EHR, administrative claims) can help in framing the research question more explicitly to guide the

¹⁴ There is a potential loss of precision when data are harmonised to a common standard. Therefore, if data are not coming in the coherent representation for an intended analysis process, some attention must be put on precision degradation, when assessing specific variables (later in this document).

564 choice of a specific data source. This process can be further facilitated by registering data sources in
565 relevant online repositories, such as the HMA-EMA Catalogue of Real-World Data Sources ¹⁵.

566 **Step 3 – Detailed fitness-for-use assessment**

567 Once the DQ of a RWD source is determined acceptable at an overall level, a specific RWD source
568 inspection is required. To do so, one must first:

- 569 1. Articulate the research question and the relevant design elements such as study population, study
570 sampling (e.g., case-control, cohort), treatment/exposure group, comparator group, primary and
571 secondary outcome(s), length of follow-up, data lag time, confounders.
- 572 2. Operationalise data elements into variables depending on the specificities of the research question
573 to get a better understanding of the disease area (e.g., rate of evolution of standard-of-care, i.e.,
574 how frequently the standard-of-care changes for a given indication, time-to-disease progression).
 - 575 • Where possible, pre-specify the importance of the quality of data elements in the protocol –
576 this assessment should be done in anticipation of the analysis methods (e.g., sample size
577 calculations, use of time-to-event endpoints, sensitivity analyses, statistical adjustment for
578 measurement error), which will impact what is considered acceptable for missing data or
579 errors. While not part of the quality assessment itself, anticipating methods is important to
580 provide context for performing a quality assessment.

581 After this phase, the qualification of the RWD source can be performed by assessing if the data quality
582 of the variables of interest is adequate for the intended analysis. This entails assessing the
583 extensiveness (e.g.: completeness) of the required design elements as well as the coherence,
584 reliability, and timeliness of those elements.

585 Note that the fitness-for-use assessment could be done based on metrics and metadata that are
586 reported for an overall RWD source or could be performed on the final (sub)dataset selected for the
587 study (e.g., specific data cut/sub-population/aggregation of RWD sources).

588 In general, all summary metrics may change when a subset of a population is considered (e.g.: the
589 precision of “age” may change if a subset of a population focusing on paediatric patients is
590 considered). While this is rare for accuracy and timeliness, extensiveness is often affected: for an
591 identified data (sub)set of interest, a fit-for-use assessment also entails seeing if the sample size of
592 planned patient population is enough to guarantee robust evidence, and whether data are
593 representative of the target population when relevant. Coherence in particular needs to be re-assessed
594 each time a new data source is introduced in an analysis.

595 Generally, the RWD source should be chosen to match the research question, rather than adjusting the
596 research question to fit the RWD source. It is important to note that, in some cases, the metrics and
597 characterisation can lead to changes in the study design to accommodate limitations in the data
598 (iterative process). For example, if a rare disease is insufficiently captured in a RWD source or in the
599 patients of interest included in the RWD source, but a broader concept that is also of interest is well-
600 captured, the study may focus on the broader concept instead. In contrast, in a causal study, if
601 important confounders are not captured in the RWD source, it may be necessary to find an alternative
602 RWD source to conduct the study.

¹⁵ <https://catalogues.ema.europa.eu/catalogue-rwd-sources>

603 **5.2. Framework for detailed fitness-for-use assessment**

604 This framework is inspired by The Structured Process to Identify Fit-For-Purpose Data (SPFID) [14].
 605 However, it differs from it in that the aim of this RW-DQF is not to exhaustively look for different RWD
 606 sources and rank them comparatively for their fitness for purpose. It rather provides a guideline to
 607 assess if a data source is suitable for regulatory use.

608 Table 7 provides a template to be filled in during the Step 3 of the fit-for-purpose assessment.

609 **Table 7 - Fitness-for-use assessment to be filled in during the suitability**
 610 **assessment of a RWD source**

Design elements <i>(to be pre-specified)</i>	Operational definition <i>(to be pre-specified)</i>	Data elements for valid capture <i>(to be pre-specified)</i>	Criticality of the quality of the element, including justification where relevant <i>(to be pre-specified)</i>	Suggested extensiveness assessment <i>(to be filled in during assessment)</i>	Suggested assessment of other quality dimensions <i>(to be filled in during assessment)</i>	Suggested substantiation by documentation <i>(to be filled in during assessment)</i>
Study population	Inclusion criteria <i>Criterion 1</i> ... <i>Criterion n</i> Exclusion criteria <i>Criterion 1</i> ... <i>Criterion n</i>	Data elements required for I/E criteria	Low/Medium/High	Completeness metrics	Reliability metrics (Precision)	As relevant
	Cohort size	Minimum cohort size	Low/Medium/High	N/A (to be assessed on a research question basis)	N/A	N/A
	Representativeness of the target population	Population characteristics for which similarity to those of the studied sample is important	Low/Medium/High	Coverage metrics	N/A	As relevant
Treatment/exposure		Data elements required	Low/Medium/High	Completeness metrics	Reliability metrics	As relevant
	Newly treated population size	Minimum number	Low/Medium/High	N/A (to be assessed on a research question basis)	N/A	N/A
Comparator group (if relevant)		Data elements required	Low/Medium/High	Completeness metrics	Reliability metrics	As relevant
	Size of comparator sample	Minimum number	Low/Medium/High	N/A (to be assessed on a research question basis)	N/A	N/A
Key endpoint	Key endpoint 1 ... Key endpoint n	Data elements required	Low/Medium/High	Completeness metrics (overall and over time)	Reliability metrics Coherence metrics Timeliness metrics	As relevant

Confounders (if relevant)	Confounder 1 ... Confounder n	Data elements required	Low/Medium/High	Completeness metrics (overall and over time)	Reliability metrics Coherence metrics	As relevant
Follow-up time (if relevant)		Minimum follow-up	Low/Medium/High	Coverage metric on follow-up time	Timeliness metrics (overall or for specific variables if relevant)	As relevant
Lag time		Maximum lag time	Low/Medium/High	N/A	Timeliness metrics (overall or for specific variables if relevant)	N/A

611

612 **5.3. Illustrative example for detailed fitness-for-use assessment**

613 An example is provided here for a Chronic Lymphocytic Leukaemia External Comparator study based
614 on multi-site EHR. Note that this is purely an illustrative use case to demonstrate how to use the
615 framework for step 3 (See
616 Table 8).

617 **Table 8 - Detailed fitness-for-use assessment for a Chronic Lymphocytic Leukaemia**
618 **study**

Design elements (pre-specified)	Operational definition (pre-specified)	Data elements for valid capture (pre-specified)	Criticality of the quality of the element, including justification where relevant (pre-specified)	Extensiveness assessment (filled in during feasibility assessment)	Other quality assessment (filled in during feasibility assessment)	Documentation (filled in during feasibility assessment)
Study population	Age >18 years at time of enrolment	Date of birth (month)	High	100% of patients have DOB available, or age at registration to the RWD source	100% of patients have DOB captured in the same month and year format (MM/YYYY)	Documentation on the RWD source's target population age range and format of age/DOB capture (if available)
	Confirmed diagnosis of CLL	Physician diagnosis (ICD 10 code or equivalent)	High	100% of patients have a CLL diagnosis	100% of diagnoses have been mapped to ICD-10	Documentation on mapping of different coding systems to ICD-10
		Lab results	Medium	40% of patients have confirmatory lab results	100% of lab results are within a plausible range	Documentation on consistency of lab assessments

						across different sites
	Known 17p deletion status	17p deletion status	Low (possibility of using <i>probabilistic bias methods using published prevalence of 17p deletion and its association with selected endpoints to derive subgroup estimates</i>)	70% of patients have known 17p deletion status	Time lag (6 months) between 17p deletion availability and initial diagnosis date	N/A
	Cohort size	5000 patients	Medium	6000 patients after application of I/E criteria of interest	N/A	N/A
	Representativeness vs target emulation population (from RCT)	Average age in acceptable +/- range compared to target population	High (<i>bias towards worse outcomes if older population</i>)	Average age is 83 vs 82 in target population	N/A	N/A
Treatment/ exposure	Received a BTKi	Treatment information (BTKi)	High	BTKis are in the list of drugs covered by this database, and 90% of patients have at least 1 record of treatment	95% of cancer treatment information has a date after diagnosis of CLL. 90% of treatment records pass uniqueness checks	N/A
	Number of newly treated	300 patients receiving BTKi after confirmed diagnosis of CLL	High	315 patients	N/A	N/A
Comparator group (if relevant)	Received best supportive care	Absence of anti-cancer treatment	High	Explicit negation of treatment received only	N/A	N/A

				for 20% of patients		
	Number in comparator	300 patients who are on best supportive care	Medium	270 patients	N/A	N/A
Key endpoint	Overall Response Rate	Response per criteria at intervals	High (<i>consistency of response assessment essential for primary outcome</i>)	85% of patients with at least one assessment 40% with 3 assessments or more	N/A	Documentation detailing re-assessment of response by adjudication committee for homogeneous assessment
		Treatment regimen and/or cycle start date	Medium	90% of patients with start date available	Variable only available at month level for 50% of records	N/A
	Overall Survival	Date of death	High	20% of patients with known death have date of death	Statistical checks for reliability of linkage to death registry passed for 100% of patients	Linkage process documentation Traceability documentation
		Treatment regimen and/or cycle start date	Medium	90% of patients with start date available	Variable only available at month level for 50% of records	N/A
	Number of participants with AE	AE capture across patients during follow-up period	High	30% of patients with AE data available	N/A	Documentation detailing method for AE capture, and which AEs are to be captured
Confounders (if relevant)	Sex	Sex, reported as male or female	High	100% of patients have sex information available	100% of patients reported with a pregnancy are classified as	Documentation detailing method for sex capture

					females 100% of patients have sex captured as one ASCII character (0 or 1)	
	Cancer stage	Stage I-IV	Medium	80% of patients have at least one stage record 30% have stage at each line of therapy	Distribution of staging observed to be different pre- and post-2017 (due to update in guidelines)	Documentation of internal guidelines for consistent stage assessment
Follow-up time (if relevant)	6 months	6 months follow-up	Medium	80% of patients have >6 months follow up	N/A	Documentation of internal guidelines for typical follow-up time (if available)
Lag time	2 years	2 years maximum	Medium	There is on average a 3-year lag time to perform data curation and linkage	N/A	Documentation of the RWD source about lag time

619 DOB: date of birth; CLL: Chronic Lymphocytic Leukaemia; N/A: not assessed; I/E: in- and exclusion criteria; BTKi:
620 Bruton tyrosine kinase inhibitor; AE: adverse event.

621 When assessing the fitness-for-use of the RWD source, this table can provide guidance in making a
622 final decision on the suitability of a dataset for a given study, or prompt changes in the
623 method/protocol when necessary to leverage available data.

624 **5.4. Toward a generalisation of question-specific aspects**

625 The maturity model for “question-specific determinants” in the EMRN DQF suggests the definition of
626 typical use cases that could be the basis for some more standardised approach to DQ acceptance
627 processes (e.g.: more automated domain-based quality requirement and test packages, as well as
628 standardised reporting). As a step in this direction, the design elements to be pre-specified (outlined in
629 Tables 7 and 8) (i.e., study population, treatment/exposure, comparator group, key outcomes,
630 confounders, follow-up time, lag time) are relevant for different study questions to guide the specific
631 DQ assessment. For example, confounders can be considered a relevant design element across various
632 study types (e.g., clinical management and unmet needs, drug utilisation, disease epidemiology, etc),
633 while comparator group(s) may only be relevant for comparative effectiveness studies. These typical
634 design elements, which might also impact regulatory actions, are important but not prescriptive.

635 **5.5. Providing supporting information for RWD in regulatory submissions**

636 Where RWD are used for regulatory submissions, the data source characteristics allowing for DQ
637 assessment should be provided, with the adequate granularity for regulatory assessment. Relevant
638 characteristics include:

- 639 • Information on the standard quality management practices routinely applied to the data, as
640 well as the processes and methods behind the generation of data at the source. This includes
641 details on the level of automation and the use of computerised systems and can be relevant to:
 - 642 ○ Data cleaning, extraction, or transformation steps.
 - 643 ○ Data quality checks to detect logical inconsistencies and erroneous, missing or out-of-
644 range values.
 - 645 ○ Any remedial actions taken at the RWD source level. Information on the data collection
646 or recording process and any selection mechanisms involved (e.g., inclusion of specific
647 patients, capturing of specific clinical data).
- 648 • Measures taken to improve the completeness of data elements (e.g. data collection
649 prerequisites for reimbursement);
- 650 • Any linkage performed on the data, including details on the data elements used for linking and
651 the linkage methodology;
- 652 • Information on whether patient-level data or only aggregate data are available.

653 This information can be made publicly available for transparency by the data holder by registering the
654 RWD source in the HMA-EMA Catalogue of RWD Sources. The Catalogue is a repository of metadata
655 collected from RWD sources and contains information on the systems and processes behind data
656 capture, as well as descriptors of the data. It is intended to capture the extent of variety of existing
657 RWD sources and facilitate fit-for-purpose assessments [21].

658 To allow for an adequate DQ assessment, it is important that the information published in the
659 Catalogue of RWD Sources remains up to date, with the last update occurring within the past 12
660 months. This provision may be included in the contractual agreement between the MAH or Applicant
661 and the data holder, as relevant.

662 **6. Concluding remarks**

663 This document is an extension of the generic EMRN DQF for medicines regulation, focusing on RWD
664 specificities. The RW-DQF provides background on how DQ can impact the use of RWD to generate
665 RWE for regulatory assessment. It further provides guidance for the characterisation of systems and
666 processes underpinning data and their impact, key metrics to assess different aspects of DQ within a
667 dataset as well as guidance to use metrics to assess the suitability of a dataset by a fitness-for-use
668 assessment in relation to a specific research question. These parts provide actionable and focused
669 recommendations for assessing DQ of RWD with the goal of improving the usefulness of RWE for
670 regulatory purposes.

671 **7. References**

- 672 1. Data Quality Framework for EU medicines regulation. 2022. Available from:
673 [https://www.ema.europa.eu/en/documents/regulatory-procedural-guideline/data-quality-
framework-eu-medicines-regulation_en.pdf](https://www.ema.europa.eu/en/documents/regulatory-procedural-guideline/data-quality-
674 framework-eu-medicines-regulation_en.pdf).

- 675 2. Reflection paper on use of real-world data in non-interventional studies to generate real-world
676 evidence. 2024. Available from: [https://www.ema.europa.eu/en/documents/scientific-](https://www.ema.europa.eu/en/documents/scientific-guideline/reflection-paper-use-real-world-data-non-interventional-studies-generate-real-world-evidence_en.pdf)
677 [guideline/reflection-paper-use-real-world-data-non-interventional-studies-generate-real-world-](https://www.ema.europa.eu/en/documents/scientific-guideline/reflection-paper-use-real-world-data-non-interventional-studies-generate-real-world-evidence_en.pdf)
678 [evidence_en.pdf](https://www.ema.europa.eu/en/documents/scientific-guideline/reflection-paper-use-real-world-data-non-interventional-studies-generate-real-world-evidence_en.pdf).
- 679 3. Cave, A., X. Kurz, and P. Arlett, Real-World Data for Regulatory Decision Making: Challenges
680 and Possible Solutions for Europe. Clin Pharmacol Ther, 2019. **106**(1): p. 36-39.
- 681 4. European Commission. Can we use data for another purpose? 2022. Available from:
682 [https://ec.europa.eu/info/law/law-topic/data-protection/reform/rules-business-and-](https://ec.europa.eu/info/law/law-topic/data-protection/reform/rules-business-and-organisations/principles-gdpr/purpose-data-processing/can-we-use-data-another-purpose_en)
683 [organisations/principles-gdpr/purpose-data-processing/can-we-use-data-another-purpose_en](https://ec.europa.eu/info/law/law-topic/data-protection/reform/rules-business-and-organisations/principles-gdpr/purpose-data-processing/can-we-use-data-another-purpose_en).
684 5. RWE in regulatory assessment and decision-making processes. A report on the experience with
685 regulatory-led RWD studies.
686 [https://www.ema.europa.eu/en/documents/presentation/presentation-report-experience-](https://www.ema.europa.eu/en/documents/presentation/presentation-report-experience-regulatory-led-rwd-studies-s-prilla-ema_en.pdf)
687 [regulatory-led-rwd-studies-s-prilla-ema_en.pdf](https://www.ema.europa.eu/en/documents/presentation/presentation-report-experience-regulatory-led-rwd-studies-s-prilla-ema_en.pdf). 27 June 2023.
- 688 6. European Health Data Space Data Quality Framework, Deliverable 6.1 of TEHDAS EU 3rd
689 Health Program (GA: 101035467). [https://tehdas.eu/app/uploads/2022/05/tehdas-european-](https://tehdas.eu/app/uploads/2022/05/tehdas-european-health-data-space-data-quality-framework-2022-05-18.pdf)
690 [health-data-space-data-quality-framework-2022-05-18.pdf](https://tehdas.eu/app/uploads/2022/05/tehdas-european-health-data-space-data-quality-framework-2022-05-18.pdf).
- 691 7. Rudrapatna, V.A. and A.J. Butte, Opportunities and challenges in using real-world data for
692 health care. J Clin Invest, 2020. **130**(2): p. 565-574.
- 693 8. Dusetzina, S.B., Tyree, S., Meyer, A.M., Linking Data for Health Services Research: A
694 Framework and Instructional Guide. Rockville (MD). Chapter 4 An Overview of Record Linkage
695 Methods. Agency for Healthcare Research and Quality (US). Available from:
696 <https://www.ncbi.nlm.nih.gov/books/NBK253312/>, 2014.
- 697 9. Big data use for public health: publication of Big Data Steering Group workplan 2022-25.
698 European Medicines Agency. 2022. Available from: [https://www.ema.europa.eu/en/news/big-](https://www.ema.europa.eu/en/news/big-data-use-public-health-publication-big-data-steering-group-workplan-2022-25)
699 [data-use-public-health-publication-big-data-steering-group-workplan-2022-25](https://www.ema.europa.eu/en/news/big-data-use-public-health-publication-big-data-steering-group-workplan-2022-25).
- 700 10. Kent, S., et al., Common Problems, Common Data Model Solutions: Evidence Generation for
701 Health Technology Assessment. PharmacoEconomics, 2021. **39**(3): p. 275-85.
- 702 11. Rivera, D.R., et al., The Friends of Cancer Research Real-World Data Collaboration Pilot 2.0:
703 Methodological Recommendations from Oncology Case Studies. Clin Pharmacol Ther, 2022.
704 **111**(1): p. 283-292.
- 705 12. The European Network of Centres for Pharmacoepidemiology and Pharmacovigilance (ENCEPP).
706 Guide on Methodological Standards in Pharmacoepidemiology (Revision 11). EMA/95098/2010.
707 Available from: https://encepp.europa.eu/encepp-toolkit/methodological-guide_en.
- 708 13. Guideline on registry-based studies - Scientific guideline. European Medicines Agency. 2021.
709 Available from: [https://www.ema.europa.eu/en/guideline-registry-based-studies-scientific-](https://www.ema.europa.eu/en/guideline-registry-based-studies-scientific-guideline)
710 [guideline](https://www.ema.europa.eu/en/guideline-registry-based-studies-scientific-guideline).
- 711 14. Gatto, N.M., et al., The Structured Process to Identify Fit-For-Purpose Data: A Data Feasibility
712 Assessment Framework. Clin Pharmacol Ther, 2022. **111**(1): p. 122-134.
- 713 15. Kahn, M.G., et al., A Harmonized Data Quality Assessment Terminology and Framework for the
714 Secondary Use of Electronic Health Record Data. EGEMS (Wash DC), 2016. **4**(1): p. 1244.
- 715 16. Schmidt, C.O., et al., Facilitating harmonized data quality assessments. A data quality
716 framework for observational health research data collections with software implementations in R. BMC Med Res Methodol, 2021. **21**(1): p. 63.
- 717 17. NESTcc. Data Quality Framework, A report of the Data Quality Subcommittee of the NEST
718 Coordinating Center - An initiative of MDIC. 2020 [February 14th, 2022]. Available from:
719 <https://nestcc.org/nestcc-data-quality-framework>.
- 720 18. REQueST Tool and its vision paper [Internet]. EUnetHTA. 2019. Available from:
721 <https://www.eunetha.eu/request-tool-and-its-vision-paper/>.
722
- 723 19. Documentation: Getting started [observational health data sciences and informatics. Available
724 from: https://www.ohdsi.org/web/wiki/doku.php?id=documentation:getting_started.
- 725 20. Observational Health Data Sciences, Informatics. Chapter 15 Data Quality. 2021. Available
726 from: <https://ohdsi.github.io/TheBookOfOhdsi/DataQuality.html>.
- 727 21. Good Practice Guide for the use of the Metadata Catalogue of Real-World Data Sources. 2022.
728 Available from: [https://www.ema.europa.eu/en/documents/regulatory-procedural-](https://www.ema.europa.eu/en/documents/regulatory-procedural-guideline/good-practice-guide-use-metadata-catalogue-real-world-data-sources_en.pdf)
729 [guideline/good-practice-guide-use-metadata-catalogue-real-world-data-sources_en.pdf](https://www.ema.europa.eu/en/documents/regulatory-procedural-guideline/good-practice-guide-use-metadata-catalogue-real-world-data-sources_en.pdf).
- 730 22. Compassionate use. European Medicines Agency. 2024. Available from:
731 [https://www.ema.europa.eu/en/human-regulatory-overview/research-](https://www.ema.europa.eu/en/human-regulatory-overview/research-development/compassionate-use)
732 [development/compassionate-use](https://www.ema.europa.eu/en/human-regulatory-overview/research-development/compassionate-use).
- 733

734 **Definitions**

Abbreviation	Definition
CDM	Common data model
DQ	Data Quality
DQF	Data Quality Framework
EHDS	European Health Data Space
EHR	Electronic Health Record
EMRN	European Medicines Regulatory Network
ENCePP	The European Network of Centres for Pharmacoepidemiology and Pharmacovigilance
ETL	Extract, Transform, Load
I/E criteria	Inclusion/Exclusion criteria
KPI	Key Performance Indicator
QMS	Quality Management System
RWE	Real-World Evidence
RCT	Randomised Clinical Trial
RWD	Real-World Data
SLA	Service Level Agreement
SOP	Standard Operating Procedure
TEHDAS	Towards European Health Data Space

735 **Glossary**

736 The detailed definitions and concepts, with accompanying examples are found in the EMRN DQF [1].
 737 However, to facilitate the reading of this document, a glossary addressing frequently used terms is
 738 provided below:

Definition	Explanation
Coherence	Coherence (also referred to as Consistency) is defined as the dimension that expresses how different parts of an overall dataset are consistent in their representation and meaning.
Data linkage	Data linkage is the process of bringing information from different data sources together for the same person/identifier or entity to create a new, richer dataset. Data linkage allows researchers to exploit and enhance existing data sources without the time and cost associated with primary data collection. Linked data can be used to supplement studies by creating population-level cohorts with longer follow-up and can answer questions that require large sample sizes

Definition	Explanation
	(e.g., for rare diseases) or whole population coverage (e.g., for pandemic response planning).
Extensiveness	Extensiveness is defined as the dimension capturing the amount of data available.
Foundational determinants	A characterisation of the systems and processes underpinning data generation and manipulation that have an impact on DQ.
Intrinsic determinants	DQ aspects that can be observed only on the basis of a given dataset, without requirement for information about how the data were captured, or about its primary/intended use.
Maturity model	Provide guidance as to how determinants (foundational, intrinsic, and question-specific) can be characterised in successive levels of maturity.
Plausibility metrics	Indicators of plausibility that can be used as proxy to detect errors. When some combination of information is unlikely (or impossible) to happen in the Real-World this reveals accuracy issues. For example, a weight of a person exceeding 300 kg is possible, but the weight of many or all persons in a dataset exceeding that value is implausible and likely revealing some errors in the measurement or the processing of the data.
Primary use of data	Primary use of (electronic) health data is the processing of personal health data for the provision of health services to assess, maintain or restore the state of health of the person it belongs to, including the prescription, dispensation and provision of medicinal products and medical devices, as well as for relevant social security, administrative or reimbursement services.
Question-specific determinants	Aspects of DQ that cannot be assessed independently of a research question.
Relevance	Relevance is defined as the extent to which a dataset presents the data elements useful to answer a given research question.
Reliability	Reliability is defined as the dimension that covers how closely the data reflect what they are directly measuring.
Representativeness	Representativeness is defined as the data having the same characteristics as the whole it is meant to represent.
RWD end users	People getting access to and using RWD for secondary purposes, such as using RWD from multiple RWD sources as external comparators to a clinical trial arm, in submissions to regulators and payers/HTAs, using RWD from multiple RWD sources to assess the RW safety of a treatment across geographies and ethnicities, etc.
RWD holder	People owning and or holding the RWD

Definition	Explanation
RWD practitioners	People involved in the RWD collection or recording process such as researchers, data analysts and data custodians.
RWD submitter	RWD submitters are the RWD end users or stakeholders that get access to RWD data and use it for secondary purpose to answer a research question.
Secondary use of data	Secondary use of (electronic) health data is the processing of health data for other purposes rather than primary use such as national statistics, education/teaching, scientific research etc. The data used may include personal health data initially collected in the context of primary use, but also electronic health data collected for the purpose of secondary use.
Quality Management System	A QMS is a formalised approach adopted by an organisation that documents processes, procedures, and responsibilities for achieving quality policies and objectives (e.g., Good Clinical Practices, Good Laboratory Practices or Good Manufacturing Practice. It achieves these quality objectives through quality planning, quality assurance, quality control and quality improvement. Standards like the ISO 9000 family define QMS across industries, while more specific QMS have been developed for specific industry or products.