



eunethta
EUROPEAN NETWORK FOR HEALTH TECHNOLOGY ASSESSMENT

EUnetHTA21 - Individual Practical Guideline Document

D4.3.1: DIRECT AND INDIRECT COMPARISONS

Version 0.2, 26.07.2022

Template version 04 February 2022

1 Document history and contributors

Version	Date	Description
V0.1	27/04/2022	First draft
V0.2	26/07/2022	Second draft

3 Disclaimer

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

11 Participants

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

12

The work in European Network for Health Technology Assessment (EUnetHTA) 21 is a collaborative effort. While the agencies in the Hands-on Group will be actively writing the deliverable, the entire EUnetHTA 21 consortium is involved in its production throughout various stages. This means that the Committee for Scientific Consistency and Quality (CSCQ) will review and discuss several drafts of the deliverable prior to validation. Afterwards the Consortium Executive Board (CEB) will endorse the final deliverable prior to publication.

18 Conflict of interest

The authors declare that there is no conflict of interest.

20 Copyright

All rights reserved.

22 TABLE OF CONTENTS

23	1 INTRODUCTION.....	5
24	1.1 <i>Definitions</i>	5
25	1.2 <i>Relevant articles in Regulation (EU) 2021/2282</i>	5
26	2 SCOPE AND OBJECTIVE OF THE GUIDELINE.....	6
27	3 GENERAL CONSIDERATIONS.....	7
28	3.1 <i>Initial feasibility questions</i>	7
29	3.2 <i>Assessment of exchangeability</i>	7
30	3.2.1 <i>Assessment of similarity</i>	7
31	3.2.2 <i>Assessment of homogeneity</i>	9
32	3.2.3 <i>Assessment of consistency</i>	10
33	3.3 <i>Possible approaches when the assumptions are violated</i>	12
34	4 METHODS APPLICABLE TO DIRECT OR INDIRECT COMPARISONS	14
35	4.1 <i>Methods for direct comparisons</i>	14
36	4.1.1 <i>Standard approaches</i>	14
37	4.1.2 <i>Application of the Knapp–Hartung method</i>	14
38	4.1.3 <i>Direct comparisons with very few studies</i>	14
39	4.2 <i>Indirect comparisons</i>	15
40	4.3 <i>Evidence synthesis of time-to-event data</i>	17
41	4.3.1 <i>Assessment of the proportional hazards assumption</i>	17
42	4.3.2 <i>(Network) meta-analysis of restricted mean survival time</i>	17
43	4.3.3 <i>(Network) meta-analysis with flexible survival time models</i>	18
44	5 Assessment of population-adjusted methods	20
45	5.1 <i>General considerations: is population adjustment for indirect comparisons appropriate?</i>	20
46	5.2 <i>Assessing covariate selection (all population-adjusted methods)</i>	20
47	5.3 <i>Additional considerations for outcome regression approaches</i>	21
48	5.4 <i>Additional considerations for matching-adjusted indirect comparisons</i>	23
49	5.5 <i>Dealing with unanchored MAICs and STCs: additional challenges</i>	24
50	5.6 <i>Interpretation and use of population-adjusted results</i>	24
51	6 ASSESSMENT OF COMPARISONS BASED UPON NON-RANDOMISED EVIDENCE	26
52	6.1 <i>General considerations</i>	26
53	6.2 <i>Propensity scores</i>	26
54	6.2.1 <i>Checking that the assumptions of propensity score matching and/or weighting are</i>	
55	<i>valid</i>	26
56	6.2.2 <i>Interpreting results of propensity score</i>	27
57	7 FURTHER RELEVANT DOCUMENTS (UNDER DEVELOPMENT).....	29
58	8 FUTURE RECOMMENDATION.....	29
59	9 References	30
60		

61 **LIST OF ABBREVIATIONS**

Abbreviation	Meaning
AgD	Aggregate data
AIC	Akaike information criterion
ATE	Average treatment effect
ATT	Average treatment effect among treated
BIC	Bayesian information criterion
CA	Collaborative assessment
DIC	Deviance information criterion
DSL	DerSimonian-Laird
EMA	European Medicines Agency
EU	European Union
EUnetHTA	European Network for Health Technology Assessment
FP	Fractional polynomials
HR	Hazard ratio
HTD	Health technology developer
IPD	Individual patient-level data, also known as individual patient data or individual participant data
IQWiG	Institut für Qualität und Wirtschaftlichkeit im Gesundheitswesen
ITC	Indirect treatment comparison
JCA	Joint clinical assessment
KH	Knapp-Hartung
MAIC	Matching-adjusted indirect comparison
MD	Mean difference
ML-NMR	Multilevel network meta-regression
MS	Member State
NMA	Network meta-analysis
OR	Odds ratio
PH	Proportional hazards
PICO	Population, intervention, comparator, outcome
RCT	Randomised controlled trial
RMST	Restricted mean survival time
RR	Risk ratio
RoB	Risk of bias
ROBINS-I	Risk of bias in non-randomised studies - of interventions
SAP	Statistical analysis plan
STC	Simulated treatment comparison
SUCRA	Surface under the cumulative ranking curve

62

63

64 1 INTRODUCTION

65 The European Network for Health Technology Assessment (EUnetHTA) 21 Methodological Guideline
66 D4.3.2 *Direct and Indirect Comparisons* is an update of a previous EUnetHTA guideline (2013, updated
67 in 2015); it describes the currently available method for direct and indirect comparisons regarding their
68 underlying assumptions, strengths, and weaknesses, and specifies the appropriateness of methods to
69 the data situation. This Practical Guideline is intended for assessor/co-assessors and gives additional,
70 more-detailed advice for use in practice. As indicated in the criteria for selection of assessors and co-
71 assessors, it is expected that statistical expertise will be available in the assessment team (see the
72 EUnetHTA 21 Guideline D5.3.1 *Procedural Guideline for Appointing Assessors and Co-Assessors for*
73 *JCA/CA*).

74 1.1 Definitions

75 The terms used in this document might be used with a slightly different meaning in other contexts. Below,
76 we define the terms as they are used in this guideline.

77 **Direct comparison:** comparison of treatments either by means of a single comparative study or a
78 pairwise meta-analysis or other method for synthesis of comparative studies without indirect
79 comparisons;

80 **Effectiveness:** describes how well a treatment works in practice; includes efficacy and safety;

81 **Exchangeability:** if patients from one treatment group are substituted to another, the same treatment
82 effect is expected; contains the components similarity, homogeneity, and, in the case of indirect
83 comparisons, consistency;

84 **Health technology:** treatment, intervention, and health technology are terms for any health technology
85 that can be assessed;

86 **Indirect comparison:** evidence synthesis in which inference about the relative effectiveness of two
87 treatments is made without the use of trials comparing both treatments head-to-head; indirect
88 comparisons are also made when more general methods of network meta-analysis are applied, even
89 when direct evidence for the comparison of interest is available;

90 **Network meta-analysis (NMA):** generalisation of meta-analysis to include more-complex evidence
91 networks, which can include both direct evidence and indirect evidence; NMA incorporates other terms
92 used in the literature to describe the synthesis of both direct and indirect evidence, such as mixed
93 treatment comparisons and indirect treatment comparisons;

94 **Population-adjusted method for indirect comparisons:** method for indirect comparisons with
95 adjustment for imbalances in effect modifiers between studies and additionally, in the case of
96 disconnected networks, for imbalances in prognostic variables between studies.

97 1.2 Relevant articles in Regulation (EU) 2021/2282

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

- 99
- 100 • Article 9: Joint Clinical Assessment (JCA) reports and the dossier of the Health Technology Developer (HTD);
 - 101 • Article 18: Preparation of the joint scientific consultations outcome document.

102 2 SCOPE AND OBJECTIVE OF THE GUIDELINE

103 This Practical Guideline describes how to deal in practice with evidence syntheses in JCA reports and
104 provides guidance for assessors and co-assessors dealing with submitted results of direct and indirect
105 treatment comparisons from HTDs. Each Section of this Guideline contains a list of requirements that
106 should be reported in the JCA reports in cases in which a direct or indirect treatment comparison was
107 submitted. It is not the objective of this Guideline to make explicit recommendations about whether a
108 submitted direct and indirect treatment comparison should be accepted by the Member States (MSs).
109 Each MS should be enabled to decide on the validity of direct or indirect treatment comparisons itself
110 based on the JCA report, which should include all methodological details needed to do so.

111 In the EUnetHTA 21 Methodological Guideline *Direct and Indirect Comparisons*, the methods for
112 evidence syntheses are summarised, and general guidance is provided on which method(s) are
113 appropriate in a particular situation. This Practical Guideline gives more practical advice for
114 assessors/co-assessors within the framework described in the Methodological Guideline. Often, when
115 using evidence synthesis methodology, some assumptions will be made, which might affect the certainty
116 of results. The aim of this Guideline is to enable HTA assessors and developers to identify potential
117 issues and reduce bias and uncertainty as much as possible. However, we recognise that there is an
118 element of subjectivity in the assessment of many assumptions and that decisions might vary between
119 MSs. There might be exceptional circumstances in which methods of evidence synthesis will need to be
120 applied despite uncertainty or doubt as to their validity. We believe that these should be kept to a
121 minimum and only used in circumstances in which there is a lack of other options to produce an estimate
122 of relative treatment effect. In these scenarios, the HTD has to provide convincing evidence to support
123 the claim that the corresponding results still produce meaningful estimates of relative treatment
124 effectiveness, and acknowledge any additional sources of bias and uncertainty. The argument that the
125 required data to apply preferable methods are not available, is insufficient to demonstrate the validity of
126 the results coming from less appropriate methods with high uncertainty.

DRAFT

127 3 GENERAL CONSIDERATIONS

128 JCAs can use results from multiple trials, which are combined through evidence synthesis. A rigorous
129 systematic review of the relevant literature with explicit inclusion and exclusion criteria is a prerequisite
130 before conducting any evidence synthesis. Evidence synthesis can allow researchers to obtain a more-
131 robust estimate of the treatment effect and, in the case of indirect treatment comparisons (ITCs), provide
132 relative treatment effects for interventions that have not been studied in the same trial. However, it is
133 important that the selection of trials and modelling choices is made with caution and is rigorously
134 examined by assessors in collaboration with healthcare professionals and statisticians. Importantly,
135 according to the European Union (EU) regulation, assessors must ensure that estimates are obtained
136 by pooling relative treatment effects from each trial (i.e., compared with an appropriate comparator) and
137 no inference is based on pooling the absolute effect of a particular treatment in a trial (i.e., regarding the
138 mean outcome in one group only). The rest of this Section details how to assess whether trials are
139 sufficiently similar to be combined, the main modelling choices to consider and scrutinise, and the
140 inferences that can or cannot be made based on the methods and data used.

141 3.1 Initial feasibility questions

142 For direct and indirect comparisons by means of evidence syntheses, the aspects of the population,
143 intervention, control, outcome (PICO) framework and the study design of the included studies have to
144 be examined. Depending on the research goal, the patient population of interest, the intervention, and
145 the control are prespecified and studies have accordingly been searched and selected. Only studies
146 relevant for the given research question and fitting to the PICO scheme should be included in the
147 evidence synthesis. Here, we assume that all studies included in a considered evidence synthesis are
148 relevant for the research question and the corresponding PICO.

149 However, patient characteristics, such as distributions of age, sex, disease duration, measurement, and
150 operationalisation of the outcome of interest, and features of the experimental design still need to be
151 assessed in detail. Additional aspects, such as year and region of study conduct or forms of treatment
152 application, also have to be assessed if they potentially represent possible effect modifiers (see following
153 Sections).

154 Requirements for reporting:

- 155 • A determination that the studies included in the evidence synthesis match the established PICO
156 based on all information described above.

157 3.2 Assessment of exchangeability

158 The fundamental assumption of exchangeability, which is required for (network) meta-analysis, is
159 operationalised by assessing the properties of similarity, homogeneity, and, in the case of indirect
160 comparisons, consistency. We emphasise here that these three properties are not, strictly speaking,
161 distinct assumptions, because a failure of homogeneity or consistency is often the result of an imbalance
162 in effect modifiers between studies (i.e., a violation of similarity). However, in many cases, not all effect
163 modifiers will be known or reported across all studies and, therefore, assessment of homogeneity and
164 consistency (if relevant) could detect an imbalance in unknown effect modifiers that would not be
165 identified through assessment of similarity alone. In situations in which few studies are available for one
166 or more pairwise comparisons, statistical tests might be underpowered to detect violations of
167 homogeneity or consistency and, therefore, the assessment of exchangeability will depend entirely on
168 the similarity of the included studies in terms of observed characteristics. Thus, assessors should be
169 aware that such assessments cannot explore the potential impact of unknown effect modifiers.

170 3.2.1 Assessment of similarity

171 The similarity assumption states that all studies considered are comparable with respect to possible
172 effect modifiers across all interventions. This is tested by means of the PICO scheme (see above). The
173 PICO scheme chosen and the resulting inclusion and exclusion criteria apply to all studies included in

174 the evidence synthesis. For similarity, the following aspects should always be evaluated to identify
175 possible effect modifiers [7].

- 176 **1. Study and patient characteristics** (including duration of follow-up): a list of potential effect
177 modifiers should be drawn up *a priori*. The basis for this can be not only clinical considerations,
178 but also findings from other studies on the therapeutic indication. The following characteristics
179 are generally relevant: age, sex, disease severity, region, and study duration. Only those factors
180 that are identified as potential effect modifiers should be included in this list. Effect modifiers
181 can be identified through a literature search, input from healthcare professionals and, other
182 methods;
- 183 **2. Characteristics of the intervention** (e.g., dosage or application, concomitant treatments): an
184 intervention should be considered as a comparator even if it has not yet been granted European
185 Medicines Authority (EMA) marketing authorisation.

186 The evaluation of similarity should also consider methodological factors that should not differ
187 substantially between studies. Consideration of the observed values of relevant outcomes has also been
188 shown to be helpful in evaluating similarity.

- 189 **3. Characteristics of outcomes** (e.g., definitions of outcomes): an *a priori* definition of what is
190 considered sufficiently similar for each characteristic will usually be difficult. It will often also
191 depend on what is present in the studies included;
- 192 **4. Observed values of relevant outcomes at baseline:** an examination of the observed values
193 of relevant outcomes at baseline can provide information on the similarity of the individual
194 studies, especially the study arms in which the comparator is used. However, to determine
195 similarity, it is not a standard prerequisite that the observed values have to be identical, because
196 the distribution of prognostic variables might well differ between studies. Nevertheless, extreme
197 differences that even lead to floor or ceiling effects regarding the range of possible outcome
198 values should not exist. If the corresponding information is not available at baseline, the values
199 recorded during the course of the study or at the time of analysis can be used instead.

200 It is important to note the following issues regarding effect modification:

- 201 • Not all prognostic variables are effect modifiers;
- 202 • Effect modification is a property of the relative effect between a pair of treatments. As such, it is
203 possible for a variable to modify the relative treatment effect of A versus B, but not the effect of
204 treatment A versus treatment C. This could occur, for example, when A is placebo, B is a therapy
205 targeting a particular genetic mutation [e.g., an epidermal growth factor receptor (EGFR)
206 tyrosine kinase inhibitor (TKI)], and C is another active treatment that does not specifically target
207 this mutation (e.g., chemotherapy): in this case, the presence of this genetic mutation in an
208 individual could be an effect modifier for A versus B but not A versus C. More generally, a
209 variable could be an effect modifier for both A versus B and A versus C, but the magnitude and
210 even the direction of this effect could differ between comparisons (e.g., if patients responded
211 less well to chemotherapy in the presence of this genetic mutation);
- 212 • The status of a variable as an effect modifier, and the magnitude and direction of this effect, is
213 specific to the scale on which the treatment effect is measured. For example, in a hypothetical
214 placebo-controlled study of an influenza vaccine, female participants experience a reduction in
215 risk from 10% to 5% and male participants from 6% to 3%, with vaccination compared with
216 placebo. On the relative risk scale, sex is not an effect modifier [relative risk (RR)=0.5 in both
217 groups], but it is on the risk-difference scale (-5% for females versus -3% for males).

218 It is essential that the process used to identify relevant effect modifiers is comprehensive and
219 transparently reported. This process should include a comprehensive review of the literature and
220 consultation of healthcare professionals with knowledge of the disease area. The set of all potentially
221 relevant effect modifiers should be reported in the submission.

222 The assessment of similarity should include a quantitative analysis of the impact on all observed patient
223 covariates. However, statistical tests for effect modification using subgroup data from clinical trials (e.g.,
224 testing for the significance of interaction terms) will often be underpowered and suffer from issues with
225 multiplicity. Given that the risk of both type 1 and type 2 errors is typically high, statistical tests for effect
226 modification should not be used in isolation to justify the selection of covariates as potential effect
227 modifiers [19,21]. The assessor should also obtain opinions from healthcare professionals to assess
228 whether there are missing effect modifiers.

229 After assessment of all these aspects, a decision has to be made about whether all studies considered
230 in the evidence synthesis are comparable with respect to possible effect modifiers across all
231 interventions (sufficient similarity) or not (insufficient similarity).

232 Requirements for reporting:

- 233 • Description of methodology used to identify potential effect modifiers and whether it sufficiently
234 captures all possible effect modifiers;
- 235 • Assessment of the list of all potential effect modifiers identified and whether this list is likely to be
236 complete; where possible, estimates of the magnitude and direction of the interaction effects;
- 237 • Description of any likely missing effect modifiers and the magnitude of their effect ;
- 238 • The final conclusion about whether the assumption of sufficient similarity is expected to hold or not,
239 with reasoning.

240 **3.2.2 Assessment of homogeneity**

241 The homogeneity assumption states that there is no meaningful heterogeneity between the effect
242 estimates of the individual studies of each possible direct comparison. Even if studies are sufficiently
243 similar, it is still possible that the data show meaningful heterogeneity. Heterogeneity can be caused by
244 unknown effect modifiers and also by factors initially judged to be sufficiently similar or not judged to be
245 potential effect modifiers. To test the homogeneity assumption for a pairwise comparison, at least two
246 direct studies must be available for this comparison in principle, although typically at least five studies
247 are required for a reliable assessment [9]. If only one study is available for each pairwise comparison,
248 the homogeneity assumption cannot be tested. However, this does not prevent the performance of an
249 indirect comparison. The heterogeneity between the studies has to be assessed to determine whether
250 a pooling of the results is meaningful at all and to choose between the fixed-effect and random-effects
251 approach for the evidence synthesis. It is important to use statistical methods as well as design features
252 of the included studies to assess heterogeneity.

253 Two widely used statistical approaches to assess heterogeneity are given by the statistical test based
254 on the Q statistic (Q-test) [8,53] and the heterogeneity measure I^2 [8,25]. As a rough guide for the
255 interpretation of I^2 , the following overlapping categories were proposed [9]:

- 256 • 0–40%: might not be important;
- 257 • 30–60%: might represent moderate heterogeneity;
- 258 • 50–90%: might represent substantial heterogeneity;
- 259 • 75–100%: Considerable heterogeneity.

260 However, the importance of observed I^2 values depends on the magnitude and direction of treatment
261 effects and the strength of evidence for heterogeneity (p -value from the Q-test, uncertainty of the I^2 , or
262 number of studies) [9].

263 One easy and objective criterion to decide whether the studies should not be pooled is given by the
264 statistical significance of the Q-test ($p < 0.05$). However, the current data situation should always be
265 considered when interpreting the results of the Q-test. On the one hand, the Q-test suffers from low
266 power, especially in the situation of few studies [8], which means that a non-significant Q-test does not
267 necessarily indicate that there is no relevant heterogeneity. On the other hand, in the case of a large
268 number of studies, the Q-test might be statistically significant although only low heterogeneity is shown

269 in the forest plot. In such instances, the I^2 measure can help to describe the amount of heterogeneity.
270 For example, if $I^2 < 50\%$, it might be useful to decide that there is no substantial heterogeneity even if
271 the Q-test is statistically significant. In any case, a graphical inspection of the forest plot is advisable in
272 addition to the use of the Q-test and the heterogeneity measure I^2 for the assessment of heterogeneity
273 [8].

274 After assessing heterogeneity, it must be determined whether there is meaningful heterogeneity
275 between the effect estimates of the individual studies of each possible direct comparison (insufficient
276 homogeneity) or not (sufficient homogeneity). If it can be decided that there is sufficient homogeneity
277 and it is meaningful to pool the included studies, it has to be determined whether a fixed-effect or a
278 random-effects model should be used for the evidence synthesis. In the case of indirect comparisons,
279 the assessment results of the consistency assumption also have to be considered (see below). A fixed-
280 effect model assumes a common treatment effect in all studies, which is implausible in many situations
281 and requires rigorous justification from the HTD. Therefore, the standard approach is to use the random-
282 effects approach. However, if there is a marked consistency of the PICO and design properties of all
283 studies, for instance in the case of evidence syntheses when only few studies are available, a fixed-
284 effect model might be appropriate. An example of where the fixed-effect model can regularly be
285 assumed to be valid is the situation of two studies with identical design, which can be found after drug
286 approval. In general, however, the random-effects model is the appropriate choice, although its use
287 might not be feasible in practice (see Section 4.1.3).

288 Requirements for reporting:

- 289 • The complete evaluation of whether the analyses provided to support the homogeneity assumption
290 (including the forest plots, the p -values for the heterogeneity test, and the I^2 values) for all pairwise
291 comparisons are sufficient to demonstrate that it is likely to hold;
- 292 • The final conclusion of whether the assumption of sufficient homogeneity holds or not with reasoning
293 (including sensitivity analyses);
- 294 • The final conclusion of whether it is meaningful to pool the included studies, with reasoning;
- 295 • The decision of whether a fixed-effect or random-effects approach is adequate, with reasoning.

296 **3.2.3 Assessment of consistency**

297 **General remarks**

298 Under the consistency assumption, the same treatment effect is estimated through both the direct and
299 indirect pathways for a particular contrast in the network. Inconsistency is a form of heterogeneity that
300 is linked to the structure of the network but concerns the contrasts between treatments. Thus,
301 inconsistency is between-trial variation comparing different treatment contrasts, and heterogeneity is
302 between-trial variation within treatment contrasts.

303 The choice of methods to assess consistency depends on the network structure. Not all methods are
304 suitable for networks of any complexity, but simpler methods are preferred where they are suitable. The
305 methods used to test for inconsistency should be clearly identified and justification provided for this
306 choice, with reference to the network structure.

307 Although any indirect comparison relies on the consistency assumption, it cannot be tested in networks
308 without a loop structure. Therefore, the first step is to examine the network diagram for loops. It is also
309 important to identify multi-arm trials, because these represent a loop that is consistent by definition.

310 In practice, NMAs often contain too few studies and sparse data to assess inconsistency adequately.
311 Failure to detect inconsistency does not imply that the evidence is consistent. Statistical detection of
312 inconsistency requires more data than are required to establish a treatment effect. Inconsistency can
313 be caused by imbalance in the distributions of effect modifiers in the direct and indirect evidence,
314 commonly factors such as age, severity, and line of treatment, which might be confounded with each
315 other. Therefore, the check for inconsistency should be done alongside the assessment of similarity and
316 heterogeneity in the NMA.

317 To minimise the risk of drawing incorrect conclusions, more empirical indicators are also suggested.
318 Empirical assessment of heterogeneity and the between-trials variation in trial baseline can be used to
319 assess the risk of inconsistency. Comparison of events and responses in the placebo arms might be
320 useful in this context, although we emphasise that, although differences between placebo arms might
321 indicate an imbalance in prognostic variables across studies, this need not result in inconsistency unless
322 these variables are also effect modifiers.

323 ***Bucher method for single loops***

324 The Bucher method is a two-stage method for testing consistency, in which the first step is to synthesise
325 each pair-wise contrast and the second is to test whether the direct and indirect evidence are in conflict.
326 The estimate of inconsistency comes from subtracting the direct and indirect estimates and referring the
327 null hypothesis of no inconsistency to the normal distribution. This test can be applied on three
328 independent sources of data but not on multi-arm trials, because the effect estimates in multi-arm trials
329 are correlated.

330 The Bucher method can also be extended to networks with multiple loops calculating the statistic
331 referring to a chi-square distribution. However, repeated use of the Bucher test in large complex
332 networks with multiple loops can be problematic and, instead, an inconsistency model could be applied
333 for assessing consistency in complex networks. Furthermore, the use of the Bucher method to test for
334 consistency is not advisable when random-effects models are used to synthesise one or more of the
335 pairwise comparisons [12].

336 ***Inconsistency models: Bayesian NMA***

337 The principle of the inconsistency model is to assume no consistency, that all contrasts in the network
338 are unrelated, and that the relative treatment effects are estimated directly from all contrasts (unrelated
339 mean effect). In a consistency model, effects of all included treatments are estimated relative to the
340 reference treatment. To test consistency, the deviance and deviance information criterion (DIC) statistics
341 of the consistency and inconsistency models are compared. Plots of the posterior mean deviance of the
342 individual data points in the inconsistency model against the corresponding posterior mean deviance in
343 the consistency model can help identify loops in which inconsistency is present [16].

344 Further assessment of inconsistency will be a comparison of the posterior estimates of the treatment
345 effect between the consistency and inconsistency models and assessing whether credible intervals
346 overlap.

347 ***Node-splitting methods: Bayesian***

348 The node-splitting method [15] can be applied to any contrast in any network of different complexities in
349 which there is both direct and indirect evidence. In this method, the information contributing to the
350 estimates of a parameter (a so-called 'node') is split into evidence that is direct only and indirect, which
351 is based on the remaining evidence in the network meta-analysis. The indirect estimates in the node-
352 splitting method use not only the indirect evidence of a specific loop, but also the whole evidence base
353 in the network. Comparison of the residual deviance and DIC and the estimate of the heterogeneity
354 parameter for random-effect models, of the full NMA, and of the model with a split node can then be
355 used to assess potential inconsistency between the evidence for a particular node. Reduction of these
356 parameters in the node-split model can be an indicator of inconsistency. The assumption that a split
357 results in equal treatment effects for direct and indirect evidence can be tested in the same way as an
358 inferential hypothesis. Therefore, p -values can be calculated to indicate that the hypothesis of equal
359 treatment effects for direct and indirect evidence can be rejected. Although a smaller p -value would
360 indicate inconsistency, interpretation of these p -values is context dependent and no formal framework
361 for the required significance level exists.

362 Requirements for reporting:

- 363 • Methods used to test for inconsistency and justification for this choice with reference to the network
364 structure; the report should highlight whether the methods used are likely to be appropriate;
- 365 • Criteria used to determine whether a meaningful violation of consistency has been detected;
- 366 • Summary of the results of statistical tests and/or models used to investigate consistency, stating
367 whether these indicate the presence of inconsistency, and describing the extent of the inconsistency
368 and resulting uncertainty in these results;
- 369 • In cases in which inconsistency is detected, description of the possible sources of inconsistency in
370 terms of effect modifiers and, if possible, estimates of the magnitude of effect modification;
- 371 • If methods have been used to explore the qualitative aspects of node splits, resulting p -values
372 should be reported as well as an explanation of the assumptions underlying the analysis;
- 373 • The final conclusion of whether the assumption of sufficient consistency holds, with reasoning.

374 **3.3 Possible approaches when the assumptions are violated**

375 If at least one of the components of the exchangeability assumption is not valid for a considered data
376 situation, the following consequences are possible.

- 377 1. **Splitting into subgroups:** if dissimilarity is shown for a potential effect modifier or heterogeneity
378 is shown that can be explained by the effect modifier, it might be useful to divide the entire study
379 pool into several subpools and draw separate conclusions (e.g., for men and women). The
380 limitations of subgroup analyses based upon aggregated data should be taken into account [14];
- 381 2. **Use of (network) meta-regression:** potential effect modifiers can be included as covariates in
382 a (network) meta-regression model. This requires a sufficient number of data points (= number
383 of studies) so that all parameters can be estimated in the model. The limitations and
384 assumptions of meta-regression based upon aggregated data should be taken into account
385 [4,6,24];
- 386 3. **Exclusion of studies:** in the case that only very few studies are responsible for dissimilarity or
387 heterogeneity, whether it is useful to exclude these studies from the analysis should be
388 discussed. However, in this case, sensitivity analyses should also be performed in which at least
389 some of these studies are included (see below);
- 390 4. **Sensitivity analyses:** if a clearly useful procedure is not possible, sensitivity analyses at least
391 should be performed that allow assessment of the impact of the violated assumptions (e.g.,
392 consideration of study results with unexplained heterogeneity in two separate study pools with
393 homogeneous study results). Given that there are many areas of uncertainty regarding the 'right'
394 methods for meta-analysis, sensitivity analyses are also an important aid in all further decisions
395 in the process to estimate their impact on the results;
- 396 5. **Population-adjusted indirect comparisons:** when there is a suspected violation of the
397 similarity assumption via one of more observed (patient-level) effect modifiers, it might be
398 possible to apply population-adjusted methods, such as matching-adjusted indirect comparison
399 (MAIC), simulated treatment comparison (STC), or multilevel network meta-regression (ML-
400 NMR), to obtain indirect estimates of treatment effects. However, these methods have
401 numerous limitations and might not generate results that are applicable to the research question
402 (see Section 5).

403 The options described above could lead to the formation of new networks and study pools (e.g., for two
404 different subgroups) and, thus, to a separate performance of a direct or indirect comparison. In this case,
405 subsequent testing of the assumptions in the respective new networks and study pools is necessary.

406 Requirements for reporting:

- 407
- 408 • The complete evaluation results regarding potential effect modification;
 - 409 • Approach to, and reasoning for, handling dissimilarity and heterogeneity;
 - 410 • The complete results of all sensitivity analyses;
 - 411 • If the entire study pool was split into several subpools:
 - 412 ○ A complete description of the subpools;
 - 413 ○ The complete evaluation results of the similarity and homogeneity assumptions of all pairwise comparisons within all subpools (see Sections 3.2.1 and 3.2.2).

DRAFT

414 4 METHODS APPLICABLE TO DIRECT OR INDIRECT COMPARISONS

415 4.1 Methods for direct comparisons

416 4.1.1 Standard approaches

417 Standard approaches for meta-analyses according to the fixed-effect model (with the assumption of a
418 common effect in all included studies) are given by the inverse variance method for continuous data and
419 the Mantel–Haenszel method for binary data [3]. Other useful methods are available in special
420 situations, such as rare events (see below).

421 Given that the assumption of a common effect in all included studies can be implausible, the standard
422 model for meta-analyses is usually given by the random-effects model. In accordance to the
423 recommendations of the Cochrane Collaboration, the Knapp–Hartung (KH) method should be used with
424 the Paule–Mandel estimator for the heterogeneity parameter for frequentist meta-analyses with five or
425 more studies [55].

426 With sufficient justification, other methods for meta-analysis can be used in special situations. In the
427 situation of binary data with rare events, the Peto method [9] can be applied. However, this method
428 should not be used when treatment effects are large and the trial arm sizes are unbalanced [5,54]. In
429 the situation of many double-zero studies (i.e., no observed events in both treatment arms), the beta-
430 binomial model can be applied [31,33]. This model allows the inclusion of double-zero studies and
431 contains a random effect for the baseline risk. Nevertheless, the treatment–study interaction is included
432 as a fixed effect, which means that the standard beta-binomial model is a fixed-effect meta-analytic
433 model. As a general alternative to frequentist methods, a Bayesian approach can be used for meta-
434 analysis provided that the required prior distributions are available and can be justified [52].

435 In general, the choice of methods for direct comparisons must be justified. This includes, but is not
436 limited to, justification for the use of a fixed-effect model over a random-effects model, the choice of
437 informative, non-informative, or vague priors (Bayesian), and baseline risk adjustment models.
438 Additionally, any subgroup or meta-regression analysis for different levels of identified effect modifiers
439 must be described and justified. Further considerations must be given to the number and heterogeneity
440 of the included studies, number of events (rare versus common events), scale [odds ratio (OR), RR,
441 hazard ratio (HR), or mean difference (MD)], quality of evidence etc. when assessing the
442 appropriateness of the method and model choices.

443 4.1.2 Application of the Knapp-Hartung method

444 For direct comparisons based on the random-effects model, the general standard approach is given by
445 the KH method. In general, this method holds the type 1 error even in the case of few studies [55].
446 However, in homogeneous data situations, the standard error of the estimated treatment effect
447 according to the KH method might be arbitrarily small. In this case, the calculated confidence interval is
448 misleadingly narrow [3]. To avoid such misleading results, a simple *ad hoc* variance correction was
449 proposed [30]. In practice, a check is required to decide whether the use of the *ad hoc* variance
450 correction is required. A comparison with the confidence interval calculated by means of the
451 DerSimonian-Laird (DSL) method [10] should be used for this purpose. If the confidence interval of the
452 KH method is narrower than that of the DSL method, the use of the *ad hoc* variance correction is required
453 [29,58]. However, the application of the KH method with *ad hoc* variance correction can reduce the
454 power of the KH method. Thus, in the case of very few studies, the KH method can lead to non-
455 informative results (see Section 4.1.3).

456 4.1.3 Direct comparisons with very few studies

457 Meta-analyses with fewer than five studies are problematic in most cases. First, a reliable assessment
458 of heterogeneity is frequently not possible and, therefore, the choice between the fixed-effect and the
459 random-effects model is difficult. Second, the standard random-effects KH approach has frequently very
460 low power. The power might be so low that the KH confidence interval is wider than the union of all
461 confidence intervals of the included studies [48]. In such cases, the KH method is not useful because
462 the results are non-informative and, thus, alternative approaches are required.

463 In general, a random-effects model should be applied even for meta-analyses with very few studies.
464 However, the chance that the assumption of the fixed-effect approach is valid is greater in the case of
465 very few studies compared with situations with a large number of studies. Especially in the situation with
466 only two studies, it might be justified to apply the fixed-effect model by default. This means that the fixed-
467 effect model should always be applied when there are only two studies, unless there are clear reasons
468 against its use.

469 In the situation in which a random-effects model is indicated (i.e., 2 studies and clear reasons against
470 the fixed-effect model and in the case of three or four studies without clear reasons in favour of the fixed-
471 effect model), the first approach should be to use the random-effects model by means of the KH method
472 (with or without *ad hoc* variance correction). However, because of the low power, a comparison with the
473 DSL method is helpful to find a valid conclusion. If both methods (DSL and KH) yield the same result
474 regarding statistical significance, the corresponding conclusion can be drawn for the treatment effect. If
475 the estimated treatment effect resulting from the DSL method is statistically significant but that of the
476 KH method is not, the situation is less clear. In this case, a qualitative summary of the study results can
477 be performed. If there are at least two statistically significant studies in the same direction and most of
478 the available evidence supports this direction, the conclusion of a significant effect can be drawn,
479 although this effect cannot be quantified. To explore whether most of the evidence supports this
480 direction, the study weights of the performed random-effects model according to the KH method can be
481 used. As an example, clear thresholds for the required study weights are proposed in the methods paper
482 of the Institute for Quality and Efficiency in Health Care (Institut für Qualität und Wirtschaftlichkeit im
483 Gesundheitswesen, IQWiG) [27]. More details about the model choice for direct comparisons with very
484 few studies are given elsewhere [48].

485 Alternatively, a random-effects Bayesian meta-analysis with weakly informative prior distribution for the
486 heterogeneity parameter might be useful in the case of very few studies, because external heterogeneity
487 information decreases the problem of estimating heterogeneity with insufficient data [44]. A clear
488 rationale is required for the choice of the prior distributions, because this choice can have substantial
489 effects on the final results in Bayesian meta-analyses with very few studies [50]. Moreover, the impact
490 of the chosen prior distribution should be explored in sensitivity analyses.

491 Requirements for reporting:

- 492 • Assessment of whether the assumptions of the chosen method for meta-analysis are justified;
493 in the case of deviations from standard meta-analytic approaches, a thorough justification for
494 the chosen approach should be given and assessed in the JCA;
- 495 • Assessment of the forest plot with point estimates and confidence intervals of all included
496 studies, the p -value of the heterogeneity test, the I^2 value, and the pooled effect estimate with
497 confidence interval; in the case of a random-effects model, an assessment of the prediction
498 intervals;
- 499 • Determination of whether a fixed- or random-effects model is appropriate;
- 500 • In the case of a Bayesian meta-analysis, an assessment of the chosen prior distributions with
501 justification and sensitivity analyses (see also Section 4.2 for reporting requirements for
502 Bayesian approaches); a clear indication of the extent to which the estimated treatment effects
503 are sensitive to the choice of prior distribution; where informative prior distributions are used,
504 the justification for this should be assessed;
- 505 • In the case of a qualitative summary of the study results, a description of the chosen approach
506 with criteria used for the decision whether there is an overall effect (e.g., thresholds for study
507 weights).

508 **4.2 Indirect comparisons**

509 In this section, important domains in assessment of the credibility of indirect comparison methods are
510 described.

511 The first step in the assessment of the statistical analysis is to consider whether the method used is
512 correct for the network of evidence. The Bucher indirect treatment comparison is appropriate for a
513 network comprising two treatments indirectly compared through a common comparator. The Bucher
514 method can also be applied in star-shaped and ladder networks, but then as multiple pairwise
515 comparisons. Multi-arm trials can only be included as pairwise comparisons, but the generated effect
516 estimates are correlated and the corresponding standard errors are inappropriate. This correlation will
517 be problematic if the aim is to use the estimates in a decision model because the method assumes
518 independence between pairwise comparisons. In cases with several different pairwise comparisons, a
519 network meta-analysis encompassing all this evidence should be considered. Frequentist and Bayesian
520 methods are equally applicable. Naive comparisons should not be used because they do not preserve
521 randomisation. Unanchored and disconnected evidence networks cannot be analysed with these
522 methods (see Section 6).

523 If the method for analysis is deemed appropriate (assumptions met), the appropriateness of the model
524 used must be validated. This includes, but is not limited to, justification for the use of a fixed-effect model
525 over a random-effects model, the choice of informed, uninformed, or vague priors (Bayesian), and
526 baseline risk adjustment models. Additionally, any subgroup or meta-regression analysis for different
527 levels of identified effect modifiers must be described and justified. Further considerations must be given
528 to the number and heterogeneity of studies informing each contrast, number of events (rare versus
529 common events), scale (OR, RR, HR, or MD), quality of evidence, and so on, when assessing the
530 appropriateness of the method and model choices. In networks with large discrepancies in the number
531 of studies informing each contrast, the estimated treatment effect might change from being significant
532 in a clinical trial study to non-significant in the evidence synthesis.

533 In addition to the methods of NMA described here and in the EUnetHTA 21 Methodological Guideline
534 D4.3.2 *Direct and Indirect Comparisons*, many other approaches have been proposed in the literature,
535 such as the original method of Lumley [32] and the 'arm-based' NMA introduced by Hong *et al.* [26].
536 However, many of these methods make different fundamental assumptions to those described in this
537 document and are, in general, unlikely to be suitable for use in JCAs [13,46,57]. If such an analysis is
538 presented in a JCA, it is essential that assessors carefully examine the underlying assumptions and
539 assess their plausibility, as well as the relevance of the results obtained.

540 Requirements for reporting:

- 541 • Determination of whether pooling of the studies is meaningful, and justification for this determination
542 (will be informed by the assessment of exchangeability);
- 543 • Assessment of whether the chosen method for the network meta-analysis is appropriate given the
544 evidence base (including assumptions regarding the variances of the effects);
- 545 • Assessment of the graphical and tabular presentations of the evidence network, including the
546 information on the number of randomised controlled trials (RCTs) per contrast;
- 547 • Assessment of the separate results from direct and indirect comparisons, including measures of
548 uncertainty; where both direct and indirect estimates for a particular comparison are available, any
549 discrepancies between them should be highlighted;
- 550 • If possible, assessment of rankograms [surface under the cumulative ranking curve (SUCRA),
551 cumulative probability curves, and probability of being the best treatment];
- 552 • In the case of a Bayesian NMA, an assessment of the following issues:
 - 553 ○ The chosen prior distributions with justification and sensitivity analyses;
 - 554 ○ Plots of the posterior mean deviance of individual data points for the original model versus the
555 inconsistency model;
 - 556 ○ Convergence of the Markov chains

557 **4.3 Evidence synthesis of time-to-event data**

558 **4.3.1 Assessment of the proportional hazards assumption**

559 A (network) meta-analysis of HRs requires that the proportional hazards (PH) assumption holds for all
560 pairwise comparisons in the network. This means that the assumption must be tested for all included
561 studies, which requires either access to individual patient-level data (IPD) for all studies or the
562 construction of pseudo-IPD from digitised Kaplan–Meier curves (e.g., by using the algorithm proposed
563 by Guyot [22]). The PH assumption is also required for comparisons for which there is no direct evidence
564 available; this cannot be assessed directly.

565 Failure of the PH assumption occurs when the HR between treatment arms is non-constant, which can
566 be interpreted as a time-varying treatment effect. An example of this is the delayed effect on survival
567 observed in studies of immunotherapies in the treatment of advanced cancers [34]. When the PH
568 assumption fails, the average HR will vary according to the length of follow-up, which can differ across
569 studies in the network. Furthermore, the HR obtained from the Cox model might be biased as an
570 estimate of this average because of censoring (unless a suitable adjustment is performed to account for
571 this) [47]. Therefore, if the PH assumption is deemed to be implausible for one or more comparisons in
572 the network, then (network) meta-analysis of HRs should not be carried out. In this scenario, there are
573 two alternative approaches that may be undertaken:

- 574 • (Network) meta-analysis of restricted mean survival times (see Section 4.3.2);
- 575 • (Network) meta-analysis of flexible survival models [fractional polynomials (FPs) or piecewise
576 exponential models] (see Section 4.3.3).

577 Requirements for reporting:

- 578 • Assessment of whether the PH assumption has been thoroughly evaluated in the submission, with
579 particular reference to the following criteria:
 - 580 ○ Log-cumulative hazard plots for all studies (the lines representing the intervention and
581 comparator should be parallel if PH holds);
 - 582 ○ Plots of Schoenfeld residuals (these should show no trend over time if PH holds);
 - 583 ○ Results of any statistical tests used to assess the PH assumption;
 - 584 ○ Any opinions from healthcare professionals received on the plausibility of the PH assumption;
585 for example, if a delayed treatment effect is expected, then PH might not hold.

586 **4.3.2 (Network) meta-analysis of restricted mean survival time**

587 When the PH assumption does not hold, it is possible to carry out a (network) meta-analysis of restricted
588 mean survival time (RMST) [45]. This involves the selection of a relevant time-point for follow-up and
589 then calculation of the area under the Kaplan–Meier curve between randomisation and this time.
590 Relative treatment effects are then computed as either the difference or ratio of RMSTs between
591 treatment arms. These effects can then be synthesised in a fixed or random-effects meta-analysis using
592 methods previously described.

593 When RMST is used, a key consideration is the choice of follow-up time, because different choices can
594 produce different results. Possible values are limited by the available data, and some higher values
595 might be more uncertain because of the limited numbers at risk. Therefore, it might be necessary to
596 consider the duration of follow-up of the included studies to select an appropriate time-point. It is
597 important that prespecified criteria for selecting the base case follow-up time are clearly reported, and
598 that a range of follow-up times be presented in sensitivity analysis.

599 **4.3.3 (Network) meta-analysis with flexible survival time models**

600 The use of flexible models for the hazard function allows (network) meta-analysis to be carried out
601 without the assumption of PHs. These methods require the use of IPD or, more commonly, pseudo-IPD,
602 whereby published survival curves for the endpoints of interest are scanned and digitalised (e.g., by
603 using the algorithm proposed by Guyot [22]).

604 FP (network) meta-analysis involves modelling time-dependent hazard rates for each intervention
605 separately (as linear combinations of positive and negative powers of time), allowing for a wide range
606 of different-shaped hazards. Treatment effects comprise multiple correlated parameters, and can be
607 synthesised using fixed or random-effects models. A similar approach is possible using restricted cubic
608 spline models [20].

609 Evidence synthesis using FP requires selection of the most appropriate model for the hazard rates; that
610 is, the most appropriate combination of powers of the time variable. This can be assessed using
611 measures of statistical fit [e.g., Akaike information criterion (AIC), Bayesian information criterion (BIC),
612 or, in a Bayesian framework, DIC], visual fit to the observed hazards and survival functions, and/or
613 clinical plausibility. Assessors should be aware that the use of different FP models can lead to different
614 conclusions regarding relative treatment effects; therefore, sensitivity analysis is important.

615 Piecewise exponential models also allow for a relaxation of the proportional hazards assumption. With
616 this approach, the follow-up period for all treatments is split into a fixed number of pieces, and the hazard
617 rate for each intervention in the network is assumed to be constant within each piece. Treatment effects
618 estimated using this method comprise piecewise HRs; thus, the PH assumption is required within each
619 piece, but not over the entire follow-up period. Such an assumption might be plausible in situations in
620 which a delayed treatment effect is expected (e.g., immunotherapies in oncology), but where PH is
621 expected to hold thereafter. These piecewise hazard ratios can be incorporated into an (network) meta-
622 analysis in the usual way, using fixed- or random-effects models.

623 To carry out piecewise exponential (network) meta-analysis, it is necessary to choose the number and
624 location of the cut-points of the pieces. This can be done using visual and statistical fit to the observed
625 data, and external opinion might also be helpful. Furthermore, the plausibility of the PH assumption
626 within each piece should be assessed using the methods described previously. Assessors should again
627 be aware that choosing different numbers and locations of cut-points can alter the estimated treatment
628 effect; thus, sensitivity analysis is important.

629 A limitation that applies to both FPs and piecewise exponential models is that the estimated treatment
630 effects they produce are multidimensional and not easily interpretable. There is no obvious way to
631 perform statistical inference in this setting (i.e., testing for statistically significant treatment effects). The
632 usual method for addressing this is to compare either restricted mean survival time or extrapolated mean
633 survival time, obtained from the chosen models. This is usually carried out in a Bayesian framework, in
634 which posterior distributions of the model parameters are used to obtain posterior estimates of
635 (restricted/extrapolated) mean survival times. When extrapolation is used, the plausibility of long-term
636 extrapolations should also be considered as part of the model selection process. When restricted means
637 are used, consideration must be given to the chosen time-point.

638
639
640
641
642
643
644
645
646
647
648
649

Requirements for reporting:

- For flexible parametric models: assessment of model choice with reference to measures of statistical fit and any other information used to inform this choice (e.g., clinical opinion);
- For RMST: assessment of the rationale for the choice of follow-up time; sensitivity of the results to this choice should be assessed;
- Comparison of observed and modelled HRs over time (e.g., table of the HRs at different time points and/or the plot of HR over time to indicate whether the chosen method is appropriate);
- Comparison of the survival time distribution implied by the chosen (best-fitting or most plausible) model along with the alternative models and the study KM data; evaluation of visual fit to the observed data;
- Where multiple model choices are comparable in terms of fit and/or plausibility, the results obtained from these alternative models should be compared and assessed.

DRAFT

650 5 Assessment of population-adjusted methods

651 5.1 General considerations: is population adjustment for indirect comparisons 652 appropriate?

653 Population-adjusted methods are used in the context of an ITC or more general NMA, in which there is
654 concern that the similarity assumption might not hold. These methods aim to adjust for this imbalance
655 to obtain an unbiased estimate of the relative treatment effect in the scenario in which IPD is available
656 for one or more trials in the network, and only aggregate data (AgD) for others. MAIC and STC should
657 not be used when full IPD is available for all studies; IPD network meta-regression is generally the
658 appropriate method to adjust for covariate imbalances in this case.

659 The most common examples of population-adjusted methods are MAIC, STC, ML-NMR, and other
660 mixed IPD and aggregate data regression methods. The MAIC method reweights patients in the IPD
661 study to match the characteristics of the AgD study, whereas STC and ML-NMR fit outcome regression
662 models to the IPD studies, which can be extrapolated to other populations. A consideration when
663 selecting among MAIC, STC, and ML-NMR is that MAIC adjustments are only applicable to the
664 population of the AgD study, whereas STC and ML-NMR can be used to extrapolate treatment effects
665 to any population with known covariate values for the effect modifiers. Both MAIC and STC are limited
666 to simple networks with two studies, whereas ML-NMR can be applied to any connected network.

667 When assessing a population-adjusted indirect comparison, the problem of multiplicity arising from
668 'researcher degrees of freedom' must be considered. Indeed, the number of methods and potential
669 covariate combinations available to the modeller raises the possibility of selecting the method that
670 produces the most favourable results for the intervention under assessment. For this reason, these
671 methods are often more suitable as an exploratory analysis rather than as the primary analysis. In the
672 case of anchored comparisons, it should be demonstrated that bias will be reduced by the use of a
673 population-adjusted methods.

674 MAIC and STC methods are also sometimes used in the case of disconnected networks; in this context,
675 absolute outcomes (rather than relative effects) are adjusted and, therefore, adjustment must account
676 for all potential confounders in addition to effect modifiers. The EUnetHTA 21 Methodological Guideline
677 *Direct and Indirect Comparisons* details the many issues regarding population adjustment
678 methodologies for unanchored ITCs. By describing these methods here, we are not endorsing them,
679 and once again reiterate that estimates arising from unanchored ITCs are unreliable.

680 Requirements for reporting:

- 681 • Assessment of the justification for population adjustment as a means of estimating treatment
682 effectiveness;
- 683 • A complete description of the method and/or model(s) used for population adjustment and
684 estimation of the treatment effects, and an assessment of the appropriateness of this choice.

685 5.2 Assessing covariate selection (all population-adjusted methods)

686 The validity of all population-adjusted methods depends on the inclusion of all effect modifiers as
687 covariates in the relevant model. These should be identified using the methods described in Section
688 3.2.1, ideally before conducting the analysis. In the case of unanchored comparisons, all prognostic
689 variables must also be included (see Sections 5.5 and 6).

690 In the case of both MAIC and STC, only effect modifiers of the relative effect being estimated in the IPD
691 trial are needed to carry out the adjustment. However, interpretation of the results also requires
692 knowledge of effect modifiers for the AgD trial. In the case of more-complex networks of evidence (e.g.,
693 using ML-NMR), knowledge of effect modifiers for all pairwise comparisons is typically needed.

694 Covariates that are initially balanced (or approximately balanced) between study populations at baseline
695 should not be omitted because the adjustment procedure could create an imbalance where none existed
696 before. Methods of covariate selection based upon statistical significance or model fit are of limited use
697 for STC and MAIC, given that, when limited IPD is available, these methods will typically be
698 underpowered to detect relevant effects.

699 When effect modifiers are omitted for any reason, population-adjusted treatment effects obtained using
700 the methods described here will necessarily be biased. The magnitude of this bias depends on both the
701 magnitude of effect modification associated with the missing covariate(s) and the extent of the imbalance
702 between treatment groups in terms of this characteristic (after adjustment). It is often unknown whether
703 covariates are missing or which covariates these might be. When multiple relevant-effect modifiers are
704 missing, the combined impact becomes difficult to predict. Assessors should highlight the potential for
705 residual bias in the resulting estimate and give an indication of the size and direction of that bias where
706 possible.

707 Assessors should be aware that, when population-adjusted indirect comparisons are carried out despite
708 relevant covariates being unavailable, bias in the estimated treatment effects could be increased as a
709 result of adjustment compared with the results of a standard NMA. Consider an example in which the
710 relevant effect modifiers are background statin use and history of cardiovascular disease and where, in
711 one study, there is a strong positive association between these two variables (e.g., because of statin
712 therapy being initiated following a cardiovascular event), but there is no such association in the other
713 study (e.g., because of statin being used for primary prevention of cardiovascular disease among these
714 patients). In such a scenario, adjustment for background statin use but not cardiovascular disease in a
715 MAIC could increase the imbalance in the proportion of patients with a history of cardiovascular disease
716 across studies.

717 To account for the risk of bias (RoB) because of missing or unknown effect modifiers, it is possible to
718 perform statistical inference on the estimated treatment effect by testing against a shifted null
719 hypothesis; that is, a null hypothesis of some non-zero relative treatment effect of a magnitude large
720 enough to account for any plausible bias arising from missing covariates. If this is done, the shifted
721 hypothesis to be tested should be prespecified and its magnitude clearly justified.

722 **Requirements for reporting:**

- 723 • Assessment of the methodology used to identify relevant-effect modifiers;
- 724 • Assessment of the adequacy of the set of included effect modifiers to generate an unbiased estimate
725 of the treatment effect;
- 726 • When relevant-effect modifiers have not been included in the assessment model, a quantification of
727 the potential magnitude and likely direction of the resulting bias;
- 728 • If shifted hypothesis testing has been used, an assessment whether this is sufficient to account for
729 the likely magnitude of residual bias arising from missing covariates.

730 **5.3 Additional considerations for outcome regression approaches**

731 The STC and ML-NMR methods involve fitting an outcome regression model (e.g., a generalised linear
732 model or Cox PH model) to the available IPD to obtain an estimate of the outcome at each level of the
733 included covariates. The chosen model must estimate treatment effects on the same scale as the
734 indirect treatment comparison. For example, if the indirect comparison is to be carried out on the log OR
735 scale, an appropriate choice for the outcome regression model would be a generalised linear model with
736 a binomial likelihood and logit link. It is not appropriate to use logistic regression to adjust absolute risks
737 in each arm and then carry out the indirect treatment comparison on the risk-difference or log-risk ratio
738 scales.

739 A fundamental assumption of STC and ML-NMR is that the effect of the covariates is additive on the
740 outcome measure scale (i.e., that the functional form of the outcome regression model is appropriate).

741 For example, if a Cox PH model is used for the treatment effect (log hazard ratio), then the effect of the
742 covariates is assumed to be linear on the log hazard ratio scale (i.e., PHs). In the case of the IPD study,
743 this should be assessed and reported using standard model diagnostics (e.g., analysis of residuals).
744 External data could also be helpful here; for example, the effect of LDL cholesterol levels on
745 cardiovascular event rates has been characterised as approximately linear on a log-rate scale [18].

746 In the case of anchored STCs, the inclusion of additional prognostic variables (that are not also effect
747 modifiers) in the outcome model will not reduce bias, but could improve precision of the estimated
748 treatment effect and, therefore, can be considered. In this case, standard measures of model fit, such
749 as AIC/BIC, residual deviance, and so on, can be used to select these additional covariates.

750 The STC and ML-NMR methods can generate estimates of the treatment effect in any target population
751 by substituting the relevant mean covariate values into the outcome regression model. This can be
752 useful if the population of interest differs from the trial populations. However, the validity of these
753 estimates is unknown outside the range of covariate values included in the IPD study; extrapolation
754 beyond this region might not generate meaningful estimates of the treatment effect. For example, if the
755 age range of the IPD study population is 40–55 years, it would not be appropriate to use STC/ML-NMR
756 to extrapolate treatment effects to a population with a mean age of 60 years, because the relationship
757 between age and treatment effect cannot be assessed outside the range of the IPD study. More
758 generally, treatment effects for covariate combinations that are not well represented in the IPD study
759 will be uncertain. Therefore, the degree of overlap in baseline covariates should be reported and
760 assessed by, for example, plotting the distributions of baseline characteristics in the IPD trial(s) together
761 with the mean and confidence intervals from the AgD trials.

762 The usual approach to STC involves substituting mean covariate values from the AgD population into
763 the outcome regression model, which estimates the conditional treatment effect at this level of the
764 covariates (i.e., the predicted individual-level response) [39,42]. However, the summary effect estimate
765 from the AgD study is typically a marginal treatment effect (i.e., population average) or, in some cases,
766 a conditional effect but typically adjusted for a different set of covariates. As a result, STC, conducted
767 using the substitution of mean covariate values, combines incompatible effect estimates, potentially
768 leading to bias in the estimation for both estimands (conditional and marginal) when the outcome
769 regression model is nonlinear and there are invalid standard errors in all cases [39,42]. To overcome
770 this, approaches to STC targeting marginal treatment effects have been proposed [28,43]. These
771 approaches generally require additional assumptions to estimate the joint covariate distribution from the
772 AgD study; therefore, the plausibility of these assumptions should be assessed. The JCA report should
773 clarify the STC approach used and the target estimand.

774 Additional assumptions and/or data requirements for ML-NMR depend on the targeted treatment effect.
775 In a simple network of one IPD study and one AgD study, the (population average) marginal treatment
776 effect in the AgD population can be estimated with no further assumptions beyond those required for
777 STC [35]. However, estimation of conditional treatment effects, marginal effects in any other population,
778 or any application of ML-NMR in a network with two or more AgD studies requires the estimation of
779 additional treatment–covariate interactions. To achieve this, the available data must include, for each
780 treatment in the network, either full IPD from at least one study investigating that treatment or enough
781 AgD studies investigating that treatment to estimate the relevant interactions. If such data are not
782 available, then the 'shared effect modifier' assumption is required (see Section 5.6) for certain treatment
783 classes within the network and, therefore, the plausibility of this assumption must be assessed [35,38].
784 Furthermore, specification of the joint covariate distributions for the AgD studies is required, which
785 typically necessitates additional assumptions.

786 Requirements for reporting:

- 787 • An assessment of the model fit and appropriateness of the outcome regression model to capture
788 the effect of covariates (including treatments) on outcomes;
- 789 • An assessment of the covariate overlap between the IPD study (or studies) and the populations to
790 which relative treatment effects are adjusted (e.g., the AgD study or studies);
- 791 • For STC, a description of the method used to estimate outcomes (e.g., substitution of mean
792 covariate values, simulation, or numerical integration) and the treatment effect that is targeted by
793 the chosen approach; assessment of whether the estimands that have been combined are
794 compatible, highlighting any potential for bias;
- 795 • For ML-NMR, clear statement as to whether the available data are sufficient to estimate treatment-
796 covariate interactions; statement of any additional assumptions (e.g., shared effect modifier) that
797 have been made to estimate the model;
- 798 • An assessment of the method used to estimate the joint covariate distributions in the AgD studies,
799 if required (applies to ML-NMR and certain approaches to STC).

800 **5.4 Additional considerations for matching-adjusted indirect comparisons**

801 When MAIC is used to carry out population adjustment, the principal concern is whether the weighted
802 pseudo-population has the same distribution of effect modifiers (anchored and unanchored
803 comparisons) and prognostic variables (unanchored only) as the target population. These distributions
804 should be reported and their similarity assessed; if nontrivial differences exist for one or more variables
805 after matching, then the results of the MAIC will likely be biased. The use of hypothesis tests for the
806 equality of means after matching is not recommended as a method to decide whether sufficient balance
807 has been achieved, because multiple tests are typically required (increasing the risk of type 1 error) and
808 statistical power might be low (type 2 error).

809 The distribution of weights should be examined to assess the extent of overlap between the two
810 populations. The approximate effective sample size should also be reported. If this is considerably
811 smaller than the original sample size, then statistical power will be reduced accordingly. The presence
812 of extreme weights and/or low ESS also indicates that the target population of the MAIC is considerably
813 different from the source population. This could be problematic in the context of a JCA because it is
814 likely that the source population is of greater interest to the assessment than the target population (see
815 also Section 5.6).

816 The assessor should ensure that an appropriate method, such as robust standard errors or
817 bootstrapping, has been used to estimate the confidence interval associated with the treatment effect.
818 Failure to do so will result in confidence intervals that are artificially narrow and do not capture the full
819 extent of (statistical) uncertainty in the estimated treatment effect.

820 Requirements for reporting:

- 821 • Assessment of covariate balance achieved after matching, and of potential impact of any residual
822 imbalance on the results (if this can be estimated);
- 823 • Assessment of the distribution of weights and effective sample size after matching to assess the
824 extent of overlap between the two populations;
- 825 • Statement as to whether the reported confidence interval for the treatment effect appropriately
826 captures the additional uncertainty arising from reweighting (e.g., whether the confidence interval
827 has been estimated using an appropriate method, such as robust standard errors or bootstrapping).

828

829 **5.5 Dealing with unanchored MAICs and STCs: additional challenges**

830 Population-adjusted methods for indirect comparisons are also used when considering disconnected
831 networks. Comparing treatments in an unanchored network is essentially a comparison of absolute
832 effects rather than of relative effects, which is not the goal of the JCA. The validity of the results depends
833 on all relevant prognostic variables (as well as effect modifiers) being included as covariates in the
834 relevant model, which is unlikely to be satisfied in practice. In general, this will substantially increase the
835 amount of adjustment required. The process used to identify prognostic variables is analogous to that
836 described previously for effect modifiers and should be reported transparently in the submission.

837 Differences in patient characteristics are typically more likely to affect absolute outcomes than they are
838 relative outcomes, which means that more covariates must be included in the adjustment model to
839 obtain an unbiased estimate the treatment effect. For example, if two hypothetical treatments, A and B,
840 aimed at lowering blood pressure were to be compared in an unanchored comparison, then adjustment
841 would need to be carried out for all covariates potentially affecting blood pressure, such as age, sex,
842 smoking status, race, geographical location, body mass index, diabetes status, and many others that
843 might not have been recorded. By contrast, an anchored comparison of an A and B via a common (e.g.,
844 placebo) comparator would only require adjustment for covariates affecting response to treatment.

845 There will inevitably be differences in the trials other than patient characteristics. Interventions will be
846 administered under different conditions and endpoints might be recorded in different ways (e.g.,
847 investigator versus independent assessment of tumour progression). Again, these differences typically
848 have a greater impact on unanchored comparisons compared with anchored comparisons, because
849 absolute effects are being compared. An assessor should assess these carefully, using opinions from
850 healthcare professionals again if required, to decide whether it is appropriate to undertake an
851 unanchored MAIC or STC.

852 In summary, although these methods are often presented as the only way of quantifying a relative
853 treatment effect, this does not mean that the method will be of sufficient standard to confidently estimate
854 a relative treatment effect. A better, although still problematic, option is the use of methods for the
855 analysis of non-randomised data, which require access to the full IPD information (see Section 6).

856 Requirements for reporting:

- 857 • An assessment of the methodology used to identify all relevant prognostic variables ;
- 858 • An assessment of the appropriateness of carrying out an unanchored indirect comparison, with
859 reference to data availability, definitions of outcomes, comparability of study characteristics, and
860 other considerations; if full IPD were available for all studies, then this should be clearly highlighted
861 because, under this scenario, other IPD-based methods (e.g., propensity score matching) would
862 likely be more appropriate;
- 863 • An assessment of whether the set of included covariates is likely to be sufficient to generate an
864 unbiased comparison of outcomes; quantification of the magnitude and direction of potential bias
865 arising from missing prognostic variables in the analysis.

866 **5.6 Interpretation and use of population-adjusted results**

867 Population adjustment estimates treatment effects in the population of the AgD study, which might not
868 be generalisable outside of that population. However, in the context of JCAs, it is likely that estimation
869 of the treatment effect in the population of the 'source' (IPD) study is of interest. Phillippo *et al.* [37]
870 highlight this issue in relation to two MAIC analyses of the same two trials comparing secukinumab and
871 adalimumab to placebo as treatments for ankylosing spondylitis. The relative treatment effect differed
872 depending on which trial the IPD was taken from, which is explained by patient differences in the target
873 studies. Although STC and ML-NMR have the advantage that treatment effects from the IPD study can
874 be extrapolated to any population, none of the available population adjustment methods can estimate
875 the relevant (indirect) treatment effect outside of the population of the target AgD study unless additional
876 assumptions are made [37].

877 In some situations, it might be reasonable to ‘transpose’ effect estimates obtained from anchored
878 population-adjusted indirect comparisons to other populations, such as that of the source trial. Doing so
879 requires the additional ‘shared effect modifier’ assumption proposed by Phillippo *et al.* [36]. This
880 assumption applies to a set of active treatments and states that, relative to a common comparator: (i)
881 the covariates that are effect modifiers and (ii) the change in treatment effect for each effect modifier
882 (i.e., the magnitude and direction of the interaction terms), are the same for all active treatments in this
883 set. When this holds, the relative effect between any pair of treatments in this set will be the same in
884 any population, which means, in particular, that treatment effects obtained from population-adjusted
885 indirect comparisons can be transposed to the population of the source (IPD) trial or indeed any other
886 relevant population. The shared-effect modifier assumption is more likely to hold for treatments with a
887 similar mode of action (e.g., an ITC of two angiotensin-converting enzyme inhibitors) than for those in
888 different classes (e.g., in an ITC of an angiotensin-converting enzyme inhibitor versus an angiotensin
889 receptor blocker). Strong biological and/or clinical justification must be provided to justify its use in a
890 JCA.

891 Different population adjustment methods target different estimands. The MAIC method targets the
892 marginal treatment effect (population average effect over the AgD population), whereas STC, performed
893 using substitution of mean covariate values, targets the conditional treatment effect at the specified level
894 of the covariates (individual-level treatment effect for the ‘average’ patient in the AgD population). In its
895 most general form, ML-NMR can target either estimand in any target population [40-42].

896 To incorporate results of an MAIC or STC into a wider NMA, it is necessary to assume similarity of effect
897 modifiers across the network after adjustment; in other words, the distribution of effect modifiers across
898 all studies in the network is similar to that of the target study rather than of the source study.

899 Population adjustment aims to reduce bias arising from an imbalance of effect modifiers (or prognostic
900 variables for unanchored comparisons) but does so at the cost of increased variance. The result is a
901 loss of precision when estimating treatment effects (i.e., wider confidence intervals) or, equivalently, a
902 loss of statistical power. When inference is made on the basis of population-adjusted comparisons,
903 assessors should take into account that these comparisons are typically underpowered.

904 Requirements for reporting:

- 905 • A clear description of the population in which the treatment effect has been estimated, and its
906 relevance to the assessment question; any limitations should be clearly outlined, and potential
907 biases arising from population differences should be reported (including an assessment of the
908 likely magnitude and direction of any bias, if possible);
- 909 • Clear statement as to whether the ‘shared-effect modifier’ assumption is required to estimate
910 the treatment effect in the target population; if this assumption is invoked, the biological and/or
911 clinical basis for this assumption should be scrutinised and the strengths and limitations clearly
912 described;
- 913 • A comparison between the population-adjusted estimates of treatment effects with those
914 obtained from standard methods of (network) meta-analysis; if the magnitude, direction, and/or
915 precision of these effects differ considerably, then assessors should discuss likely explanations
916 for this (e.g., covariate adjustment, loss of effective sample size, or underlying assumptions);
- 917 • The target estimand of the chosen population-adjustment method, that is, marginal (population-
918 average relative treatment effect) or conditional (individual level treatment effect for the
919 ‘average’ patient) relative effects, and its relevance to the assessment question.

920 **6 ASSESSMENT OF COMPARISONS BASED UPON NON-RANDOMISED EVIDENCE**

921 **6.1 General considerations**

922 All commonly encountered sources of evidence outside of RCTs are non-randomised (i.e., single-arm
923 trials, cohort studies, case-control studies, other observational studies, and the use of historical
924 controls). Any such study has much greater potential to include material bias in the estimate of treatment
925 effect compared with an appropriate RCT, and this is likely to carry through when combining evidence
926 from these sources. A key concern is that the underlying assumption of exchangeability is unlikely to
927 hold because there is a very high risk of confounding bias, meaning that the association between
928 intervention and outcome differs from its causal effect.

929 Therefore, treatment comparisons based upon non-randomised evidence require careful consideration
930 of its validity. The inclusion and exclusion criteria for each study should be carefully examined, because
931 these criteria are typically more restrictive for clinical trials than for observational studies, leading to
932 potential violations of the positivity assumption (e.g., individuals with very poor prognosis are often
933 excluded from clinical trials but not from cohort studies). The potential for unmeasured confounding
934 arising from ‘volunteer bias’ should also be considered when interpreting the results: willingness to
935 participate in a clinical trial might be associated with several prognostic variables that might be
936 unmeasured, such as access to medicine, socioeconomic status, location, educational attainment, and
937 overall health status. The RoB in the results should be assessed by using an appropriate tool, such as
938 Risk of Bias in Non-Randomised Studies – of Interventions (ROBINS-I) [49] (see also the EUnetHTA 21
939 Practical Guideline *Validity of Clinical Studies* for practical guidance on the assessment of the validity
940 of individual studies).

941 In some cases, it might be that the lack of randomisation can be compensated for by rigorous adjustment
942 for confounding. However, this requires that all confounders and effect modifiers relevant for adjustment
943 are measured and that the model and covariate selection strategies for adjustment are prespecified and
944 based upon transparent criteria [23]. The requirement of all confounders and effect modifiers being
945 measured is unlikely to be met, given that unknown modifiers and confounders are assumed to be
946 always present. These adjustment methods require access to the full IPD information. Aggregated data
947 alone are not sufficient to reliably estimate treatment effects. A statistical analysis plan (SAP) is required
948 to describe the methods planned to adjust for confounding.

949 Requirements for reporting:

- 950 • Assessment of the inclusion and exclusion criteria for the relevant non-randomised data;
- 951 • Assessment of the RoB and the validity of the results of all included trials;
- 952 • Comparison of baseline characteristics of all included trials;
- 953 • An assessment of the methodology used to identify all relevant prognostic variables and effect
954 modifiers;
- 955 • An assessment of the SAP with the methods used to adjust for confounding;
- 956 • An assessment of whether the set of included covariates is likely sufficient to generate an
957 unbiased comparison of outcomes; quantification of the magnitude and direction of potential
958 bias arising from missing prognostic variables and effect modifiers in the analysis.

959 **6.2 Propensity scores**

960 **6.2.1 Checking that the assumptions of propensity score matching and/or weighting are** 961 **valid**

962 An important method to adjust for confounding in non-randomised studies is by using propensity scores.
963 This method requires careful planning of all possible modelling options in the form of an SAP [56]. As
964 mentioned in EUnetHTA 21 Methodological Guideline on *Direct and Indirect Comparisons*, three

965 assumptions must be met when using non-randomised data and propensity scores or another method
966 to adjust for confounding: positivity, overlap, and balance. In the JCA context, the assessor/co-assessor
967 must check and report the validity of these assumptions.

968 **Checking the positivity assumption**

969 The positivity assumption means that patients in both groups must be theoretically eligible for both
970 treatments of interest. In randomised evidence, positivity is guaranteed by randomisation. In non-
971 randomised evidence, the positivity assumption concerns the probability of receiving treatment, but this
972 probability needs to be modelled (e.g., by propensity score) because of the absence of randomisation.
973 Suspicion of violation for positivity assumption (e.g., inclusion of patients in one treatment group, with
974 contraindication to the other treatment group, or patients with a propensity score equal or close to zero
975 or one) should be systematically reported.

976 **Checking the overlap assumption**

977 Sufficient overlap means that the distribution of patients among the different propensity scores must be
978 similar. To allow this assumption to be checked, propensity score distribution (using histograms or
979 density plot), among samples if applicable (i.e., whole population, matched population, and/or population
980 created by weighting), should be reported in the JCA and discussed. The overlap depends on the
981 matching performed and the techniques used [17]. In the case of trimming, if a large proportion of the
982 sample is lost after trimming regions of non-overlap, then it could indicate insufficient overlap [11]. When
983 trimming is performed, the selected population should be described in detail to investigate whether it
984 sufficiently represents the original research population.

985 **Checking the balance assumption**

986 The populations in the compared groups must be sufficiently balanced after adjustment for confounding.
987 The achieved balance must be assessed before and after matching, weighting, or stratification. Absolute
988 standardised differences between the treatment groups should be used to compare the balance for each
989 covariate [2]. Cut-offs for acceptable absolute standardised difference vary (0.1–0.25) [51]. Therefore,
990 the final conclusion regarding the balance assumption would be left to MS for absolute standardised
991 differences <0.25; if any absolute standardised difference is ≥ 0.25 , violation of the balance assumption
992 should be stated. Doubly robust methods combining propensity scores and outcome regression can be
993 used to reduce bias arising from residual covariate imbalance after matching or weighting.

994 **Checking the inferential goal**

995 The inferential goal (i.e., target of inference) determines, in part, the choice of a specific propensity score
996 method. The most common estimands are the average treatment effect (ATE) and the average
997 treatment effect among treated (ATT). Adequation between inferential goal and chosen estimand (ATT
998 or ATE) should be evaluated by the assessor/co-assessor. Adequation between propensity score
999 method (matching, stratification, adjustment using the propensity score, or weighting) and the chosen
1000 estimand should also be assessed (e.g., matching primarily estimates ATT) [1].

1001 Requirements for reporting:

- 1002 • An assessment of the SAP with the propensity score methods used to adjust for confounding;
- 1003 • An assessment of the required assumptions of sufficient positivity, overlap, and balance;
- 1004 • The final decision whether the assumptions of positivity, overlap, and balance hold, with
1005 reasoning based on the analyses submitted by the HTD.

1006 **6.2.2 Interpreting results of propensity score**

1007 In the JCA report, when propensity score methods are used, qualitative analysis must be first performed
1008 by the assessor/co-assessor, assessing the validity of assumptions (see Section 6.2.1). If underlying

1009 assumptions are considered to be violated, this must be explicitly reported before quantitative results
1010 are interpreted, because they would be biased.

1011 Quantitative results (effect estimates with confidence intervals) should be presented for both crude and
1012 propensity score analyses. Results of sensitivity analyses should always be presented to evaluate the
1013 robustness of results. Quantitative results assess the degree of statistical association, but a statistically
1014 significant association does not necessarily imply a causal relationship. The JCA report should be factual
1015 and the assessor/co-assessor is not supposed to conclude on causality.

1016 Requirements for reporting:

- 1017 • An assessment of the models used for confounder adjustment and estimation of the treatment
1018 effects and whether any limitations exist with regard to model choice;
- 1019 • A clear description of the population in which the treatment effect has been estimated, and its
1020 relevance to the assessment question; any limitations should be clearly outlined, and potential
1021 biases arising from population differences should be reported (including an assessment of the
1022 likely magnitude and direction of any bias if possible);
- 1023 • An assessment of the effect estimates with confidence intervals for the crude data and after
1024 adjustment;
- 1025 • An assessment of the results of all sensitivity analyses;
- 1026 • If shifted hypothesis testing has been used, an assessment whether this is sufficient to account
1027 for the likely magnitude of residual bias arising from missing covariates.

1028

1029 7 FURTHER RELEVANT DOCUMENTS (UNDER DEVELOPMENT)

1030 • **EUnetHTA 21 Methodological Guideline D.4.3.2: *Direct and Indirect Comparisons***

1031 A methodological guideline for assessors and co-assessors that describes the currently
1032 available usual methods for direct and indirect comparisons regarding their underlying
1033 assumptions, strengths, and weaknesses.

1034 • **EUnetHTA 21 Practical Guideline D.4.2.1: *Scoping Process***

1035 A practical guideline for assessors and co-assessors that describes the methods and principal
1036 steps of the scoping process.

1037 • **EUnetHTA 21 Practical Guideline D4.5.1: *Applicability of Evidence: Practical Guideline*** 1038 ***on Multiplicity, Subgroup, Sensitivity and Post-hoc Analyses***

1039 A practical guideline for assessors and co-assessors that describes how to consider
1040 complementary analyses and how to handle multiplicity issues.

1041 • **EUnetHTA 21 Practical Guideline D4.6.1: *Validity of Clinical Studies***

1042 A practical guideline for assessors and co-assessors that describes possible approaches and
1043 specific instructions for action when assessing the certainty of results coming from individual
1044 studies, regardless of whether they are RCTs or other types of study.

1045 • **EUnetHTA 21 Guideline D5.3.1: *Procedural Guideline for Appointing Assessors and Co-*** 1046 ***Assessors for JCA/CA***

1047 A guidance document that describes the minimum selection criteria for the appointment of
1048 assessors and co-assessors of JCAs/collaborative assessments (CAs) for pharmaceuticals and
1049 MDs.

1050 8 FUTURE RECOMMENDATION

1051 This document was written in combination with the EUnetHTA 21 Methodological Guideline *Direct and*
1052 *Indirect Comparisons*. The Hands-on Group recommends that consideration be given to merge these
1053 documents at a later date under the direction of the future Methodological Subgroup of the HTA
1054 Regulation Coordination Group.

1055 Given that the methods used for direct and indirect comparisons might be updated more frequently
1056 compared with other methodological areas, it is advisable that these guidelines be reviewed under the
1057 direction of the future Methodological Subgroup every 3 years.

1058

1059 9 References

- 1060 1 Ali MS, Prieto-Alhambra D, Lopes LC *et al.* Propensity score methods in health technology
1061 assessment: principles, extended applications, and recent advances. *Front Pharmacol.*
1062 2019;10:973.
- 1063 2 Austin PC. Balance diagnostics for comparing the distribution of baseline covariates between
1064 treatment groups in propensity-score matched samples. *Stat Med.* 2009;28(25):3083-107.
- 1065 3 Bender R, Friede T, Koch A *et al.* Methods for evidence synthesis in the case of very few studies.
1066 *Res Synth Methods* 2018;9(3):382-92.
- 1067 4 Berlin JA, Santanna J, Schmid CH *et al.* Individual patient- versus group-level data meta-
1068 regressions for the investigation of treatment effect modifiers: ecological bias rears its ugly head.
1069 *Stat Med.* 2002;21(3):371-87.
- 1070 5 Brockhaus AC, Grouven U, Bender R. Performance of the Peto odds ratio compared to the usual
1071 odds ratio estimator in the case of rare events. *Biom J.* 2016;58(6):1428-44.
- 1072 6 Cooper NJ, Sutton AJ, Morris D *et al.* Addressing between-study heterogeneity and inconsistency
1073 in mixed treatment comparisons: application to stroke prevention treatments in individuals with non-
1074 rheumatic atrial fibrillation. *Stat Med.* 2009;28(14):1861-81.
- 1075 7 Cope S, Zhang J, Saletan S *et al.* A process for assessing the feasibility of a network meta-analysis:
1076 a case study of everolimus in combination with hormonal therapy versus chemotherapy for
1077 advanced breast cancer. *BMC Med.* 2014;12:93.
- 1078 8 Cordero CP, Dans AL. Key concepts in clinical epidemiology: detecting and dealing with
1079 heterogeneity in meta-analyses. *J Clin Epidemiol.* 2021;130:149-51.
- 1080 9 Deeks JJ, Higgins JPT, Altman DG *et al.* Analysing data and undertaking meta-analyses. In: Higgins
1081 JPT, Thomas JJC, Cumpston M *et al.* (eds). *Cochrane Handbook for Systematic Reviews of*
1082 *Interventions*, 2nd edition. Hoboken, NJ: Wiley; 2019, pp. 241-84.
- 1083 10 DerSimonian R, Laird N. Meta-analysis in clinical trials. *Control Clin Trials.* 1986;7(3):177-88.
- 1084 11 Desai RJ, Franklin JM. Alternative approaches for confounding adjustment in observational studies
1085 using weighting based on the propensity score: a primer for practitioners. *BMJ.* 2019;367:l5657.
- 1086 12 Dias S, Ades A, Welton N *et al.* *Network Meta-Analysis for Decision Making.* Chichester, UK: Wiley;
1087 2018.
- 1088 13 Dias S, Ades AE. Absolute or relative effects? Arm-based synthesis of trial data. *Res Synth*
1089 *Methods.* 2016;7(1):23-8.
- 1090 14 Dias S, Sutton AJ, Welton NJ *et al.* Evidence synthesis for decision making 3: heterogeneity –
1091 subgroups, meta-regression, bias, and bias-adjustment. *Med Decis Making.* 2013;33(5):618-40.
- 1092 15 Dias S, Welton NJ, Caldwell DM *et al.* Checking consistency in mixed treatment comparison meta-
1093 analysis. *Stat Med.* 2010;29(7-8):932-44.
- 1094 16 Dias S, Welton NJ, Sutton AJ *et al.* *NICE DSU Technical Support Document 4: Inconsistency in*
1095 *Networks of Evidence Based Upon Randomised Controlled Trials.* London, UK: National Institute
1096 for Health and Care Excellence; 2011.
- 1097 17 Faria R, Hernandez Alava M, Manca A *et al.* NICE DSU Technical Support Document 17: The Use
1098 of Observational Data to Inform Estimates of Treatment Effectiveness for Technology Appraisal:
1099 Methods for Comparative Individual Patient Data. London, UK: National Institute for Health and Care
1100 Excellence; 2015.
- 1101 18 Ference BA, Ginsberg HN, Graham I *et al.* Low-density lipoproteins cause atherosclerotic
1102 cardiovascular disease. 1. Evidence from genetic, epidemiologic, and clinical studies. A consensus
1103 statement from the European Atherosclerosis Society Consensus Panel. *Eur Heart J.*
1104 2017;38(32):2459-72.
- 1105 19 Fisher DJ, Carpenter JR, Morris TP *et al.* Meta-analytical methods to identify who benefits most
1106 from treatments: daft, deluded, or deft approach? *BMJ* 2017;356:j573.
- 1107 20 Freeman SC, Carpenter JR. Bayesian one-step IPD network meta-analysis of time-to-event data
1108 using Royston-Parmar models. *Res Synth Methods.* 2017;8(4):451-64.
- 1109 21 Greenland S. Tests for interaction in epidemiologic studies: a review and a study of power. *Stat*
1110 *Med.* 1983;2(2):243-51.
- 1111 22 Guyot P, Ades AE, Ouwens MJ *et al.* Enhanced secondary analysis of survival data: Reconstructing
1112 the data from published Kaplan-Meier survival curves. *BMC Med Res Methodol.* 2012;12:9.

- 1113 23 Hernán MA, Robins JM. Using big data to emulate a target trial when a randomized trial is not
1114 available. *Am J Epidemiol* 2016;183(8):758-64.
- 1115 24 Higgins JPT, Thompson SG. Controlling the risk of spurious findings from meta-regression. *Stat*
1116 *Med.* 2004;23(11):1663-82.
- 1117 25 Higgins JPT, Thompson SG. Quantifying heterogeneity in a meta-analysis. *Stat Med.*
1118 2002;21(11):1539-58.
- 1119 26 Hong H, Chu H, Zhang J *et al.* A Bayesian missing data framework for generalized multiple outcome
1120 mixed treatment comparisons. *Res Synth Methods.* 2016;7(1):6-22.
- 1121 27 IQWiG. *General Methods, Version 6.1.* Cologne, Germany: Institute for Quality and Efficiency in
1122 Health Care; 2022.
- 1123 28 Ishak KJ, Proskorovsky I, Benedict A. Simulation and matching-based approaches for indirect
1124 comparison of treatments. *Pharmacoeconomics* 2015;33(6):537-49.
- 1125 29 Jackson D, Law M, Rücker G *et al.* The Hartung-Knapp modification for random-effects meta-
1126 analysis: a useful refinement but are there any residual concerns? *Stat Med.* 2017;36(25):3923-34.
- 1127 30 Knapp G, Hartung J. Improved tests for a random effects meta-regression with a single covariate.
1128 *Stat Med.* 2003;22(17):2693-710.
- 1129 31 Kuss O. Statistical methods for meta-analyses including information from studies without any events
1130 – add nothing to nothing and succeed nevertheless. *Stat Med.* 2015;34(7):1097-116.
- 1131 32 Lumley T. Network meta-analysis for indirect treatment comparisons. *Stat Med.* 2002;21(16):2313-
1132 24.
- 1133 33 Mathes T, Kuss O. A comparison of methods for meta-analysis of a small number of studies with
1134 binary outcomes. *Res Synth Methods.* 2018;9(3):366-81.
- 1135 34 Mick R, Chen TT. Statistical challenges in the design of late-stage cancer immunotherapy studies.
1136 *Cancer Immunol Res.* 2015;3(12):1292-8.
- 1137 35 Phillippo DM. *Calibration of Treatment Effects in Network Meta-Analysis using Individual Patient*
1138 *Data, PhD Thesis.* Bristol, UK: Bristol Medical School; 2019.
- 1139 36 Phillippo DM, Ades AE, Dias S *et al.* Methods for population-adjusted indirect comparisons in health
1140 technology appraisal. *Med Decis Making.* 2018;38(2):200-11.
- 1141 37 Phillippo DM, Ades AE, Dias S *et al.* *NICE DSU Technical Support Document 18: Methods for*
1142 *Population-Adjusted Indirect Comparisons in Submission to NICE.* London, UK: National Institute
1143 for Health and Care Excellence; 2016.
- 1144 38 Phillippo DM, Dias S, Ades AE *et al.* Multilevel network meta-regression for population-adjusted
1145 treatment comparisons. *J R Stat Soc Ser A Stat Soc.* 2020;183(3):1189-210.
- 1146 39 Phillippo DM, Dias S, Ades AE *et al.* Assessing the performance of population adjustment methods
1147 for anchored indirect comparisons: a simulation study. *Stat Med.* 2020;39(30):4885-911.
- 1148 40 Phillippo DM, Dias S, Ades AE *et al.* Target estimands for efficient decision making: Response to
1149 comments on "Assessing the performance of population adjustment methods for anchored indirect
1150 comparisons: a simulation study". *Stat Med.* 2021;40(11):2759-63.
- 1151 41 Remiro-Azócar A, Heath A, Baio G. Conflating marginal and conditional treatment effects:
1152 comments on "Assessing the performance of population adjustment methods for anchored indirect
1153 comparisons: a simulation study". *Stat Med* 2021;40(11):2753-8.
- 1154 42 Remiro-Azócar A, Heath A, Baio G. Methods for population adjustment with limited access to
1155 individual patient data: a review and simulation study. *Res Synth Methods* 2021;12(6):750-75.
- 1156 43 Remiro-Azócar A, Heath A, Baio G. Parametric G-computation for compatible indirect treatment
1157 comparisons with limited individual patient data. *Res Synth Methods.* Published online April 29,
1158 2022. <https://dx.doi.org/10.1002/jrsm.1565>.
- 1159 44 Röver C, Bender R, Dias S *et al.* On weakly informative prior distributions for the heterogeneity
1160 parameter in Bayesian random-effects meta-analysis. *Res Synth Methods.* 2021;12(4):448-74.
- 1161 45 Royston P, Parmar MK. Restricted mean survival time: an alternative to the hazard ratio for the
1162 design and analysis of randomized trials with a time-to-event outcome. *BMC Med Res Methodol.*
1163 2013;13:152.
- 1164 46 Salanti G, Higgins JP, Ades AE *et al.* Evaluation of networks of randomized trials. *Stat Methods*
1165 *Med Res.* 2008;17(3):279-301.
- 1166 47 Schemper M. Cox analysis of survival data with non-proportional hazard functions. *Statistician.*
1167 1992;41:455-65.

- 1168 48 Schulz A, Schürmann C, Skipka G *et al.* Performing meta-analyses with very few studies. In:
1169 Evangelou E, Veroniki AA (eds). *Meta-Research: Methods and Protocols*. New York: Humana;
1170 2022, pp. 91-102.
- 1171 49 Sterne JA, Hernán MA, Reeves BC *et al.* ROBINS-I: a tool for assessing risk of bias in non-
1172 randomised studies of interventions. *BMJ*. 2016;355:i4919.
- 1173 50 Stijnen T, Hamza TH, Ozdemir P. Random effects meta-analysis of event outcome in the framework
1174 of the generalized linear mixed model with applications in sparse data. *Stat Med*. 2010;29(29):3046-
1175 67.
- 1176 51 Stuart EA, Lee BK, Leacy FP. Prognostic score-based balance measures can be a useful diagnostic
1177 for propensity score methods in comparative effectiveness research. *J Clin Epidemiol*. 2013;66(8
1178 Suppl):S84-S90.
- 1179 52 Sutton AJ, Abrams KR. Bayesian methods in meta-analysis and evidence synthesis. *Stat Methods
1180 Med Res*. 2001;10(4):277-303.
- 1181 53 Sutton AJ, Abrams KR, Jones DR *et al.* *Methods for Meta-Analysis in Medical Research*. Chichester,
1182 UK: Wiley; 2000.
- 1183 54 Sweeting MJ, Sutton AJ, Lambert PC. What to add to nothing? Use and avoidance of continuity
1184 corrections in meta-analysis of sparse data. *Stat Med*. 2004;23(9):1351-75.
- 1185 55 Veroniki AA, Jackson D, Bender R *et al.* Methods to calculate uncertainty in the estimated overall
1186 effect size from a random-effects meta-analysis. *Res Synth Methods* 2019;10(1):23-43.
- 1187 56 Wang SV, Pinheiro S, Hua W *et al.* STaRT-RWE: structured template for planning and reporting on
1188 the implementation of real world evidence studies. *BMJ* 2021;372:m4856.
- 1189 57 White IR, Turner RM, Karahalios A *et al.* A comparison of arm-based and contrast-based models
1190 for network meta-analysis. *Stat Med* 2019;38(27):5197-213.
- 1191 58 Wiksten A, Rücker G, Schwarzer G. Hartung-Knapp method is not always conservative compared
1192 with fixed-effect meta-analysis. *Stat Med* 2016;35(15):2503-15.
- 1193