

1

2



eunetha

EUROPEAN NETWORK FOR HEALTH TECHNOLOGY ASSESSMENT

7

8

9

10

11

12

13

14

15

16

17

EUnetHTA21 - Individual Practical Guideline Document

18

19

20

**Validity of clinical studies
(Project D4.6)**

21

22

23

24

Template version XXX, XX 2022

25

26 Document history and contributors

Version	Date	Description
V0.1	23/03/2022	First draft for CSCQ and NC-HTAb review
V0.2	25/05/2022	Second draft for CSCQ and NC-HTAb review
V0.3	04/07/2022	Third draft for public consultation

27

28 Disclaimer

29 This Practical Guideline was produced under the Third EU Health Programme through a service
 30 contract with the European Health and Digital Executive Agency (HaDEA) acting under the mandate
 31 from the European Commission. The information and views set out in this Practical Guideline are those
 32 of the author(s) and do not necessarily reflect the official opinion of the Commission/ Executive Agency.
 33 The Commission/Executive Agency do not guarantee the accuracy of the data included in this study.
 34 Neither the Commission /Executive Agency nor any person acting on the Commission's/Executive
 35 Agency's behalf may be held responsible for the use which may be made of the information contained
 36 therein.

37

38 Participants

Hands-on Group	Gemeinsamer Bundesausschuss (G-BA), Germany Haute Autorité de Santé (HAS), France Institut für Qualität und Wirtschaftlichkeit im Gesundheitswesen (IQWiG), Germany
Project Management	Zorginstituut Nederland (ZIN), the Netherlands
Committee for Scientific Consistency and Quality (CSCQ)	Agencia Española de Medicamentos y Productos Sanitarios (AEMPS), Spain Austrian Institute for Health Technology Assessment (AIHTA), Austria Belgian Health Care Knowledge Centre (KCE), Belgium Gemeinsamer Bundesausschuss (G-BA), Germany Haute Autorité de Santé (HAS), France Institut für Qualität und Wirtschaftlichkeit im Gesundheitswesen (IQWiG), Germany Italian Medicines Agency (AIFA), Italy
Consortium Executive Board (CEB)	National Authority of Medicines and Health Products, I.P. (INFARMED), Portugal National Centre for Pharmacoeconomics, St. James Hospital (NCPE), Ireland National Institute of Pharmacy and Nutrition (NIPN), Hungary Norwegian Medicines Agency (NOMA), Norway The Dental and Pharmaceutical Benefits Agency (TLV), Sweden Zorginstituut Nederland (ZIN), the Netherlands

39 The work in EUnetHTA 21 is a collaborative effort. Whereas the agencies in the Hands-on Group actively wrote the deliverable,
 40 the entire EUnetHTA 21 consortium was involved in its production throughout various stages. Thus, the CSCQ has reviewed and
 41 discussed several drafts of the deliverable before validation and the CEB has endorsed the final deliverable before publication.

42 Copyright

43 All rights reserved.

44

45 **Table of Contents**

46	1 INTRODUCTION	5
47	1.1 <i>PROBLEM STATEMENT</i>	5
48	1.2 <i>SCOPE/OBJECTIVE(S) OF THE GUIDELINE</i>	5
49	1.3 <i>RELEVANT ARTICLES IN REGULATION (EU) 2021/2282</i>	5
50	2 GENERAL CONSIDERATIONS	6
51	3 CLINICAL STUDY DESIGNS	8
52	3.1 <i>TERMINOLOGY</i>	8
53	3.1.1 <i>INTERVENTIONAL STUDIES</i>	8
54	3.1.2 <i>OBSERVATIONAL STUDIES</i>	9
55	3.2 <i>CLASSIFICATION</i>	10
56	4 SPECIFIC STRENGTHS, WEAKNESSES, AND RECOMMENDATIONS REGARDING	
57	DIFFERENT DESIGNS	11
58	4.1 <i>RANDOMISED CLINICAL TRIAL: GOLD STANDARD</i>	11
59	4.2 <i>NONRANDOMISED CONTROLLED TRIAL</i>	12
60	4.3 <i>UNCONTROLLED CLINICAL TRIALS (E.G., SINGLE-ARM TRIALS)</i>	12
61	4.4 <i>COHORT STUDIES</i>	13
62	4.5 <i>CASE-CONTROL STUDIES</i>	13
63	4.6 <i>CROSS-SECTIONAL STUDIES</i>	14
64	4.7 <i>CASE-SERIES AND CASE-REPORTS</i>	14
65	5 PARTICULARITIES	14
66	5.1 <i>MASTER PROTOCOLS</i>	14
67	5.1.1 <i>PLATFORM TRIALS</i>	15
68	5.1.2 <i>BASKET TRIALS</i>	16
69	5.1.3 <i>UMBRELLA TRIALS</i>	17
70	5.2 <i>REAL-WORLD DATA AND REAL-WORLD EVIDENCE</i>	17
71	5.3 <i>REGISTRIES</i>	19
72	6 REFERENCES	20

73

74

75 **List of Abbreviations**

CEB	Consortium Executive Board
CI	Confidence interval
CSCQ	Committee for Scientific Consistency and Quality
EUnethTA	European Network of Health Technology Assessment
HTA	Health Technology Assessment
HTAb	Health Technology Assessment body
HTAR	Health Technology Assessment Regulation
HTD	Health Technology Developer
ICH	International Council for Harmonisation of Technical Requirements for Pharmaceuticals for Human Use
JCA	Joint Clinical Assessment
MAMS	Multi-arm, multi-stage trials
PICO	Population, intervention, comparator, outcome
RCT	Randomised clinical trial
RoB	Risk of bias
RWD	Real-world data
RWE	Real-world evidence
TWIC	Trial within a cohort

76

77 1 INTRODUCTION

78 1.1 Problem statement

79 One key element of Health Technology Assessment (HTA) is to assess and describe the certainty (and
80 validity) of clinical study results in an objective, reproducible, and transparent way. In 2020, the
81 European Network of Health Technology Assessment (EUnetHTA) Executive Board concluded that
82 GRADE (1) (or any other system for rating the overall quality of evidence and developing healthcare
83 recommendations) can only partially be applied within EUnetHTA because overall conclusions or
84 recommendations might interfere with the independent contextualisation and decision-making at the
85 national level (2). However, valid scientific principles are still required, not only to guide the development
86 of Joint Clinical Assessments (JCAs) at the European level, but also to support the understandability
87 and usability of these results for national decision-making.

88 1.2 Scope/Objective(s) of the Guideline

89 This Practical Guideline is dedicated to the definition, classification, and certainty of results of studies
90 leading to the statistical analysis of what is considered one data set (one sample of patients). Studies
91 that consist in evidence synthesis by pooling the results of multiple already-analysed data sets from
92 multiple samples of patients [e.g., pairwise meta-analysis, indirect comparison, or interventional studies
93 using external control (including historical control)] are not included in this Guideline. The EUnetHTA
94 21 Methodological and Practical Guidelines *Direct and Indirect Comparisons* provide recommendations
95 and guidance for the classification of these evidence syntheses. Finally, the present Guideline does not
96 offer guidance on how to assess diagnostic accuracy studies, because these studies might have a
97 conventional cross-sectional or cohort design, but still require specific assessment of internal validity
98 (3).

99 The way in which the validity of clinical studies will be assessed and interpreted for drawing conclusions
100 at a national level cannot be dissociated from the population, intervention, comparator, outcome (PICO)
101 question that will be formulated by Health Technology Assessment bodies (HTAbs) (see the EUnetHTA
102 21 Practical Guideline *Scoping Process*). For complementary elements relating to the reporting and
103 assessment of multiple hypothesis testing, subgroup, sensitivity, and post-hoc analyses, the reader is
104 referred to the EUnetHTA 21 Practical Guideline *Applicability of Evidence - Practical Guideline on*
105 *Multiplicity, Subgroup, Sensitivity, and Post-Hoc Analyses*. Additional considerations of the definition of
106 clinically relevant outcomes and endpoints, and the assessment of their validity, reliability, and
107 interpretability are discussed within the EUnetHTA 21 practical guideline *Endpoints*.

108 Note: EUnetHTA 21 practical guidelines *Direct and Indirect Comparisons* and *Endpoints* will not be
109 available at the time of the public consultation for this Guideline.

110 1.3 Relevant articles in Regulation (EU) 2021/2282

111 Articles from Regulation (EU) 2021/2282 directly relevant to the content of this practical guideline are:

- 112 • Recital (14);
- 113 • Recital (28);
- 114 • Article 8: Initiation of Joint Clinical Assessments;
- 115 • Article 9: Joint Clinical Assessment Reports and the Dossier of the Health Technology
116 Developer.

117 2 GENERAL CONSIDERATIONS

118 HTA requires the relative effectiveness of an intervention to be determined as correctly and precisely
119 as possible. Relative effectiveness is the quantification of the effect caused by the intervention relative
120 to a comparator (e.g., standard of care or placebo) on an outcome of interest. Interventions can be
121 medicinal products, medical devices, *in vitro* diagnostic medical devices, medical procedures, as well
122 as measures for disease prevention, diagnosis, or treatment. For any effectiveness assessment, it is
123 essential to examine and report the certainty of results systematically. Given that the certainty of results
124 is fundamental, this needs to be communicated alongside the numerical results. According to Article 9
125 of EU-HTA Regulation (4), it is therefore essential that a JCA contains a description of both ‘the relative
126 effects of the health technology’ and ‘the degree of certainty of the relative effects, taking into account
127 the strengths and limitations of the available evidence’.

Practical Guideline (requirement for JCA reporting)

Any effectiveness result in a JCA report must be accompanied by a description of its certainty.

128 The certainty of effectiveness results is determined by three concepts: **internal validity** [i.e., the extent
129 to which a study is free from bias (also called systematic errors), a concept analogous to Risk of Bias
130 (RoB)]; **applicability** (i.e., the extent to which study results provide a basis for generalisation to the
131 target population, a concept close to external validity and generalisability); and **statistical precision**
132 (i.e., the extent to which study results are free from random errors resulting from sampling hazards).
133 These three concepts assess three different dimensions of the certainty of results, which, for example,
134 means that shortcomings in internal validity cannot be remedied by higher statistical precision.
135 Furthermore, evidence that has high internal validity does not necessarily have high external validity.
136 Although HTA usually requires a high target certainty of results, it is necessary to assess all available
137 data, as submitted by the Health Technology Developer (HTD). Nevertheless, there might be
138 justification to not assess the evidence that ranges below a minimum level of internal validity,
139 applicability, or statistical precision in detail, if the PICO question can be sufficiently answered on the
140 basis of higher-certainty results. Furthermore, the certainty of results is independent of the medical
141 context of the PICO question. It is methodologically inappropriate, for example, to take the rareness of
142 a disease or the impossibility of blinding as an excuse to ignore or to euphemise the resulting
143 uncertainties in the clinical evidence.

144 Following international standards of evidence-based medicine, the internal validity of a study has a
145 paramount role in determining the overall certainty of the study results (i.e., if study results have a low
146 level of internal validity, the levels of statistical precision and external validity are irrelevant) (5,6). The
147 classical **hierarchy of evidence** (7) includes several types of study, from case-reports and nonclinical
148 data (level 5 evidence, the lowest level of evidence), case-control studies (level 4), retrospective (or
149 lower-quality) cohort studies (level 3), prospective (or higher-quality) cohort studies (level 2), up to
150 randomised controlled trials (RCTs; level 1, the highest level of evidence). Classification of study design
151 alone (see Section 3) is insufficient for a full assessment of internal validity (8,9), but has much practical
152 value for distinguishing between higher- and lower-quality evidence and for selecting a suitable RoB
153 assessment tool.

Practical Guideline

For internal validity, it is useful in a JCA to distinguish between different study designs.

154 RoB can be defined as any potential systematic error in clinical research that might lead to an incorrect
155 estimate of the effect of interest. RoB can be present at different levels, including: (i) the meta level
156 (e.g., publication bias in a systematic review or meta-analysis); (ii) the study level (e.g., confounding
157 bias in a cohort study); and (iii) the outcome level (e.g., information bias caused by unblinded
158 assessment of an outcome). If the type of evidence requires it, the assessment of RoB needs to be
159 level specific; however, the scope of the present Guideline is limited to bias at the study and outcome
160 levels. Furthermore, some types of bias can occur only in certain study designs, whereas other types
161 can affect all types of study. Therefore, different tools have been developed for RoB assessment in
162 different study designs (10,11). It is essential to use these standard tools (see other Guidance
163 documents).

Practical Guideline

Standard study design-specific tools should be used to assess RoB.

164 The terms '**applicability**', 'external validity', 'transferability', 'generalisability', and 'directness' are often
165 used interchangeably. In the context of an HTA report, it is most appropriate to use the term
166 'applicability' (12,13), although the term 'indirectness' can be chosen in the context of GRADE
167 methodology. The key question is how well the evidence matches the elements of the PICO question
168 and, therefore, whether it can be applied to answer that question (14). In statistical terms, the
169 applicability of clinical evidence threatens the overall certainty of results if, because of relevant effect
170 modification, the effect in the population of interest is probably different from the effects in the clinical
171 studies.

172 Limitations to the applicability of the evidence can occur if: (i) the study population (based on eligibility
173 criteria or actual recruitment) differs from the intended target population; (ii) the experimental or control
174 interventions were not performed in the way that they are applied or would be applied in the target
175 setting; or (iii) the study outcomes (e.g., surrogate outcomes) fail to offer information about the
176 outcomes of interest. For the applicability of clinically relevant evidence, effect modification has to be
177 taken into account (15). For example, if the relative effectiveness of a drug was shown to vary
178 substantially with age, the application of overall study results would be questionable. Instead, the
179 subgroup results for the corresponding age groups or other analytical techniques could supplement
180 information on relative effectiveness.

181 Given that applicability is usually less relevant and more straightforward to assess compared with
182 internal validity, it might be sufficient to assess any issues with regard to patients and interventions on
183 a case-by-case basis using qualitative descriptive methods. Most HTA agencies found this approach to
184 be preferable and do not use a specific instrument or checklist to judge the applicability of clinical
185 evidence (16). The applicability of a study can differ between European member states, not only
186 because PICO questions are often different, but also because of different healthcare settings (e.g.,
187 organisational aspects). Therefore, a final judgment on applicability can only be made at the national
188 (or even regional) level by each member state itself. Accordingly, the HTA Regulation (HTAR) mentions
189 that 'external validity' (i.e., applicability) should be assessed in a JCA, but without forestalling any
190 national judgement on applicability. To support national decision-making, specific issues in relation to
191 applicability should be described and addressed in a JCA, where necessary. This primarily includes
192 any potential mismatch between the PICO of interest and the PICO examined in a clinical study.
193 However, in the JCA, each aspect (e.g., questionable applicability because of differences in patient
194 population or control intervention) will only be commented on and briefly analysed, but without providing
195 a conclusion on applicability.

196 Issues with regard to surrogate outcomes usually require specific attention in HTA (17). However,
197 surrogacy is outside the scope of this Practical Guideline, and is addressed in the EUnetHTA 21
198 practical guideline *Endpoints*.

Practical guideline

Different aspects of applicability (primarily any PICO mismatch between assessment scope and clinical study) should be addressed in a JCA, but the final judgment on the applicability of study results must be left to the discretion of each member state.

199 **Statistical precision** is a quantitative concept that can be applied for each outcome of interest, at both
200 the meta and study level. Variation, and the uncertainty that comes with it, can occur in both primary
201 studies and evidence synthesis, and differentiation of both (within-study and between-study variability)
202 is required to better understand the underlying sources of variation. Effect estimates and other key
203 results should always be accompanied by the corresponding measures of statistical precision,
204 preferably confidence intervals (CIs) at a 95% level (18,19). To increase the transparency and
205 understandability of results, data submissions and analyses should contain counts and other types of
206 descriptive statistics, including raw data whenever useful.

207 In a single study, statistical hypothesis testing can be used to decide whether an effect was proven.
208 Statistical testing in a clinical study requires transparent and clear prespecification of hypotheses,
209 adequate handling of eventual multiplicity issues (20), full reporting of results (21), and careful

210 interpretation (22) (see also EUnetHTA 21 D4.5 Practical Guideline *Applicability of Evidence: Practical*
211 *guideline on Multiplicity, Subgroup, Sensitivity, and Post-hoc Analyses*). Data-driven statistical tests
212 provide results of low internal validity. Similarly, early stopping of clinical studies, deliberate extension
213 of recruitment, and selective reporting of results all undermine the validity of study results (23).
214 However, the rates of type I and type II errors in a clinical trial are not directly related to the validity of
215 the observed treatment effects, because these errors are relevant only when interpreting the results of
216 statistical tests (24).

217 Most comparative studies on interventions examine superiority hypotheses, but, depending on the
218 medical context, non-inferiority and equivalence are also tested. Although the type of question
219 (superiority, non-inferiority, or equivalence) is also important in HTA, common work on a JCA should
220 consider the rejection of the null hypothesis of a statistical hypothesis test against a prespecified α level
221 [which, in biomedical research, is usually set at 0.05 (5%)]. This neither represents nor predetermines
222 a conclusion of the added value of the assessed technology. Similarly, the clinical relevance of an effect
223 size, which can be assessed by comparing the effect size with a predefined threshold or by responder
224 analyses (25), needs to be judged at the national context. This point is addressed in the EUnetHTA 21
225 Practical Guideline *Endpoints*.

226 The certainty of a positive or negative effect will be higher if a very large effect size was found and the
227 accompanying 95% CI and p value safely exclude the possibility of no effect (26–28). Which effect sizes
228 can be considered very large and which p values can be accepted as sufficiently low is an unresolved
229 scientific question (29). Nevertheless, in the context of a JCA, it might be helpful to highlight such
230 situations, especially when no RCT evidence is available. For effect sizes expressed as relative risks,
231 the threshold of a relative risk superior to 5 (or inferior to 0.2) and a p value <0.01 (as an indicator of
232 sufficient precision) was proposed as a ‘rule of thumb’ (i.e., an arbitrary rule based on expert opinion)
233 (26,30). The JCA report will describe effect estimates, but without a conclusion on whether the certainty
234 of results is increased, because this is best made at the national level.

Practical Guideline (requirement for JCA reporting)

To describe statistical precision accurately, effect estimates should always be accompanied by the corresponding measures of variation, preferably CIs at a specified $1-\alpha$ level of confidence, which is 0.95 (95%) in most cases.

235 **3 CLINICAL STUDY DESIGNS**

236 **3.1 Terminology**

237 Classification and labelling of studies design can vary. For this Practical Guideline, we establish
238 standardised definitions for classifying and labelling clinical studies. These are used without prejudice
239 to the definitions that might be applied in national legislation and related regulatory guidance.

240 Studies are classified into two categories: **interventional studies** and **observational studies**. For
241 consistency, synonyms, such as ‘clinical trials’ or ‘experimental studies’ for interventional studies or
242 ‘non-interventional’, ‘non-experimental’, or ‘nonrandomised studies’ for observational studies are not
243 used in this Guideline.

244 Distinction between interventional and observational studies depends on whether the intervention under
245 assessment is assigned through the study protocol (interventional) or is given during routine clinical
246 care (observational).

247 **3.1.1 Interventional studies**

248 In interventional study, the intervention(s) (one or several) under assessment are assigned to
249 participants according to the study protocol.

250 Classification of interventional studies could be established based on the study characteristics. These
251 have already been fully defined but can sometimes vary (31,32). Therefore, it is valuable to establish
252 definitions of design characteristics for this Guideline. Study characteristics are summarised in Table

253 3.1, based on the glossary from the International Council for Harmonisation of Technical Requirements
 254 for Pharmaceuticals for Human Use (ICH) E9 Statistical Principles for Clinical Trials (33) or EU Clinical
 255 Trials Register (34). Note that Table 3.1 is intended to combine definitions and not to be used as a
 256 reporting template.

257 **Table 3.1. Interventional study characteristics**

Characteristic	Definition
1. Control	
Controlled (or comparative)	The study compares the effect of one or multiple treatments of interest to one or multiple comparators
2. Randomisation	
Randomised	A form of controlled allocation whereby patients are randomly assigned to one of the treatment groups
3. Blinding	
Blind	When people (patients and/or investigators and/or outcome assessors and/or statisticians) do not know which intervention is being given
4. Design	
Single arm	A trial in which all patients receive the same intervention
Parallel	Two or more interventions are evaluated concurrently in separate groups of patients
Cross-over	Comparison of two (or more) interventions in which patients are switched to the alternative treatment after a specified period (therefore, each patient receives each treatment)
Factorial	Two or more treatments are evaluated simultaneously through the use of varying combinations of those treatments
5. Objective	
Superiority	Trial with the primary objective of showing that the response to the treatment(s) of interest is clinically superior to that of a comparator
Non-inferiority	Trial with the primary objective of showing that the response to the treatment(s) of interest is not clinically inferior to that of a comparator. This is usually demonstrated by showing that the true treatment difference is unlikely to cross a threshold of an acceptable non-inferiority margin
Equivalence	Trial with the primary objective of showing that the response to two or more treatments differs by an amount that is clinically negligible. This is usually demonstrated by showing that the true treatment difference is likely to lie between a lower and an upper equivalence margin of clinically acceptable differences

258 **3.1.2 Observational studies**

259 In observational studies, there is no forced change in routine care and neither is the usual decision for
 260 intervention affected by an observational study. Given that observational studies are performed based
 261 on routine healthcare, this suggests that they allow the assessment of relative effectiveness of only
 262 those interventions that are already used in medical practice, rather than of new ones.

263 **Descriptive or analytical**

264 Observational studies can be either **descriptive**, that is, without a control group (case-series and cross-
 265 sectional studies) or **analytical** (case-control and cohort studies) with a control group. Analytical studies
 266 provide a measure of the association between exposure (notably interventions) and outcome of interest.
 267 In a case-series, changes over time can be analysed (i.e., before and after the introduction of the
 268 treatment of interest); however, under usual circumstances, such before–after changes are unlikely to
 269 assess interventional effects. It is generally important to remember that association does not necessarily
 270 imply causality. Analytical studies, such as cohort and case-control design, can be useful when
 271 randomisation is deemed unethical or unfeasible.

272 **Prospective and retrospective**

273 The collection of the data from those studies can be done prospectively or retrospectively. **Prospective**
 274 studies measure exposures before the occurrence of the outcome of interest, whereas **retrospective**
 275 studies measure exposure after the occurrence of the outcome of interest.

276 Retrospective data are usually collected from existing data sources. Thus, retrospective studies can be
 277 quicker to complete compared with prospective studies, but are limited by the availability of the existing

278 data. Furthermore, there can be a high risk of recall bias if the determination of exposure status relies
279 on recall or records only. In that case, the fundamental assumption that cause precedes effect can be
280 violated, which implies that the study of causality between exposure and outcome of interest is
281 unfeasible.

282 By contrast, prospective observational studies might be more time consuming to perform, but the patient
283 follow-up is standardised, and the availability of data that can be collected is not determined before the
284 conduct of the study. Furthermore, if a prospective study is designed to ensure that exposition precedes
285 outcome, the aforementioned fundamental assumption for causality can be assumed.

286 **Cohort study**

287 **Cohort studies**, also known as incidence studies, longitudinal studies, follow-up studies, or prospective
288 studies, are studies following a group of subjects (a cohort) with a common exposure or intervention
289 over time, but without having experienced the outcome of interest at enrolment. Patients are followed
290 during a specified period, and data on outcomes of interest are collected in a prospective manner.

291 In this Guideline, cohort studies are always considered as comparative in that a cohort study follows up
292 two or more groups from exposure to outcome. As previously mentioned, this chosen
293 classification/definition is used without prejudice to definitions that might be applied elsewhere.

294 Sometimes, a cohort study data set can serve as a basis for enrolling patients into an interventional
295 study, which can be a RCT (i.e., a subset of newly included or already-included patients can be allocated
296 to one of the interventions assessed if, at a proper time, they meet the eligibility criteria for the
297 interventional study). When it happens, this design is called a 'trial within a cohort' (TWIC) (35).

298 **Case-control study**

299 **Case-control studies** are retrospective studies that enroll patients who have experienced a particular
300 outcome of interest ('cases'), compared with patients who have not experienced the outcome of interest
301 but who are representative of the study population on some controlled criterion ('controls').

302 The aim of this study design is to compare the exposure between case and controls to identify factors
303 that might contribute to be associated to the occurrence of an outcome.

304 **Cross-sectional study**

305 **Cross-sectional studies**, also known as transversal studies, measure outcomes and exposure status
306 simultaneously in a specified population to study the frequency and characteristics of an outcome at a
307 particular point in time.

308 The main use of this study design is to assess outcome and/or exposure prevalence in a population.

309 **Case study: case-report and case-series**

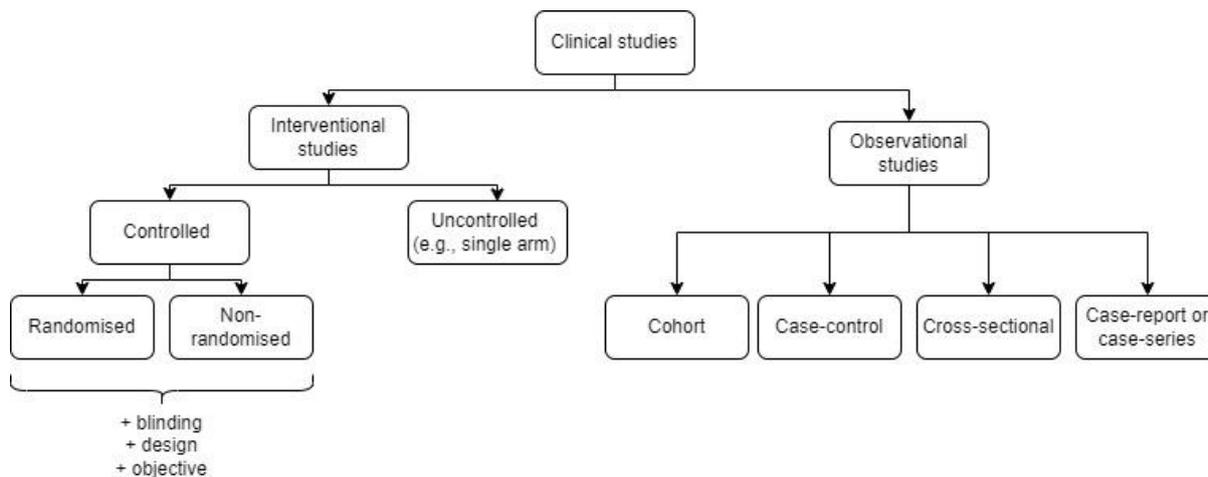
310 Case studies are descriptive studies of a single case (**case-report**) or a group of subjects with similar
311 diagnoses or intervention (**case-series**) followed over time. It provides detailed descriptions of cases
312 without the use of a control group. However, in a case-series, it is possible to compare the health status
313 of participants over time, for example, to estimate the pre–post changes induced by an intervention.
314 Given the characteristics of this design, such changes are unlikely to estimate the true effect of the
315 treatment of interest.

316 Case studies can be used to describe rare events or early trends, such as unusual manifestations of a
317 disease or unusual response to an exposure. Some case-reports in the medical literature are intended
318 to prove the feasibility of an intervention. Those study designs can not be used to assess the
319 effectiveness of an intervention. However, they can help to detect new safety signals.

320 **3.2 Classification**

321 The classification of clinical studies is presented in Figure 3.1.

322 **Figure 3.1. Classification of clinical studies**
 323



324

Practical Guideline (requirement for JCA reporting)

Classification and design characteristics for each study submitted as evidence.

325 As per its definition, this classification is used in this Guideline without prejudice to the definitions that
 326 might be applied elsewhere (36,37).

327 **4 SPECIFIC STRENGTHS, WEAKNESSES, AND RECOMMENDATIONS**
 328 **REGARDING DIFFERENT DESIGNS**

329 The JCA will report the certainty of results of the relative effectiveness of the treatment(s) of interest,
 330 taking into account the strengths and limitations of the available evidence [Article 9(1)]. As previously
 331 described, the certainty of results is determined by internal validity, applicability, and statistical
 332 precision.

333 Study design or conduct can lead to bias, impacting internal validity. Several standardised tools have
 334 been developed to evaluate RoB in various clinical study designs (13,14). They are helpful for assessing
 335 the strengths and limitations of the available evidence and should be used when performing JCA.

336 Therefore, this Guideline recommends the systematic use of Cochrane’s tools to assess RoB.

337 **4.1 Randomised clinical trials: gold standard**

338 RCTs are the gold standard for evaluating causal relationships between interventions and outcomes
 339 because randomisation eliminates much of the bias inherent to other designs (38). In brief, a proper
 340 randomisation allows the trial to be conducted under the assumption of **exchangeability** (i.e., if patients
 341 from one group were substituted to the other, the same treatment effect would be observed). This
 342 underlying assumption implies the absence of confounding bias (both on known and unknown
 343 confounders and effect modifiers). Moreover, blinding alongside with identical and standardised follow-
 344 up between each group help to maintain exchangeability over time and prevent measurement bias. As
 345 a result of randomisation and blinding, relative effectiveness assessment allows estimation of the
 346 supplementary causal effect of an intervention of interest over comparator treatment effects. Finally,
 347 rigorous follow-up and analysis of the adequate population (e.g., intention-to-treat population for a
 348 superiority RCT) help control attrition. Nonetheless, depending on numerous factors, such as the quality
 349 of the design and conduct of the study, the certainty of results of a particular RCT can be questioned
 350 and biases can arise (39,40).

351 To allow proper evaluation by member states, RoB should be assessed using **ROB-2** (10). Full
 352 guidance documents for ROB-2 could be found on the Cochrane resource website

353 (<https://methods.cochrane.org/risk-bias-2>). Given that ROB-2 assumes that overall RoB is performed
354 at the outcome level, RoB should be performed for every outcome required in the assessment scope
355 (i.e., 'O', from PICO). Indeed, although the occurrence of some biases can frequently impact internal
356 validity at the study level (i.e., irrespective of the outcomes being assessed), other biases can be
357 outcome dependent (e.g., nonblinding can affect the assessment of outcomes differently, such as
358 overall survival and quality of life measured by patient-reported outcomes). Moreover, results in the
359 JCA report will be presented primarily according to PICO questions. Therefore, it appears valuable to
360 have a RoB assessment at the outcome level.

Practical Guideline (Requirement for JCA reporting)

For outcomes with evidence coming from RCTs, assess RoB using ROB-2.

RoB should be assessed for each outcome required in the assessment scope.

The RoB of outcomes sharing the same characteristics (e.g., data collection, blinding aspects, or evaluation) can be assessed on an overall level (i.e., grouped).

RoB judgement should be provided for both each individual domain level and overall.

361 **4.2 Nonrandomised controlled trials**

362 Non-RCTs are clinical trials in which participants are allocated to intervention under assessment or
363 reference intervention using methods that are not random. Allocation could be based, for example, on
364 investigator's choice, participant's choice, or calendar dates. They allow direct estimation of relative
365 effects between interventions. However, such nonrandom allocation breaks the underlying assumption
366 of exchangeability and, therefore, is likely lead to confounding bias. Thus, the estimated association
367 between intervention and outcome is likely to be biased and will differ from its true causal effect.

368 There are different methods that can be used to control for confounding (i.e., allowing if properly
369 conducted, **conditional exchangeability**, e.g., design-based methods, such as stratification or
370 matching, or modeling-based methods, such as adjustment or models of causal inference (e.g.,
371 propensity scores or g-computation)) within the trial. However, although known and unknown, measured
372 and unmeasured confounding factors and effect modifiers are fully controlled through randomisation,
373 any other method for controlling confounding bias when allocation was not randomised requires
374 exhaustivity (i.e., all relevant confounders and effect modifiers must be known and adequately
375 measured within the trial), an unverifiable underlying assumption.

376 To allow proper evaluation by Member States, RoB should be assessed using **ROBINS-I**. Full guidance
377 documents for ROBINS-I could be found here using the Cochrane resource website
378 (<https://sites.google.com/site/riskofbiastool/welcome/home/current-version-of-robins-i>). As for ROB-2,
379 RoB assessment using ROBINS-I must be performed at the outcome level.

Practical Guideline (Requirement for JCA reporting)

For outcomes with evidence coming from non-RCTs, assess RoB using ROBINS-I.

RoB should be assessed for each outcome required in the assessment scope.

The RoB of outcomes sharing the same characteristics (e.g., data collection, blinding aspects, or evaluation) can be assessed on an overall level (i.e., grouped).

RoB judgement should be provided for both each individual domain level and overall.

380 **4.3 Uncontrolled clinical trials (e.g., single-arm trials)**

381 Unlike comparative clinical trials, uncontrolled trials do not allow relative effectiveness assessment (i.e.,
382 supplementary effect over comparator treatment effect). In terms of strengths and weaknesses, they
383 can be considered mostly akin to case-series. However, a difference with case-series is that the
384 treatment is delivered as part of a study intervention. Therefore, patients in a single-arm trial can receive
385 a treatment in a more-standardised manner and with a more-rigorous follow-up compared with those
386 from a case-series. In the context of HTA, uncontrolled clinical trials are of very limited value for
387 estimating treatment effectiveness.

388 Given the lower importance of uncontrolled trials for relative effectiveness assessment and HTA, it is
389 deemed unnecessary to propose any formal rules for assessing RoB of single-arm trials. Some tools
390 have been developed in the past (41–44), but RoB of uncontrolled studies appears to be affected by
391 only a few specific aspects of internal validity, such as the consecutiveness of recruitment, the
392 prespecification of sample size and analyses, and the blinded assessment of outcomes. Nevertheless,
393 RoB of an uncontrolled study is very unlikely to be changed by formal RoB assessment; thus, this work
394 appears dispensable.

Practical Guideline

Evidence coming from uncontrolled trials is of very limited value for performing relative effectiveness assessment.

Although the (partial) use of some tools for RoB assessment is possible, the overall conclusion on the (very limited) internal validity of uncontrolled studies is very unlikely to be changed by RoB assessment. Therefore, RoB assessment is not required.

395 **4.4 Cohort studies**

396 Cohort studies can be used when allocation of an intervention in a controlled manner is deemed
397 unethical or unfeasible. Compared with interventional studies, they can allow larger sample sizes and
398 longer follow-up, improving statistical precision or the detection of long-term adverse events (45). They
399 can also help to investigate the effectiveness of interventions when used in routine healthcare on a
400 sample of patients with less-stringent eligibility criteria compared with an interventional study, which
401 could enhance applicability.

402 Given that the intervention is not randomised between patients, the underlying assumption of
403 exchangeability cannot hold, which is very likely to lead to confounding bias. Thus, without the proper
404 use of an appropriate method for controlling for confounding (see Section 4.2), the estimated
405 association between exposure and outcome of interest will most likely differ from its true causal effect.

Practical Guideline (Requirement for JCA reporting)

For outcomes with evidence coming from cohort studies, assess RoB using ROBINS-I.

RoB should be assessed for each outcome required in the assessment scope.

The RoB of outcomes sharing the same characteristics (e.g., data collection, blinding aspects, or evaluation) can be assessed on an overall level (i.e., grouped).

RoB judgement should be provided for both each individual domain level and overall.

406 **4.5 Case-control studies**

407 A case-control study design is useful to examine rare outcomes (e.g., in rare diseases), and multiple
408 factors affecting one outcome can be studied.

409 In case-control studies, patients are enrolled based on the occurrence of outcome and exposures are
410 investigated in a retrospective manner. Thus, they are at high risk of selection bias. The selection of a
411 control group is very likely to not allow verification of the exchangeability assumption. It leads to the
412 same issues as described before for non-RCTs and cohort studies regarding confounding bias (see
413 Section 4.2). Moreover, case-control studies are also likely to lead to a measurement bias (e.g., recall
414 bias), because exposure is measured after the onset of the disease or outcome. Moreover, because
415 data are collected in a retrospective manner, it is uncertain that the exposure of interest precedes the
416 occurrence of the outcome of interest, which can lead to violation of a fundamental rule of causation
417 (exposure must precede effect).

418 Finally, this study design is not suited for rare exposures and for studying more than one outcome.

Practical Guideline (Requirement for JCA reporting)

420 For each outcome with evidence coming from a case-control study, assess RoB using ROBINS-I.

421 RoB should be assessed for each outcome required in the assessment scope.

422 The RoB of outcomes sharing the same characteristics (e.g., data collection, blinding aspects, or
423 evaluation) can be assessed on an overall level (i.e., grouped).

424 RoB judgement should be provided for both each individual domain level and overall.

425 **4.6 Cross-sectional studies**

426 A cross-sectional study design is useful to investigate multiple outcomes and exposures
427 simultaneously.

428 This type of study estimates association but cannot be used to study the cause–effect relationship or
429 causality because there is no temporality; thus, it is not possible to distinguish whether the exposure
430 preceded or followed the outcome. Therefore, it is deemed unnecessary to propose any formal tool for
431 assessing RoB of cross-sectional studies.

432 **Practical guideline**

433 Evidence coming from cross-sectional studies is of very limited value for performing relative
434 effectiveness assessment.

435 No RoB assessment using a standardised tool is required for cross-sectional studies.

436 **4.7 Case-series and case-reports**

437 These studies allow the generation of hypotheses, such as identifying unexpected effects (adverse or
438 beneficial) and describing unusual syndromes that could later be studied using study designs with a
439 higher certainty of results.

440 These studies are only descriptive and are rarely used to test hypotheses or establish causal effects.
441 Any effect estimate generated from a study lacking a control group is only a pre–post change, thus the
442 interpretation of such change as a causal effect requires the very unlikely assumption that no change
443 would have occurred without the intervention. Furthermore, case-reports generate selection bias and
444 lack external validity because of low representativeness. Therefore, it is deemed unnecessary to
445 propose any formal tool for assessing RoB of case-series and case-reports.

446 **Practical Guideline**

447 Evidence coming from case-series and case-reports is of very limited value for performing relative
448 effectiveness assessment.

449 No risk of bias assessment using a standardised tool is required for case-series and case-reports.

450 **5 PARTICULARITIES**

451 Different specificities will be introduced in this section. Indeed, although specificities are methodological
452 concepts that are now prevalent when discussing the design of clinical studies, they cannot be strictly
453 classified according to the principles described earlier in the document (see Section 3). These
454 particularities can be compatible with many features of the aforementioned designs (e.g., some can be
455 compatible with the principles of RCTs). Nonetheless, their definitions, strengths, and weaknesses need
456 to be highlighted separately because they can justify looking for specific methodological points of
457 attention.

458 **5.1 Master protocols**

459 ‘**Master protocol**’ refers to the use of an overarching logistic, design protocol allowing the investigation
460 of multiple hypotheses or interventions in one or multiple diseases (46,47). The master protocol
461 proposes a common infrastructure establishing uniformity and standardisation of procedures in
462 designing and assessing different interventions. Usually, the concept of a master protocol encompasses
463 three subtypes: **platform trials** [also called **multi-arm, multi-stage trials (MAMS)**], **basket trials**, and
464 **umbrella trials** (48).

465 5.1.1 Platform trials

466 Platform trials allow, for a particular disease, the comparison, either simultaneously and/or sequentially,
467 of multiple interventions with a common control group (48). Sometimes the different interventions can
468 also be compared with each other. The master protocol defines the overall infrastructure and sets the
469 overarching principles of the design, but specific addendum protocols are created when a new
470 intervention is assessed. Given that the assessment of certain interventions can be stopped or,
471 alternatively, added to the trial, platform trials can be considered mainly as *adaptive trials* (49). The
472 intervention that is used as a control can also evolve over time if the standard of care is updated
473 following the start of the platform trial. Platform trials are mainly phase 3 RCTs (i.e., a confirmatory
474 assessment of effectiveness), but they sometimes start as phase 2 trials (i.e., an exploratory
475 assessment of effectiveness, which can be uncontrolled), and the switch from phase 2 to phase 3 is
476 conducted under the same master protocol (this is sometimes called a 'seamless' design) (47). In that
477 case, the most promising interventions based on the results of the phase 2 trial are retained for the
478 phase 3 trial. Therefore, the follow-up of some patients from a phase 2 trial can be extended to the
479 phase 3 trial (providing they meet the phase 3 eligibility criteria).

480 Methodologically, the main strength of platform trials is their adaptive nature. Thus, they can be
481 considered as more 'disease focused' compared with more commonly used traditional RCTs because
482 they can provide a more efficient assessment of multiple interventions in a manner that can be
483 potentially perpetual with the possibility to be adapted to both scientific discoveries provided by the trial
484 and external discoveries (48). Thus, platform trials offer the potential to generate comparative evidence
485 for multiple treatments that are simultaneously in clinical development and could reduce the need to
486 use indirect comparison methods, such as network meta-analyses, for assessing the relative
487 effectiveness of multiple interventions.

488 In itself, platform trials are not a specific type of methodological design *per se*. Therefore, platform trials
489 can provide the same certainty of results as well-performed RCTs providing they are conducted in
490 conformity with the same methodological principles. Nonetheless, because of their flexibility, several
491 specific points of attention must be considered. First, platform trials can sometimes start as phase 2
492 trials. Thus, it is important that the criteria to select interventions that are going to phase 3 are clearly
493 defined (e.g., the criteria for defining sufficient presumption of effectiveness). Moreover, because
494 patients from phase 2 can participate in phase 3 of the trial, it is necessary that these patients still meet
495 the eligibility criteria for phase 3. Second, because the inclusion of new patients in the control group
496 can occur over long periods, the contemporaneity of the control group in relation to the assessment of
497 some interventions can be brought into question and the relevance of the intervention proposed within
498 the control group must be scrutinised. Third, although blinding of patients and investigators is possible,
499 it requires the use of multiple dummies, which can be difficult to achieve when there are multiple
500 treatments with different pharmaceutical formulations that are assessed simultaneously. Thus,
501 numerous platform trials are conducted in an open manner. Fourth, multiple interim analyses are usually
502 performed as well as multiple comparisons between groups. Thus, there is a risk of an inflated type 1
503 error rate if not properly managed. Therefore, assessment of the quality of these analyses (interim
504 analyses and multiples groups comparisons) should follow the guidelines proposed in the EUnetHTA
505 Practical Guideline *Applicability of Evidence: Practical Guideline on Multiplicity, Subgroup, Sensitivity
506 and Post-hoc Analyses*. Finally, it is imperative that the rules for adding new interventions into the trials
507 are explicit and justified.

Practical Guideline (Requirement for JCA reporting)

If a platform trial is an RCT then, in general, RoB needs to be assessed according to the principles described in Section 4.2.

If a platform trial is not an RCT then, in general, RoB needs to be assessed according to the principles described in the Section herein corresponding to the design of the study.

Specific points for attention

If the platform trial starts as a phase 2 trial, then the rules to select interventions that are going to phase 3.

If the platform trial starts as a phase 2 trial, do the patients that were retained from phase 2 to phase 3 meet the eligibility criteria for phase 3? (Yes/No)

Design considerations when adding new intervention(s) (criteria, process, or timing) to the trial.

The potential modifications of the intervention of the control group.

The results of interim analyses and multiple comparisons in accordance with the EUnetHTA practical guideline *Applicability of Evidence: Practical Guideline on Multiplicity, Subgroup, Sensitivity and Post-hoc Analyses*.

508 **5.1.2 Basket trials**

509 Basket trials aim to assess a targeted intervention across multiple diseases (47,50). Eligibility of patients
510 is based on a unifying criterion, which is a specific mechanism of action of the treatment of interest with
511 prognostic and/or predictive value (e.g., a specific molecular alteration or a common pathological
512 process). Therefore, the targeted intervention is supposed to produce a beneficial effect for all patients
513 because it targets a common process. Therefore, basket trials pool patients with diseases that are
514 classified as different in terms of usual nosography (e.g., cancers from different primary organs or
515 different cardiovascular diseases). Basket trials are mainly used in oncology for assessing the
516 effectiveness of interventions designed to target specific molecular alterations, but other medical areas
517 can be concerned by the use of such trials (47,50).

518 The main strength of basket trials is their potential ability to generate evidence of effectiveness
519 regarding interventions targeting a specific risk factor with prognostic and/or predictive value, therefore
520 generating evidence for multiple diseases in one trial (50). Nonetheless, the ability of basket trials to
521 provide such certainty of results relies on multiple assumptions and conditions (51).

522 In itself, a basket trial is not a type of methodological design *per se*. Therefore, the certainty of results
523 provided by such a trial is mainly dependent on its design. Although basket trials can be RCTs, most
524 are currently uncontrolled trials and, therefore, do not provide a higher certainty of results compared
525 with single-arm trials (47). Randomisation and relative effectiveness assessment in the context of
526 basket trials can be difficult because they investigate multiple diseases and, therefore, can require
527 multiple control interventions (47). Second, the hypothesis that the effect of the targeted intervention
528 will be beneficial, on average, for each 'cohort' of patients (e.g., the first cohort is patients with breast
529 cancer, the second one is patients with lung cancer, etc.) relies on the assumption of homogeneity of
530 between-cohorts effects (51). This assumption cannot be proven by analysing the data of the conducted
531 basket trial. Indeed, there is the possibility of performing an interaction statistical hypothesis test
532 between the intervention and the different cohorts (51). However, even if the test does not reject the
533 null hypothesis of homogeneity of effects, it does not experimentally prove homogeneity because the
534 test can be nonsignificant as a result of a lack of power. Thus, the plausibility of this assumption relies
535 mainly on the basis of biological arguments of the mechanisms of actions or on the proximity to other
536 situations in which the hypothesis of homogeneity has been accepted or proven. Third, the specific
537 effect of the targeted intervention in a specific 'cohort' (e.g., patients with breast cancer only) can suffer
538 from a lack of statistical precision because it can be expected that some cohorts will have a low number
539 of patients given that the occurrence of the targeted risk factor can be rare. Finally, eligibility criteria
540 often rely on the screening of a specific molecular alteration or biomarker. Therefore, inclusion within a
541 basket trial often relies on the results of a companion test and, therefore, the performance of the test
542 (sensitivity, specificity, predictive values, or probability reports, calibration, and discriminatory capacity
543 for biomarkers measured on a continuum) must be known and must be of an acceptable level (51).
544 Moreover, the test must be available for all potentially eligible patients.

Practical Guideline (Requirement for JCA reporting)

If a basket trial is an RCT then, in general, RoB needs to be assessed according to the principles described in Section 4.2.

If a basket trial is not an RCT then, in general, RoB needs to be assessed according to the principles described in the Section corresponding to the design of the study.

Specific points for attention

Rationale for the plausibility of the hypothesis of homogeneity of effects.

If the eligibility of patients within the basket trial relies on the results of a companion test, its performance, availability, and methods used for detection (e.g., on which tumor sample the test is performed).

If an interaction test for homogeneity of effect was performed, its method and result.

Results of effectiveness within each 'cohort' of patients with appropriate statistical estimates.

545 **5.1.3 Umbrella trials**

546 Umbrella trials, which are also mostly used in oncology, aim to assess multiple targeted interventions
547 for what is considered a single disease according to usual nosography (47,50). Patients with a single
548 disease are included (e.g., advanced breast cancer) and are stratified into subgroups based on a
549 biomarker or risk factor with a prognostic and/or predictive value. Thus, the single disease is split into
550 multiple subtypes with eligibility for each intervention group defined by the mechanism of action of each
551 treatment. Each intervention group receives a different targeted intervention that is supposed to have a
552 beneficial effect that is better suited for the specific subgroup of patients for which it is proposed.

553 The main strength of umbrella trials is their ability to propose targeted therapies that have the potential
554 to be better suited for subgroups of patients of a same disease, which can ultimately enhance the
555 development of stratified medicine (48).

556 As for any other types of master protocol, umbrella trials are not a type of methodological design *per*
557 *se*. Therefore, the certainty of results provided by an umbrella trial is mainly dependent on its design.
558 Akin to basket trials, although umbrella trials can be RCTs, most are currently uncontrolled trials and,
559 therefore, do not provide a higher certainty of results compared with single-arm trials (47). Nonetheless,
560 randomisation and relative effectiveness assessment can be considered easier to achieve in the context
561 of an umbrella trial compared with a basket trial, because the existing standard of care (or placebo, if
562 there is no established care) for the disease being studied can be used as a common control for all the
563 subgroups (47). As for basket trials, inclusion often relies on the search of a specific molecular alteration
564 or biomarker. Therefore, the performance of the companion test in allowing the detection and
565 quantification of the prognostic and/or predictive factor of interest must be known and of an acceptable
566 level, and the test available for all potentially eligible patients (see Section 5.1.2).

Practical Guideline (Requirements for JCA reporting)

If an umbrella trial is an RCT then, in general, RoB needs to be assessed according to the principles described in Section 4.2.

If an umbrella trial is not an RCT then, in general, RoB needs to be assessed according to the principles described in the Section corresponding to the design of the study.

Specific points of attention

If the eligibility of patients within the umbrella trial relies on the results of a companion test, its performances, availability, and methods used for detection (e.g., on which tumor sample the test is performed).

567 **5.2 Real-world data and real-world evidence**

568 **Real-world data (RWD)** is an umbrella term encompassing the use of various types of data that share
569 the common property they have been generated in the context of routine healthcare [e.g., electronic
570 health records, medical claims and billing data, administrative healthcare databases, patient-generated
571 data (including in-home-use settings) and data produced from various sources (such as electronic
572 devices) that can inform on health status] (52–54). Therefore, the term excludes data collected explicitly
573 for research purposes. In relation to the concept of RWD, **real-world evidence (RWE)** is a term defining
574 clinical evidence of a health technology derived from the analysis of RWD for a given research question.
575 RWD can be used to generate RWE for different purposes: generating hypotheses for testing in future
576 RCTs, assessing trial feasibility, informing prior probability distributions for Bayesian statistical models,
577 identifying patient baseline characteristics or prognostic and predictive factors, describe usage of a
578 health technology in real-world setting, and assessing the effectiveness and/or safety of health
579 technologies (e.g., for new indications of already-used technologies or for documenting long-term
580 follow-up).

581 Although 'RWD' is used to describe data generated in the context of routine healthcare, such data can
582 be used for various purposes in the context of clinical research. Thus, RWD can be coupled with data
583 generated for clinical research purposes. Indeed, a specific source of RWD can be used as a basis for

584 conducting a RCT in which the collection of necessary data can exclusively come from a set of RWD,
585 or as a primary source complemented by data specifically collected for the clinical study (i.e., a
586 secondary source). These types of study are sometimes considered part of what are called 'pragmatic
587 trials' (55). When only a subset of newly included patients within the collection of a specific RWD (e.g.,
588 a cohort of patients with data collected from electronic health records) are randomised over time, the
589 corresponding RCT can be considered a *TWIC* (35). When the secondary source of data is collected
590 using fully remote pathways (e.g., electronic informed consent, digital assessment tools, or virtual study
591 visits), the corresponding RCT is sometimes called a '*contactless trial*' (56). RWD can also be used as
592 the only or as the primary source of data for any type of other clinical trial (e.g., single-arm trial) or
593 observational studies (e.g., cohort study). Although this is out of the scope of this Guideline, they can
594 be used as sources of data for indirect comparisons (see the EUnetHTA 21 Methodological Guideline
595 *Direct and Indirect Comparisons*), or as additional historical data borrowing for enriching data of a
596 control group in an already existing clinical trial (e.g., when the trial concerns a rare disease).

597 The use of RWD in generating evidence can be useful in multiple ways. First, their use can enhance
598 the recruitment of patients in clinical trials, especially for rare diseases (57). Second, their use can
599 enhance the level of applicability of evidence (or external validity) and/or the level of statistical precision
600 by facilitating the conduct of clinical studies on large sample of patients with less stringent inclusion
601 criteria compared with a classical clinical trial, by assessing the effectiveness and/or safety of health
602 technologies in 'real-world' settings, and by allowing studies with clinical trials with a longer follow-up
603 than usual (55).

604 Potential weaknesses in using RWD when conducting clinical studies are mainly linked to the fact that
605 a set of RWD was not primarily structured for conducting a clinical study. Thus, data validity, data
606 integrity, and data monitoring are dependent on the quality of already-existing procedures before the
607 conduct of a given clinical study (58). A related issue can be the use of certain variables from databases
608 as proxies of the characteristics they are supposed to measure in a given clinical study, which can lead
609 to measurement bias (59). For example, data about the dispensation of pharmaceutical drugs coming
610 from administrative databases can be used as a proxy for usage even though the two concepts are not
611 equivalent (even if correlated). Second, follow-up of patients included in a clinical study using RWD
612 might not be as standardised as in *de novo* clinical studies (especially if RWD are the only source of
613 data that will be used for analysis), which can result in a greater risk of attrition bias (58). Finally,
614 particular attention to the assessment of endpoints is required because there is a risk of unblinded
615 and/or decentralised adjudication of endpoint processes (60).

616 To conclude, in itself, RWD does not define a type of clinical study design and RWE can be produced
617 with varying certainty of results for a given research question. Therefore, the certainty of results that is
618 produced, especially the level of internal validity, is mainly determined by the study design of a given
619 clinical study based on the use of RWD. Especially because most clinical studies using RWD are
620 currently not RCTs, controlling for confounding bias is one of the main issues when estimating treatment
621 effectiveness. Indeed, the lack of randomisation requires the proper use of methods to control for
622 confounding bias (see Section 4.2), which rely on assumptions (e.g., the assumption of exhaustivity on
623 confounders and effect modifiers) that are, in part, unverifiable.

Practical Guideline

RWD are data generated in routine healthcare and are not primarily structured for clinical studies.

RWD can be analysed to generate RWE for many purposes, including the assessment of effectiveness and/or safety of health technologies.

Any study design can be implemented using RWD from RCT to observational studies.

RWD can enhance recruitment, applicability, and statistical precision, and can lead to longer follow-up than classical clinical trials.

Specific points of attentions

For a given clinical study, it should be reported if RWD are the sole source of data, or a primary source of data complemented by a secondary source specifically collected for research purposes (and, if so, to which specific design it corresponds).

Given the at least partial use of data that were not initially structured for clinical research, the validity and reliability of RWD for adequately answering a given research question is of particular importance, especially the potential use of proxy variables, the risk of attrition bias, and the adequate measurement of endpoints.

Requirement for JCA reporting

RWD is not a design *per se*; thus, the design of a clinical study should be described and classified according to the principles already described in this Guideline.

RoB should be assessed according to the principles already described in this Guideline.

624 **5.3 Registries**

625 Clinical registries are organised systems collecting data on a group of patients defined by a common
626 characteristic or set of characteristics, which can be the occurrence of a particular disease, condition,
627 exposure or use of a particular health technology or health-related service (61,62). After inclusion of a
628 patient into the registry, follow-up data (i.e., outcomes) are collected. Data collected within the registry
629 can then be used to conduct registry-based studies. Given that they are often a collection of
630 observational data from routine healthcare practices, data from registries can be considered as RWD
631 (54), but it could be advocated that some registries are organised systems that are explicitly devoted to
632 research purposes. Nevertheless, registry data can be used in the same way (e.g., as the sole source
633 of data or as a primary source of data) and for as many purposes as RWD. Furthermore, RCTs
634 conducted using, exclusively or in part, data from registries are often called 'registry-based RCTs
635 (60,63).

636 The strengths that were outlined for RWD-based clinical studies can be found in registry-based studies
637 (64). A particular point that can sometimes apply is the fact some registries aim toward exhaustivity.
638 This means that they aim to include the entire population of interest of patients presenting the
639 characteristic leading to their inclusion in the registry (e.g., the diagnosis of a particular disease) within
640 the boundaries of a specific geographical area (which can sometimes be at a national level). Therefore,
641 some registry-based studies can have the ability to produce true statistics in a population of interest
642 and not estimates; therefore, external validity of a registry-based study can be better than of a study
643 conducted on a sample of patients only (providing the target population of the corresponding clinical
644 study is the same as the population covered by the registry).

645 However, many weaknesses identified for RWD-based clinical studies are also present in registry-
646 based studies, but some of these aforementioned weaknesses can be mitigated depending on the
647 context. Indeed, first, registries are sometimes built around the idea of answering specific research
648 questions. Thus, registries can produce data with a structure that is more adequately suited to answer
649 specific research questions compared with other sources of RWD. Second, data validity, integrity, and
650 monitoring can be primary concerns in well-structured registries (e.g., national-level registries for a
651 particular disease) and, thus, registry-based studies can sometimes profit from data with a higher level
652 of quality regarding these aspects compared with other types of RWD, especially regarding attrition
653 bias. However, registry data should not be automatically assumed as presenting a high level of validity
654 and reliability and procedures for collection and monitoring of data should be scrutinised anyway when
655 assessing the validity of a registry-based study. Finally, the same remark can be made as for RWD:
656 registry data, in themselves, do not define a clinical study design. Therefore, certainty of results that is
657 produced using registry data, especially the level of internal validity, is mainly determined by the design
658 of a given registry-based clinical study (65).

Practical guideline

Clinical registries are organised systems collecting data information on patients based on a common characteristic (e.g., diagnosis of a disease or use of a particular health technology).

Any study design can be implemented using registries, from RCTs to observational studies.

Some registries aim to achieve exhaustivity for a population of interest in a specific geographical area. Therefore, a registry-based study can produce true statistics providing that the target population is the same as that covered by the registry.

Specific points for attention

For a given clinical study, it should be reported if a registry is the sole source of data, or a primary source of data complemented by a secondary source specifically collected for research purposes (and, if so, to which specific design it corresponds).

Requirement for JCA reporting

Registries are not a design *per se*; thus, the design of a clinical study should be described and classified according to the principles already described in this Guideline.

RoB should be assessed according to the principles already described in this Guideline.

659 6 REFERENCES

660. Guyatt GH, Oxman AD, Vist GE, et al. GRADE: an emerging consensus on rating quality of evidence
661 and strength of recommendations. *BMJ*. 2008;336(7650):924–6.
662. EUnetHTA. *Partial use of GRADE in EUnetHTA Framework*. [www.eunetha.eu/wp-](http://www.eunetha.eu/wp-content/uploads/2021/05/EUnetHTA-GRADE-framework-paper.pdf?x16454)
663 [content/uploads/2021/05/EUnetHTA-GRADE-framework-paper.pdf?x16454](http://www.eunetha.eu/wp-content/uploads/2021/05/EUnetHTA-GRADE-framework-paper.pdf?x16454) (accessed 22 Jun 2022).
664. Whiting PF, Rutjes AWS, Westwood ME, et al. QUADAS-2 Steering Group. A systematic review
665 classifies sources of bias and variation in diagnostic test accuracy studies. *J Clin Epidemiol*.
666 2013;66(10):1093–104.
667. EU. *Regulation (EU) 2021/2282 of the European Parliament and of the Council of 15 December 2021*
668 *on health technology assessment and amending Directive 2011/24/EU (Text with EEA relevance)*.
669 <http://data.europa.eu/eli/reg/2021/2282/oj/eng> (accessed 22 Jun 2022).
670. Fletcher GS. *Clinical Epidemiology: The Essentials*. Philadelphia: Lippincott Williams & Wilkins; 2020.
671. Grimes DA, Schulz KF. Bias and causal associations in observational research. *Lancet*.
672 2002;359(9302):248–52.
673. OCEBM. *OCEBM Levels of Evidence*. [www.cebm.ox.ac.uk/resources/levels-of-evidence/ocebml-](http://www.cebm.ox.ac.uk/resources/levels-of-evidence/ocebml-levels-of-evidence)
674 [levels-of-evidence](http://www.cebm.ox.ac.uk/resources/levels-of-evidence/ocebml-levels-of-evidence) (accessed 22 Jun 2022).
675. Atkins D, Eccles M, Flottorp S, et al. Systems for grading the quality of evidence and the strength of
676 recommendations I: critical appraisal of existing approaches The GRADE Working Group. *BMC Health*
677 *Serv Res*. 2004;4(1):38.
678. Djulbegovic B, Guyatt GH. Progress in evidence-based medicine: a quarter century on. *Lancet*.
679 2017;390(10092):415–23.
680. Sterne JAC, Savović J, Page MJ, et al. RoB 2: a revised tool for assessing risk of bias in randomised
681 trials. *BMJ*. 2019;366:l4898.
682. Sterne JA, Hernán MA, Reeves BC, et al. ROBINS-I: a tool for assessing risk of bias in non-randomised
683 studies of interventions. *BMJ*. 2016;355:i4919.
684. Murad MH, Katabi A, Benkhadra R, et al. External validity, generalisability, applicability and directness:
685 a brief primer. *BMJ Evid-Based Med*. 2018;23(1):17–9.
686. EUnetHTA. *Levels of Evidence - applicability of evidence for the context of a relative effectiveness*
687 *assessment Amended JA1 Guideline Final*. [www.eunetha.eu/levels-of-evidence-applicability-of-](http://www.eunetha.eu/levels-of-evidence-applicability-of-evidence-for-the-context-of-a-relative-effectiveness-assessment-amended-ja1-guideline-final-nov-2015/)
688 [evidence-for-the-context-of-a-relative-effectiveness-assessment-amended-ja1-guideline-final-nov-](http://www.eunetha.eu/levels-of-evidence-applicability-of-evidence-for-the-context-of-a-relative-effectiveness-assessment-amended-ja1-guideline-final-nov-2015/)
689 [2015/](http://www.eunetha.eu/levels-of-evidence-applicability-of-evidence-for-the-context-of-a-relative-effectiveness-assessment-amended-ja1-guideline-final-nov-2015/) (accessed 22 Jun 2022).
690. Windeler J. [External validity]. *Z Evidenz Fortbild Qual Im Gesundheitswesen*. 2008;102(4):253–9.
691. Rothwell PM. External validity of randomised controlled trials: ‘to whom do the results of this trial apply?’
692 *Lancet*. 2005;365(9453):82–93.
693. EUnetHTA JA2. *Levels of evidence: applicability of evidence for the context of a relative effectiveness*
694 *assessment*. [www.eunetha.eu/wp-content/uploads/2018/01/Levels-of-Evidence-Applicability-of-](http://www.eunetha.eu/wp-content/uploads/2018/01/Levels-of-Evidence-Applicability-of-evidence-for-the-context-of-a-relative-effectiveness-assessment_Amended-JA1-Guideline_Final-Nov-2015.pdf)
695 [evidence-for-the-context-of-a-relative-effectiveness-assessment_Amended-JA1-Guideline_Final-Nov-](http://www.eunetha.eu/wp-content/uploads/2018/01/Levels-of-Evidence-Applicability-of-evidence-for-the-context-of-a-relative-effectiveness-assessment_Amended-JA1-Guideline_Final-Nov-2015.pdf)
696 [2015.pdf](http://www.eunetha.eu/wp-content/uploads/2018/01/Levels-of-Evidence-Applicability-of-evidence-for-the-context-of-a-relative-effectiveness-assessment_Amended-JA1-Guideline_Final-Nov-2015.pdf) (accessed 22 Jun 2022).
697. EUnetHTA. *Guideline: endpoints used in relative effectiveness assessment of pharmaceuticals -*
698 *surrogate endpoints*. www.eunetha.eu/wp-content/uploads/2018/01/Surrogate-Endpoints.pdf
699 (accessed 22 Jun 2022).
700. Guyatt GH, Oxman AD, Kunz R, et al. GRADE guidelines 6. Rating the quality of evidence--imprecision.
701 *J Clin Epidemiol*. 2011;64(12):1283–93.
702. Gardner MJ, Altman DG. Confidence intervals rather than P values: estimation rather than hypothesis
703 testing. *Br Med J Clin Res Ed*. 1986;292(6522):746–50.
704. Li G, Taljaard M, Van den Heuvel ER, et al. An introduction to multiplicity issues in clinical trials: the
705 what, why, when and how. *Int J Epidemiol*. 2017;46(2):746–55.

7061. Schulz KF, Altman DG, Moher D, et al. CONSORT 2010 statement: updated guidelines for reporting parallel group randomised trials. *PLoS Med.* 2010;7(3):e1000251.
7072. Greenland S, Senn SJ, Rothman KJ, et al. Statistical tests, P values, confidence intervals, and power: a guide to misinterpretations. *Eur J Epidemiol.* 2016;31(4):337–50.
7123. Guyatt GH, Briel M, Glasziou P, et al. Problems of stopping trials early. *BMJ.* 2012;344:e3863.
7124. Altman DG, Bland JM. Absence of evidence is not evidence of absence. *BMJ.* 1995;311(7003):485.
7125. Guyatt GH, Juniper EF, Walter SD, et al. Interpreting treatment effects in randomised trials. *BMJ.* 1998;316(7132):690–3.
7126. Guyatt GH, Oxman AD, Sultan S, et al. GRADE guidelines: 9. Rating up the quality of evidence. *J Clin Epidemiol.* 2011 ;64(12):1311–6.
7127. Hultcrantz M, Rind D, Akl EA, et al. The GRADE Working Group clarifies the construct of certainty of evidence. *J Clin Epidemiol.* 2017;87:4–13.
7128. IQWiG. *General methods.* www.iqwig.de/en/about-us/methods/methods-paper/ (accessed 22 Jun 2022).
7229. Hozo I, Djulbegovic B, Parish AJ, et al. Identification of threshold for large (dramatic) effects that would obviate randomized trials is not possible. *J Clin Epidemiol.* 2022;145:101–11.
7230. Bross IDJ. Pertinency of an extraneous variable. *J Chronic Dis.* 1967;20(7):487–95.
7231. NICE. *Glossary* www.nice.org.uk/glossary?letter=c (accessed 22 Jun 2022).
7232. ClinicalTrials.gov. *Glossary of common site terms.* www.clinicaltrials.gov/ct2/about-studies/glossary (accessed 22 Jun 2022).
7263. E 9 Statistical Principles for Clinical Trials. 2006;37.
7274. EU Clinical Trials Register. *Update.* www.clinicaltrialsregister.eu/ (accessed 22 Jun 2022).
7285. Relton C, Torgerson D, O’Cathain A, et al. Rethinking pragmatic randomised controlled trials: introducing the ‘cohort multiple randomised controlled trial’ design. *BMJ.* 2010;340:c1066.
7336. Seo HJ, Kim SY, Lee YJ, et al. A newly developed tool for classifying study designs in systematic reviews of interventions and exposures showed substantial reliability and validity. *J Clin Epidemiol.* 2016;70:200–5.
7337. Grimes DA, Schulz KF. An overview of clinical research: the lay of the land. *Lancet* 2002;359(9300):57–61.
7358. Collins R, Bowman L, Landray M, et al. The magic of randomization versus the myth of real-world evidence. *N Engl J Med.* 2020;382(7):674–8.
7379. Mansournia MA, Higgins JPT, Sterne JAC, et al. Biases in randomized trials: a conversation between trialists and epidemiologists. *Epidemiology.* 2017;28(1):54–9.
7380. Lewis SC, Warlow CP. How to spot bias and other potential problems in randomised controlled trials. *J Neurol Neurosurg Psychiatry.* 2004;75:181–7.
7441. Carey TS, Boden SD. A critical guide to case series reports. *Spine.* 2003 Aug 1;28(15):1631–4.
7442. Munn Z, Barker TH, Moola S, et al. Methodological quality of case series studies: an introduction to the JBI critical appraisal tool. *JBI Evid Synth.* 2020;18(10):2127–33.
7443. Slim K, Nini E, Forestier D, et al. Methodological index for non-randomized studies (minors): development and validation of a new instrument. *ANZ J Surg.* 2003;73(9):712–6.
7464. Institute of Health Economics. *Development of a quality appraisal tool for case series studies using a modified Delphi technique.* www.ihe.ca/advanced-search/development-of-a-quality-appraisal-tool-for-case-series-studies-using-a-modified-delphi-technique (accessed 22 Jun 2022).
7495. Grimes DA, Schulz KF. Cohort studies: marching towards outcomes. *Lancet.* 2002;359(9303):341–5.
7506. Woodcock J, LaVange LM. Master protocols to study multiple therapies, multiple diseases, or both. *N Engl J Med* 2017;377:62-70.
7527. Park JJH, Siden E, Zoratti MJ, et al. Systematic review of basket trials, umbrella trials, and platform trials: a landscape analysis of master protocols. *Trials.* 2019;20(1):572.
7548. Park JJH, Harari O, Dron L, et al. An overview of platform trials with a checklist for clinical readers. *J Clin Epidemiol.* 2020;125:1–8.
7569. Bhatt DL, Mehta C. Adaptive designs for clinical trials. *N Engl J Med* 2016; 375:65-74.
7570. Park JJH, Hsu G, Siden EG, et al. An overview of precision oncology basket and umbrella trials for clinicians. *CA Cancer J Clin.* 2020;70:125–137.
7591. Lengliné E, Peron J, Vanier A, et al. Basket clinical trial design for targeted therapies for cancer: a French National Authority for Health statement for health technology assessment. *Lancet Oncol.* 2021;22:e430–4.
7622. Arlett P, Kjær J, Broich K, et al. Real-World evidence in EU Medicines regulation: enabling use and establishing value. *Clin Pharmacol Ther.* 2022;111(1):21–3.
7633. US FDA. *Framework for FDA’s Real-World Evidence Program.* Silver Spring: FAD; 2018.

7654. Concato J, Stein P, Dal Pan GJ, et al. Randomized, observational, interventional, and real-world—
766 what's in a name? *Pharmacoepidemiol Drug Saf.* 2020;29(11):1514–7.
7675. Zuidgeest MGP, Goetz I, Groenwold RHH, et al. Series: pragmatic trials and real world evidence: Paper
768 1. Introduction. *J Clin Epidemiol.* 2017;88:7–13.
7696. Nicol GE, Piccirillo JF, Mulsant BH, et al. Action at a distance: geriatric research during a pandemic. *J*
770 *Am Geriatr Soc.* 2020;68(5):922–5.
7737. Huml RA, Dawson J, Lipworth K, et al. Use of Big Data to aid patient recruitment for clinical trials
772 involving biosimilars and rare diseases. *Ther Innov Regul Sci.* 2020;54(4):870–7.
7758. Meinecke AK, Welsing P, Kafatos G, et al. Series: pragmatic trials and real world evidence: Paper 8.
774 Data collection and management. *J Clin Epidemiol.* 2017;91:13–22.
7759. Welsing PM, Oude Rengerink K, Collier S, et al. Series: pragmatic trials and real world evidence: Paper
776 6. Outcome measures in the real world. *J Clin Epidemiol.* 2017;90:99–107.
7770. Karanatsios B, Prang KH, Verbunt E, et al. Defining key design elements of registry-based randomised
778 controlled trials: a scoping review. *Trials.* 2020;21(1):552.
7791. EMA. *Guideline on Registry-Based Studies.* Amsterdam: EMA; 2021.
7802. EUnethTA. Vision paper on the sustainable availability of the proposed Registry Evaluation and Quality
781 Standards Tool (REQueST). www.eunetha.eu/request-tool-and-its-vision-paper/ (accessed 22 Jun
782 2022).
7833. Lauer MS, D'Agostino RB. The Randomized Registry Trial — The next disruptive technology in clinical
784 research? *N Engl J Med.* 2013;369(17):1579–81.
7854. Gliklich RE, Dreyer NA, Leavy MB, eds. *Registries for Evaluating Patient Outcomes: A User's Guide*
786 (3rd ed). Rockville: Agency for Healthcare Research and Quality; 2014.
7875. IQWiG [A19-43] *Development of scientific concepts for the generation of routine practice data and their*
788 *analysis for the benefit assessment of drugs according to §35a Social Code Book V – rapid report.*
789 www.iqwig.de/en/projects/a19-43.html (accessed 22 Jun 2022).
- 790

1st DRAFT