



1 15 December 2016  
2 EMA/CHMP/44762/2017  
3 Committee for Human Medicinal Products (CHMP)

## 4 Guideline on multiplicity issues in clinical trials

5 Draft

Draft agreed by Biostatistics Working Party (BSWP)	November 2016
Adopted by CHMP for release for consultation	15 December 2016
Start of public consultation	01 April 2017
End of consultation (deadline for comments)	30 June 2017

6  
7 This guideline replaces the 'Points to consider on multiplicity issues in clinical trials'  
8 (CPMP/EWP/908/99).

9  
10 Comments should be provided using this [template](#). The completed comments form should be sent to [Multiplicity\\_GL@ema.europa.eu](mailto:Multiplicity_GL@ema.europa.eu).

<b>Keywords</b>	<b><i>Multiplicity, hypothesis test, type I error, subgroup, responder, estimation, confidence interval</i></b>
-----------------	---

11  
12



## 14 Guideline on multiplicity issues in clinical trials

### 15 Table of contents

16	<b>1. Executive summary</b> .....	<b>3</b>
17	<b>2. Introduction</b> .....	<b>4</b>
18	<b>3. Scope</b> .....	<b>4</b>
19	<b>4. Legal basis and other relevant guidance documents</b> .....	<b>5</b>
20	<b>5. Adjustment of elementary hypothesis tests for multiplicity – when is it</b>	
21	<b>necessary and when is it not?</b> .....	<b>5</b>
22	5.1. Multiple primary endpoints – when no formal adjustment of the significance level is	
23	needed .....	6
24	5.1.1. Two or more primary endpoints are needed to describe clinically relevant treatment	
25	benefits .....	6
26	5.1.2. Two or more endpoints ranked according to clinical relevance .....	7
27	5.2. Analysis sets.....	7
28	5.3. Alternative statistical methods – multiplicity concerns .....	7
29	5.4. Multiplicity in safety variables .....	8
30	5.5. Multiplicity concerns in studies with more than two treatment arms.....	8
31	5.5.1. The three arm ‘gold standard’ design .....	9
32	5.5.2. Proof of efficacy for a fixed combination .....	9
33	5.5.3. Dose-response studies .....	9
34	<b>6. How to interpret significance with respect to multiple secondary</b>	
35	<b>endpoints and when can a regulatory claim be based on one of these?.....</b>	<b>10</b>
36	6.1. Secondary endpoints expressing supportive evidence .....	10
37	6.2. Secondary endpoints which may become the basis for additional claims.....	11
38	6.3. Secondary endpoints indicative of clinical benefit.....	11
39	<b>7. Reliable conclusions from a subgroup analysis, and restriction of the</b>	
40	<b>licence to a subgroup</b> .....	<b>11</b>
41	<b>8. How should one interpret the analysis of ‘responders’ in conjunction with</b>	
42	<b>the raw variables?</b> .....	<b>12</b>
43	<b>9. How should composite endpoints be handled statistically with respect to</b>	
44	<b>regulatory claims?</b> .....	<b>12</b>
45	9.1. The composite endpoint as the primary endpoint.....	13
46	9.2. Treatment should be expected to affect all components in a similar way .....	13
47	9.3. The clinically more important components should at least not be affected negatively ..	14
48	9.4. Any effect of the treatment on one of the components that is intended to be reflected in	
49	the product information should be clearly supported by the data.....	14
50	<b>10. Multiplicity issues in estimation</b> .....	<b>14</b>
51	10.1. Selection bias.....	15
52	10.2. Confidence intervals.....	15

53

## 54 **1. Executive summary**

55 This guideline is intended to provide guidance on how to deal with multiple comparison and control of  
56 type I error in the planning and statistical analysis of clinical trials.

57 In 2002 the EMA Points to Consider on Multiplicity issues in clinical trials (EMA/286914/2012) was  
58 adopted. Following the EMA Concept paper on the need for a guideline on multiplicity issues in clinical  
59 trials which was published in 2012, this guideline was developed as an update of the above mentioned  
60 Points to Consider considering new regulatory advisements, including a new section on multiplicity in  
61 estimation, accounting for new approaches in dose finding and clarifying specific issues and  
62 applications.

63 The present document should be considered as a general guidance. The main considerations for  
64 multiplicity issues encountered in clinical trials are described. Specific issues, including adjustment of  
65 elementary hypothesis tests for multiplicity, multiple primary endpoints, analysis sets and alternative  
66 statistical methods are addressed.

67 The main scope is to provide guidance on the confirmatory conclusions which are usually based on the  
68 results from pivotal Phase III trials and, to a lesser extent, on Phase II studies. The guideline mainly  
69 discusses issues in decision making for a formal proof of efficacy.

70 In clinical studies it is often necessary to answer more than one question about the efficacy (or safety)  
71 of the experimental treatment in a specific disease, because the success of a drug development  
72 programme may depend on a positive answer to more than a single question. It is well known that the  
73 likelihood of a positive chance finding increases with the number of questions posed, if no actions are  
74 taken to protect against the inflation of false positive findings from multiple statistical tests. In this  
75 context, concern is focused on the opportunity to choose favourable results from multiple analyses. It  
76 is therefore necessary that the statistical procedures planned to deal with, or to avoid, multiplicity are  
77 fully detailed in the study protocol or in the statistical analysis plan to allow an assessment of their  
78 suitability and appropriateness.

79 Various methods have been developed to control the rate of false positive findings. Not all of these  
80 methods, however, are equally successful at providing clinically interpretable results and this aspect of  
81 the procedure should always be considered. Since estimation of treatment effects is usually an  
82 important issue, the availability of confidence intervals with correct coverage that allow for consistent  
83 decision making with the primary hypothesis testing strategy may be a criterion for the selection of the  
84 corresponding multiple testing procedure.

85 Additional claims on statistically significant and clinically relevant findings based on secondary  
86 endpoints or on subgroups are formally possible only after the primary objective of the clinical trial has  
87 been achieved ('claim' is used as shorthand for a confirmatory conclusion which is then prioritised in  
88 trial reporting and used as primary basis for asserting that efficacy or safety has been established),  
89 and if the respective questions were pre-specified, and were part of an appropriately planned statistical  
90 analysis strategy.

91

92 This document should be read in conjunction with other applicable EU and ICH guidelines (see Section  
93 4).

## 94 2. Introduction

95 Multiplicity of inferences is present in virtually all clinical trials. The usual concern with multiplicity is  
96 that, if it is not properly handled, unsubstantiated claims for the efficacy of a drug may be made as a  
97 consequence of an inflated rate of false positive conclusions. For example, if statistical tests are  
98 performed on five subgroups, independently of each other and each at a significance level of 2.5%  
99 (one-sided directional hypotheses), the chance of finding at least one false positive statistically  
100 significant test increases to approximately 12%.

101 This example shows that multiplicity can have a substantial influence on the rate of false positive  
102 conclusions which may affect approval and labelling of an investigational drug whenever there is an  
103 opportunity to choose the most favourable result from two or more analyses. If, however, there is no  
104 such choice, then there can be no influence. Examples of both situations will be discussed later.  
105 Control of the study-wise rate of false positive conclusions at an acceptable level  $\alpha$  is an important  
106 principle and is often of great value in the assessment of the results of confirmatory clinical trials.

107 A number of methods are available for controlling the rate of false positive conclusions, the method of  
108 choice depending on the circumstances. Throughout this document the term 'control of type I error'  
109 rate will be used as an abbreviation for the control of the study-wise type I error in the strong sense,  
110 *i.e.* there is control on the probability to reject at least one out of several true null hypotheses,  
111 regardless of which subset of null hypotheses happens to be true.

## 112 3. Scope

113 The scope of this guideline is to provide guidance on the confirmatory conclusions which are usually  
114 based on the results from pivotal Phase III trials and, to a lesser extent, on Phase II studies. The  
115 guideline mainly discusses issues in decision making for a formal proof of efficacy. Due to the  
116 precautionary principle in safety evaluations, reducing the rate of false negative conclusions on harm is  
117 usually more important than controlling the number of false positive conclusions and rigorous  
118 multiplicity adjustments could mask relevant safety signals.

119 The principles discussed in this guideline follow the frequentist approach in statistical decision theory,  
120 where the validity of a confirmatory conclusion is defined by limiting the probability of a false positive  
121 conclusion relating to data sampling and pre-defined statistical procedures of a specific study at a pre-  
122 specified level  $\alpha$ . The CHMP Points to Consider on Application with 1. Meta-analyses and 2. One Pivotal  
123 Study (CPMP/2330/99) covers the situation when the type I error needs to be controlled at the  
124 submission level where more than one confirmatory trial is included in a submission.

125 This document does not attempt to address all aspects of multiplicity but mainly considers issues that  
126 have been found to be of importance in European marketing authorisation applications. These are:

- 127 • Adjustment of multiplicity – when is it necessary and when is it not?
- 128 • How to interpret significance with respect to multiple secondary endpoints and when can a  
129 regulatory claim be based on one of these?
- 130 • When can confirmatory conclusions be drawn from a subgroup analysis?
- 131 • How should one interpret the analysis of 'responders' in conjunction with the analysis of raw  
132 variables and how should composite endpoints be handled statistically with respect to  
133 regulatory claims?
- 134 • How should multiplicity issues be addressed in estimation?

135 There are further areas concerning multiplicity in clinical trials which, according to the above list of  
136 issues, are not the focus of this document. For example, there is a rapid advance in methodological  
137 richness and complexity regarding interim analyses, with the possibility to stop early either for futility  
138 or with a claim for efficacy, or stepwise designed studies, with the possibility for adaptive changes in  
139 the trial's next steps. However, due to the importance of the problem and the amount of information  
140 specific to this issue these aspects are discussed in the CHMP Reflection Paper on Methodological issues  
141 in Confirmatory Clinical Trials planned with an Adaptive Design (CHMP/EWP/2459/02).

142 Interpretations of evaluations of the primary efficacy variable at repeated visits per patient usually do  
143 not cause multiplicity problems, because in the majority of situations either an appropriate summary  
144 measure has been pre-specified or according to the requirements on the duration of treatment,  
145 primary evaluations are made at a pre-specified visit. Therefore potential multiplicity issues concerning  
146 the analysis of repeated measurements are not considered in this document.

147

#### 148 **4. Legal basis and other relevant guidance documents**

149 This guideline has to be read in conjunction with Directive 2001/83 as amended and other applicable  
150 EU and ICH guidance documents, especially:

151 Note for Guidance on Dose-Response Information to Support Drug Registration - CPMP/ICH/378/95  
152 (ICH E4)

153 Note for Guidance on Statistical Principles for Clinical Trials - CPMP/ICH/363/96 (ICH E9)

154 Guideline on the choice of the non-inferiority margin - CPMP/EWP/2158/99

155 Guideline on the Investigation of subgroups in confirmatory clinical trials - EMA/CHMP/539146/2013

156 Guideline on Clinical Development of Fixed Combination Medicinal Products – EMA/CHMP/281825/2015

157 Points to Consider on Application with 1. Meta-analyses and 2. One Pivotal study - CPMP/2330/99

158 Reflection paper on methodological issues in confirmatory clinical trials planned with an adaptive  
159 design - CHMP/EWP/2459/02

160

#### 161 **5. Adjustment of elementary hypothesis tests for multiplicity** 162 **– when is it necessary and when is it not?**

163 A clinical study that requires no adjustment of the significance level of elementary hypothesis tests  
164 (*i.e.* single statistical tests on one parameter only) is one that consists of two treatment groups, which  
165 uses a single primary variable, and has a confirmatory statistical strategy that pre-specifies just one  
166 single null hypothesis relating to the primary variable and no interim analysis. Although all other  
167 situations require attention to the potential effects of multiplicity, there are situations where no  
168 multiplicity concern arises, for example, having a number of primary hypotheses for a number of  
169 primary endpoints that all need to be significant so that the trial is considered successful, and all other  
170 endpoints are declared supportive. The assessor should expect to find in the protocol and analysis plan  
171 a discussion on the aspects of trial design, conduct and analysis that give rise to multiple testing and  
172 the proposed strategy for controlling the study-wise rate of false positive confirmatory conclusions.

173 Methods to control the overall type I error rate  $\alpha$  are sometimes called multiple-level- $\alpha$  tests.

174 Controlling the type I error rate study-wise is frequently done by splitting the accepted and pre-

175 specified type I error rate  $\alpha$  and by then testing the various null hypotheses at fractions of  $\alpha$ . This is  
176 usually referred to as 'adjusting the local significance level' (*i.e.* adjusting the significance level of each  
177 test). Other test procedures are available, that can be more powerful if the correlation between the  
178 test statistics are taken into account, *e.g.* the Dunnett's test on multiple comparisons to a single  
179 control. The algorithms that define how to 'spend'  $\alpha$  are of different complexity.

180 In general, more than one approach is available to correctly deal with multiplicity issues. These  
181 different methods may lead to different conclusions and for this reason the details of the chosen  
182 multiplicity procedure should be part of the study protocol and should be written up without room for  
183 choice.

## 184 **5.1. Multiple primary endpoints – when no formal adjustment of the** 185 **significance level is needed**

186 The ICH E9 guideline on statistical principles for clinical trials recommends that generally clinical trials  
187 have one primary variable. A single primary variable is sufficient, if there is a general agreement that a  
188 treatment induced change in this variable demonstrates a clinically relevant treatment effect on its  
189 own. If, however, a single variable is not sufficient to capture the range of clinically relevant treatment  
190 benefits, the use of more than one primary variable may become necessary. Sometimes a series of  
191 related objectives is pursued in the same trial, each with its own primary variable, and in other cases,  
192 a number of primary endpoints are investigated with the aim of providing convincing evidence of  
193 beneficial effects on some, or all of them. In these situations planning of the sample size becomes  
194 more complex due to the different alternative hypotheses related to the different endpoints and due to  
195 the assumed correlation between endpoints.

196 If more than one primary endpoint is used to define study success, this success could be defined by a  
197 positive outcome in all endpoints or it may be considered sufficient, if one out of a number of  
198 endpoints has a positive outcome. Whereas in the first definition the primary endpoints are designated  
199 as co-primary endpoints, the latter case is different and would require appropriate adjustment for  
200 multiplicity. More generally, in case of more than two primary endpoints, adjustment is needed if not  
201 all endpoints need to be significant to define study success, and the inability to exclude deteriorations  
202 in other primary endpoints would have to be considered in the overall benefit/risk assessment.

203 For trials with more than one primary variable the situations described in the following subsections can  
204 be distinguished. The methods described allow clinical interpretation, deal satisfactorily with the issue  
205 of multiplicity but avoid the need for any formal adjustment of type I error rates. Other methods of  
206 dealing with multiple variables, that are more complex, are possible and can be found in the literature.  
207 In general, regulatory dialogue is recommended before applying these methods.

### 208 **5.1.1. Two or more primary endpoints are needed to describe clinically** 209 **relevant treatment benefits**

210 *Statistical significance is needed for all primary endpoints. Therefore, no formal adjustment of the*  
211 *significance level of the elementary hypothesis tests is necessary.*

212 Here, interpretation of the results is most clear-cut because, in order to provide sufficient evidence of  
213 the clinically relevant efficacy, each null hypothesis on every primary variable has to be rejected at the  
214 same significance level (*e.g.* 0.05). For example, according to the CHMP Guideline on clinical  
215 investigation of medicinal products in the treatment of chronic obstructive pulmonary disease  
216 (EMA/CHMP/483572/2012), lung function would be insufficient as a single primary endpoint and should  
217 be accompanied by an additional co-primary endpoint, which should either be a symptom-based  
218 endpoint or a patient-related endpoint.

219 In these situations, there is no intention or opportunity to select the most favourable result and,  
220 consequently, the individual significance levels are set equal to the overall significance level  $\alpha$ , *i.e.* no  
221 adjustment is necessary. Even though in this situation all hypotheses can be assessed at the same  
222 type I error level, the need for a significant result for more than one primary hypothesis will reduce the  
223 power of the statistical procedure or increase the sample size that is needed for a given power. This  
224 inflation must be taken into account for a proper estimation of the sample size for the trial.

### 225 **5.1.2. Two or more endpoints ranked according to clinical relevance**

226 *No numerical adjustment of each single hypothesis test is necessary. However, no confirmatory claims*  
227 *can be based on endpoints that have a rank lower than or equal to that variable whose null hypothesis*  
228 *was the first that could not be rejected.*

229 Sometimes a series of related objectives is pursued in the same trial, where one objective is of  
230 greatest importance but convincing results in others would clearly add to the value of the treatment. A  
231 typical example is the reduction of mortality in acute myocardial infarction followed by prevention of  
232 other serious events. In such cases the hypotheses may be tested (and confidence intervals may be  
233 provided) according to a hierarchical strategy. The hierarchical order may be a natural one (*e.g.*  
234 hypotheses are ordered in time or with respect to the importance of the considered endpoints) or may  
235 result from the particular interests of the investigator. Hierarchical testing can be considered as a  
236 specific multiplicity procedure. Although such a procedure may be considered as a particular  
237 adjustment, no reduction or splitting of the single  $\alpha$  levels is necessary since the pre-defined ordering  
238 avoids any choice in the assessment. The hierarchical order for testing null hypotheses, however, has  
239 to be pre-specified in the study protocol, including a clear specification of the set of hypotheses that  
240 need to be significant before the trial is claimed successful. The effect of such a procedure is that no  
241 confirmatory claims can be based on endpoints that have a rank lower than or equal to that variable  
242 whose null hypothesis was the first that could not be rejected. Evidently, type II errors are inflated for  
243 hypotheses that correspond to endpoints with lower ranks. Note that a similar procedure can be used  
244 for dealing with secondary endpoints (see Section 6.2).

### 245 **5.2. Analysis sets**

246 Multiple analyses may be performed on the same variable but with varying subsets of patient data. As  
247 is pointed out in ICH E9, the set of subjects whose data are to be included in the main analyses should  
248 be defined in the statistical section of the study protocol. From these sets of subjects one (usually the  
249 full set) is selected for the primary analysis.

250 In general, multiple additional analyses on varying subsets of subjects or with varying measurements  
251 for the purpose of investigating the robustness of the conclusions drawn from the primary analysis  
252 should not be subjected to adjustment for type I error (in contrast, however, to the confirmatory  
253 subgroup analyses described in Section 7, see also CHMP Guideline on the Investigation of subgroups  
254 in confirmatory clinical trials (EMA/CHMP/539146/2013)). The main purpose of such analyses is to  
255 increase confidence in the results obtained from the primary analysis.

### 256 **5.3. Alternative statistical methods – multiplicity concerns**

257 Different statistical models or statistical techniques (*e.g.* parametric vs. non-parametric or Wilcoxon  
258 test versus log-rank test) are sometimes tried on the same set of data. A two-step procedure may be  
259 applied with the purpose of selecting a particular statistical technique for the main treatment  
260 comparison based on the outcome of the first statistical (pre-)test, the first one of the two steps.  
261 Multiplicity concerns would immediately arise, if such procedures offered obvious opportunities for

262 selecting a favourable analysis strategy based on knowledge of the patients' assignment to treatments.  
263 In other words, the correct type I error rate refers to the overall procedure that includes the pre-test  
264 and the selected test, and therefore such a two-test procedure does not usually control the type I  
265 error. Opportunities for choice in such procedures are often subtle, especially when these procedures  
266 use comparative treatment information, and the influence on the overall type I error is difficult to  
267 assess. Applying the same line of thought, type I error control for analyses that include model selection  
268 procedures should be based on the overall procedure. Type I error control on the basis of the finally  
269 selected model only is usually not sufficient. In addition, any *post hoc* selection of the model is not  
270 considered appropriate for a confirmatory Phase III trial.

271 In some situations the selected statistical model is based on a formal blind review, *i.e.* on the basis of  
272 the pooled data set from the different treatment groups hiding the information on the allocated  
273 treatment. It is also important in this case that there is no inflation in the type I error. Therefore, the  
274 selection of the statistical model according to the results of a blinded analysis should be properly  
275 justified with respect to type I error control and its potential impact on the treatment effect estimate  
276 as regards bias.

277 In summary, the need to change or define important key features of a study on a *post hoc* basis may  
278 question the credibility of the study and the robustness of the results with the possible consequence  
279 that a further study will be necessary. Therefore, such procedures are not recommended. Confirmatory  
280 analyses should be fully and precisely pre-defined to exclude the possibility of performing different  
281 analyses *post hoc*.

#### 282 **5.4. Multiplicity in safety variables**

283 When a safety variable is part of the confirmatory strategy of a study and thus has a role in the  
284 approval or labelling claims, it should not be treated differently from the primary efficacy endpoints,  
285 except for the situation that the observed effects go in the opposite direction and may raise a safety  
286 concern (see also Section 9.3).

287 In the case of adverse effects, p-values are of very limited value as substantial differences (expressed  
288 as relative risk or risk differences) require careful assessment and will in addition raise concern,  
289 depending on seriousness, severity or outcome, irrespective of the p-value observed. A non-significant  
290 difference between treatments will not allow for a conclusion on the absence of a difference in safety.  
291 In other words, in line with general principles, a non-significant test result should not be confused with  
292 the demonstration of equivalence.

293 In those cases where a large number of statistical test procedures are performed to serve as a flagging  
294 device to signal a potential risk caused by the investigational drug it can generally be stated that an  
295 adjustment for multiplicity is counterproductive for considerations of safety. It is likewise clear that in  
296 this situation there is no control of the type I error for a single hypothesis and the importance and  
297 plausibility of 'significant findings' will depend on prior knowledge of the pharmacology of the drug, and  
298 sometimes further investigations may be required.

#### 299 **5.5. Multiplicity concerns in studies with more than two treatment arms**

300 As for studies with more than one primary endpoint, the proper evaluation and interpretation of a  
301 study with more than two treatment arms can become quite complex. This document is not intended to  
302 provide an exhaustive discussion of every issue relating to studies with multiple treatment arms.  
303 Therefore, the following discussion is limited to the more common and simple designs. As a general  
304 rule it can be stated that control of the study-wise type I error is a minimal prerequisite for  
305 confirmatory claims.



### 306 **5.5.1. The three arm 'gold standard' design**

307 For a disease, where a commonly acknowledged reference drug therapy exists, it is often  
308 recommended (when this can be justified on ethical grounds) to demonstrate the efficacy and safety of  
309 a new substance in a three-arm study with the reference drug, placebo and the investigational drug.  
310 Ideally, though not exclusively, the aims of such a study are to demonstrate superiority of the  
311 investigational drug over placebo (proof of efficacy) and to show that the investigational drug retains,  
312 at least, most of the efficacy of the reference drug as compared to placebo (proof of non-inferiority). If  
313 study success is defined by non-inferiority to the reference product combined with superiority to  
314 placebo both comparisons must show statistical significance at the required level and no formal  
315 adjustment of the significance level for the single hypotheses tests is necessary. In some settings,  
316 however, superiority to placebo is the main criterion for approval, and the comparison to the reference  
317 is not considered to be primary. In this case study success could be based on a significant superiority  
318 to placebo only, but any additional confirmatory conclusion on non-inferiority to the reference would  
319 require a pre-specified multiplicity procedure, e.g. a hierarchical procedure testing superiority to  
320 placebo first followed by a test on non-inferiority to the reference.

### 321 **5.5.2. Proof of efficacy for a fixed combination**

322 For fixed combination medicinal products the corresponding CPMP guideline (CPMP/EWP/240/95 Rev.  
323 1) requires that "each substance of a fixed combination must have documented contribution within the  
324 combination". For a combination with two (mono) components, this requirement has often been  
325 interpreted as the need to conduct a study with the two components as monotherapies and the  
326 combination therapy in a three-arm study (or a four-arm study including placebo in some settings). In  
327 case the intended contribution of the fixed combination is to improve efficacy, such a study is  
328 considered successful if the combination is shown superior to both components; no formal adjustment  
329 of the significance level for the single hypothesis tests is necessary, because there is obviously no  
330 alternative.

331 Multiple-dose factorial designs are employed for the assessment of combination drugs for the purpose  
332 (1) to provide confirmatory evidence that the combination is more effective than either component  
333 drug alone (see ICH E4 Note for Guidance on Dose Response Information to support Drug Registration  
334 (CPMP/ICH/378/95)), and (2) to identify an effective and safe dose combination (or a range of dose  
335 combinations) for recommended use in the intended patient population. While (1) usually is achieved  
336 using global test strategies, multiplicity has to be addressed for the purpose of achieving (2).

### 337 **5.5.3. Dose-response studies**

338 Phase II dose-finding studies are usually designed to estimate the dose-response relationship, e.g.  
339 with an appropriate regression model, that could be used to reasonably estimate an appropriate dose.  
340 Usually the statistical inference should focus on estimation rather than on testing, and a procedure that  
341 selects the lowest dose that shows a statistically significant difference to placebo is often of limited  
342 value and can be misleading. Therefore, the multiplicity adjustment of the different comparisons  
343 between groups in order to control the study-wise type I error may not be required in a Phase II trial.  
344 A valuable achievement in such a trial is the demonstration of an overall positive correlation of the  
345 clinical effect with increasing dose (see ICH E4, Section 3.1). Estimates and confidence intervals of the  
346 relevant parameters in the regression models are used for an appropriate interpretation of the dose  
347 response and may be used for the planning of future studies. ICH E4 also mentions under which  
348 circumstances a dose-response study can be part of the confirmatory package and in this instance a  
349 pre-specified plan to control the type I error is of importance.

350 However, for pivotal Phase III studies that use several dose groups and aim at selecting and  
351 confirming one or several doses of an investigational drug for its recommended use in a specific patient  
352 population, control of the study-wise type I error is mandatory. Due to the large variety of design  
353 features, assumptions and aims in such studies, specific recommendations are beyond the scope of  
354 this document. There are various methods published in the relevant literature on test procedures with  
355 relevance to these studies that can be adapted to the specific aims and that provide the necessary  
356 control of the type I error.

## 357 **6. How to interpret significance with respect to multiple** 358 **secondary endpoints and when can a regulatory claim be** 359 **based on one of these?**

360 Multiple secondary endpoints are included in virtually all clinical trials. These secondary endpoints will  
361 usually be included with the objective of adding weight in support of the primary efficacy claim (see  
362 Section 6.1). On occasion the secondary endpoints will be included to support a second efficacy claim  
363 (see Section 6.2). For example a symptomatic effect may be a different claim from a disease-  
364 modifying effect, and treatment and maintenance of effect may be thought of as different claims. For  
365 the purpose of this document, and distinguishing between the two sub-sections below, a claim can be  
366 thought of as a confirmatory conclusion of therapeutic efficacy or safety in a particular treatment  
367 context. The reader should not directly relate use of the word claim with the possibility to make  
368 statements or present data in the Summary of Product Characteristics, which is governed by a  
369 separate regulatory guidance document. Instead, 'claim' is used as shorthand for a confirmatory  
370 conclusion which is then prioritised in a clinical study report, clinical overview or clinical summary, and  
371 is used as a primary basis for asserting that efficacy or safety has been established.

### 372 **6.1. Secondary endpoints expressing supportive evidence**

373 *No claims are intended; confidence intervals and statistical tests are of descriptive nature.*

374 Secondary endpoints may provide additional clinical characterisation of treatment effects but are, by  
375 themselves, not sufficiently convincing to establish the main evidence in an application for a licence or  
376 for an additional labelling claim. Here, the inclusion of secondary endpoints is intended to yield  
377 supportive evidence related to the primary objective, and no confirmatory conclusions are needed.  
378 Confidence intervals and statistical tests are of descriptive nature and no claims are intended.

379 Including secondary endpoints in a multiple testing procedure (*e.g.* a 'hierarchy') is therefore not  
380 mandated, but permits a quantification of the risk of a type I error regarding these endpoints, which  
381 may lend support that an individual result is sufficiently reliable when included in the Summary of  
382 Product Characteristics.

383 The ranking of endpoints in a hierarchy can be a source of controversy. In principle, the planning and  
384 assessment of a clinical trial should prioritise those endpoints of greatest interest from a clinical  
385 perspective, but it has become common practice to rank endpoints based on the likelihood that the  
386 individual null hypothesis can be rejected. Ideally the clinical assessment should focus on those  
387 endpoints of greater clinical importance and the sponsor runs a risk of type II error if the more  
388 clinically important endpoint is set below another endpoint in the hierarchy for which the individual null  
389 hypothesis is not rejected.

390 In the event that no formal multiple testing procedure is utilised, it can still be advantageous to specify  
391 a few key secondary endpoints in the protocol that are of greater importance for assessment since  
392 selection of positive results from an unstructured list of secondary endpoints would not generally be

393 considered to provide data that are reliable for inference or for presentation in the Summary of Product  
394 Characteristics.

## 395 **6.2. Secondary endpoints which may become the basis for additional** 396 **claims**

397 *Significant effects in these endpoints can be considered for an additional claim only after the primary*  
398 *objective of the clinical trial has been achieved, and if they were part of the confirmatory strategy.*

399 Secondary endpoints may be related to secondary objectives that become the basis for an additional  
400 claim, once the primary objective has been established (see Section 5.1.2). A possible simple  
401 procedure to deal with this kind of secondary endpoint is to proceed hierarchically; other procedures  
402 are also available. Once the null hypothesis concerning the primary objective is rejected (and the  
403 primary objective is thus established), further confirmatory statistical tests on secondary endpoints can  
404 be performed using a hierarchical order for the secondary endpoints if there is more than one. In this  
405 case, primary and secondary endpoints differ just in their place in the hierarchy of hypotheses which,  
406 of course, reflects their relative importance in the study. However, more complex methods exist to  
407 control type I error over both primary and secondary endpoints, and these could be more useful in  
408 some circumstances. Depending on the degree of complexity, regulatory dialogue is recommended to  
409 assure that the outcome of the procedure can be interpreted in clinical terms.

## 410 **6.3. Secondary endpoints indicative of clinical benefit**

411 *If not defined as primary endpoints, clinically very important endpoints (e.g. mortality) need further*  
412 *study when significant benefits are observed, but the primary objective has not been achieved.*

413 Endpoints that have the potential of being indicative of a major clinical benefit or may in a different  
414 situation present an important safety issue (e.g. mortality) may be relegated to secondary endpoints  
415 because there is an *a priori* belief that the size of the planned trial is too small (and thus the power too  
416 low) to show a benefit. If, however, the observed beneficial effect is much higher than expected but  
417 the study falls short of achieving its primary objective, this would be a typical situation where  
418 information from further studies would be needed to support the observed beneficial effect.

419 If, however, the same endpoint that may indicate a major clinical benefit exhibits a treatment effect in  
420 the opposite direction, this would give rise to safety concerns (in the example of increased mortality).  
421 A Marketing Authorisation may not be granted, regardless of whether or not this endpoint was  
422 embedded in a confirmatory scheme.

## 423 **7. Reliable conclusions from a subgroup analysis, and** 424 **restriction of the licence to a subgroup**

425 *Reliable conclusions from subgroup analyses generally require pre-specification and appropriate*  
426 *statistical analysis strategies. A licence may be restricted if unexplained strong heterogeneity is found*  
427 *in important sub-populations, or if heterogeneity of the treatment effect can reasonably be assumed*  
428 *but cannot be sufficiently evaluated for important sub-populations.*

429 In clinical trials there are many reasons for examining treatment effects in subgroups. In many  
430 studies, subgroup analyses have a supportive or exploratory role after the primary objective has been  
431 accomplished. A specific claim of a beneficial effect in a particular subgroup requires pre-specification  
432 of the corresponding null hypothesis (including the precise definition of the subgroup) and an  
433 appropriate confirmatory analysis strategy. Multiplicity issues arise if study success is defined by the  
434 demonstration of a beneficial effect of the treatment in the whole study population or in a pre-defined

435 subgroup (or in one of several subgroups). An appropriate pre-planned multiplicity adjustment is  
436 needed for an unambiguous confirmatory conclusion. The complexity of the multiplicity procedure is  
437 increased if decision making is possible at an interim time point or after the final analysis. The number  
438 of subgroups should be small, in order to efficiently apply an appropriate multiplicity procedure.

439 Considerations of power are expected to be covered in the protocol, and randomisation would generally  
440 be stratified by the most important explanatory covariates. Decision making based on subgroup  
441 analyses in general are dealt with in the CHMP guideline on the Investigation of Subgroups in  
442 Confirmatory Clinical Trials (EMA/CHMP/539146/2013).

## 443 **8. How should one interpret the analysis of 'responders' in** 444 **conjunction with the raw variables?**

445 *If the 'responder' analysis is not the primary analysis it may be used after statistical significance has*  
446 *been established on the mean level of the required primary endpoint(s), to establish the clinical*  
447 *relevance of the observed differences in the proportion of 'responders'. When used in this manner, the*  
448 *test of the null hypothesis of no treatment effect is better carried out on the original primary variable*  
449 *than on the proportion of responders.*

450 In a number of applications, for example those concerned with Alzheimer's disease or depressive  
451 disorders, it may be difficult to interpret small but statistically significant improvements in the mean  
452 level of the primary endpoint. For this reason the term 'responder' (and 'non-responder') is used to  
453 express the clinical benefit of the treatment in terms of effects seen in individual patients. There may  
454 be a number of ways to define a 'responder'/'non-responder'. The definitions should be pre-specified in  
455 the protocol and should be clinically convincing. In clinical regulatory guidelines, it is stated that the  
456 'responder' analysis should be used in establishing the clinical relevance of the observed effect as an  
457 aid to assess efficacy and clinical safety. It should be noted that in instances there is some loss of  
458 information (and hence loss of statistical power) connected with breaking down the information  
459 contained in the original variables into 'responder' and 'non-responder'.

460 In some situations, the 'responder' criterion may be the primary endpoint (e.g. CHMP guideline on  
461 clinical investigation of medicinal products in the treatment of Parkinson's disease  
462 (EMA/CHMP/330418/2012 rev. 2)). In this case it should be used to provide the main test of the null  
463 hypothesis. However, the situation that is primarily addressed here is when the 'responder' analysis is  
464 used to allow a judgement on clinical relevance, once a statistically significant treatment effect on the  
465 mean level of the primary variable(s) has been established. In this case, the results of the 'responder'  
466 analysis need not be statistically significant but the difference in the proportions of responders should  
467 support a statement that the investigated treatment induces clinically relevant effects.

468 It should be noted that a 'responder' analysis cannot rescue the negative results on the primary  
469 endpoint(s).

## 470 **9. How should composite endpoints be handled statistically** 471 **with respect to regulatory claims?**

472 *Usually, the composite endpoint is primary. All components should be analysed separately. If claims*  
473 *are based on subgroups of components, this needs to be pre-specified and embedded in a valid*  
474 *confirmatory analysis strategy. In the event that treatment does not beneficially affect all components,*  
475 *in particular where the clinically more important components are affected negatively, interpretation will*  
476 *be very difficult. Any effect of the treatment in one of the components that is proposed to be reflected*  
477 *in the product information should be clearly supported by the data.*

478 There are two types of composite endpoints. The first type, namely the rating scale, arises as a  
479 combination of multiple clinical measurements. With this type there is a longstanding experience  
480 and/or validation of its use in certain indications (*e.g.* psychiatric or neurological disorders). This type  
481 of composite variable is not discussed further in this guideline.

482 The other type of a composite variable arises in the context of survival analysis. Several events are  
483 combined to define a composite outcome. A patient is said to have the clinical outcome if s/he suffers  
484 from one or more events in a pre-specified list of components (*e.g.* death, myocardial infarction or  
485 disabling stroke). The time to outcome is measured as the time from randomisation of the patient to  
486 the first occurrence of any of the events in the list. Usually, the components represent relatively rare  
487 events, and to study each component separately would require unmanageably large sample sizes.  
488 Composite endpoints therefore often present a means to increase the percentage of patients that reach  
489 the clinical outcome, and hence increase the power of the study.

### 490 **9.1. The composite endpoint as the primary endpoint**

491 When a composite endpoint is used to show efficacy it will often be the primary endpoint. In this case,  
492 it must meet the requirements for a single primary endpoint, namely that it is capable of providing the  
493 key evidence of efficacy that is needed for a licence. It is recommended to analyse in addition the  
494 single components and clinically relevant groups of components separately, to provide supportive  
495 information. There is, however, no need for an adjustment for multiplicity provided significance of the  
496 primary endpoint is achieved. If claims are to be based on (subgroups of) components, this needs to  
497 be pre-specified and embedded in a valid confirmatory analysis strategy.

### 498 **9.2. Treatment should be expected to affect all components in a similar** 499 **way**

500 A composite endpoint must make sense from a clinical perspective. For any component that is included  
501 in the composite, it is usually appropriate that any additional component reflecting a worse clinical  
502 event is also included. For example, if it is agreed that hospitalisation is an acceptable component in a  
503 composite endpoint, it would be usual to also include components for more adverse clinical outcomes  
504 that are relevant to the clinical setting (*e.g.* non-fatal myocardial infarction and stroke) and death.  
505 Excluding such events, with an argument that no beneficial effect can be expected or that these will be  
506 captured in the safety assessment, or focussing on specific types of events (for example disease-  
507 related mortality in preference to all-cause mortality) introduces difficulties for analysis and  
508 interpretation that should be approached carefully. In this event, the primary composite should always  
509 be presented and interpreted alongside a secondary analysis in which no important clinical outcomes  
510 are excluded.

511 In the event that treatment does not beneficially affect all components of a composite endpoint, in  
512 particular where the clinically more important components are affected negatively, interpretation will  
513 be complicated and the choice of composite as the primary variable should be carefully considered. An  
514 assumption of similarly directed treatment effects on all components should be based on past  
515 experience with studies of similar type. Whilst it may often be reasonable, *a priori*, to assume that no  
516 component of a composite relating to efficacy will be adversely affected, 'net clinical benefit' endpoints  
517 are employed to investigate whether beneficial effects are offset by increased detrimental effects.  
518 Because of the assumptions made in 'weighting' the components and in the overall interpretation, such  
519 composites will not usually be appropriate primary endpoints.

520 Composite endpoints also pose particular issues in the non-inferiority or equivalence setting, and  
521 analogously in relation to assessment of safety. Adding a component that foreseeably is insensitive to  
522 treatment effects tends to decrease sensitivity of the comparison, even if it does not affect

523 unbiasedness of the estimation of the treatment difference. An increased variance is an undesirable  
524 property in non-inferiority or equivalence studies. For non-inferiority or equivalence studies the more  
525 specific component (e.g. disease related mortality) can be preferred as primary endpoint for this  
526 reason, though again both this and the more general composite including all relevant events should be  
527 considered together.

### 528 **9.3. The clinically more important components should at least not be** 529 **affected negatively**

530 If time to hospitalisation is an endpoint in a clinical study it is not generally appropriate to handle  
531 patients who die before they reach the hospital as censored. It is better practice to study a composite  
532 endpoint that includes all important clinical events as components, including death in this example.  
533 One concern with composite outcome measures from a regulatory point of view is, however, the  
534 possibility that some of the treatments under study may have an adverse effect on one or more of the  
535 components, and that this adverse effect is masked by the composite outcome, e.g. by a large  
536 beneficial effect on some of the remaining components. This concern is particularly relevant if the  
537 components relate to different degrees of disease severity or clinical importance. For example, if all-  
538 cause mortality is a component, a separate analysis of all-cause mortality should be provided to ensure  
539 that there is no adverse effect on this endpoint. Since there is no general agreement on how much  
540 evidence is needed to generate suspicion of an adverse effect, it is recommended that this issue is  
541 addressed at the planning stage. For example, the study plan could address the size of the risk of an  
542 adverse effect on the more serious components that can be excluded (assuming no treatment  
543 difference under the null hypothesis) with a sufficiently high probability given the planned sample size,  
544 and the study report should contain the respective comparative estimates and confidence intervals.

545 Non-inferiority studies will also be particularly hard to interpret if negative effects on some components  
546 are observed for the experimental drug and are outbalanced by other components of the composite.

### 547 **9.4. Any effect of the treatment on one of the components that is intended** 548 **to be reflected in the product information should be clearly supported by** 549 **the data**

550 An important issue for consideration is the claim that can legitimately be made based on a successful  
551 primary analysis of a composite endpoint. Difficulties arise if the claims do not properly reflect the fact  
552 that a composite endpoint was used, e.g. if a claim is made that explicitly involves a component with  
553 the lowest frequency amongst all components. For example, if the composite outcome is death or liver  
554 transplantation and there are only a few deaths, a claim to reduce mortality and the necessity for liver  
555 transplantation would not be satisfactory, because in this context the effect on mortality will have a  
556 weak basis. This does not mean that one should drop the component death from the composite  
557 outcome, because the outcome liver transplantation would be incomplete without simultaneously  
558 considering all disease-related outcomes that are at least as serious as liver transplantation. However,  
559 it does mean that different wording should be adopted in the product information, avoiding the  
560 implication of a demonstrated effect on mortality.

## 561 **10. Multiplicity issues in estimation**

562 Often, for the more complex procedures, clinical interpretation of the findings can become difficult. For  
563 the purpose of estimation and for the appraisal of the precision of estimates, confidence intervals are  
564 of paramount importance. Multiple confidence intervals with an adjusted confidence level or  
565 multidimensional confidence regions (covering more than one unknown parameter with a given  
566 probability for the simultaneous assessment of multiple parameters) are typically used for multiple

567 comparisons but methods for their construction that are consistent with the tests are not available or  
568 not useful for many of the complex multiple testing procedures used to control the type I error.  
569 Nevertheless, a valid statistical procedure is useful only if it allows for a meaningful and informative  
570 clinical interpretation. Confidence regions, e.g. that are uninformative in the sense that they never  
571 exclude the null hypothesis of no treatment effect in order to comply with the multiple testing  
572 procedure, would have no relevance in the assessment of the trial results.

### 573 **10.1. Selection bias**

574 Multiple comparisons may lead to a bias in estimation which is defined by the difference between the  
575 mean estimation and the parameter to be estimated. For example, in a situation where several  
576 treatment groups are compared to placebo the strategy that chooses the treatment with the largest  
577 difference to placebo as the treatment that should be marketed will, on average, lead to an  
578 overestimation of the corresponding treatment effect. If selection is made not on the basis of the  
579 treatment effect it may still be based on an endpoint that is correlated with efficacy.

580 Whereas the term selection bias often relates to the bias resulting from a specific patient or subgroup  
581 selection, selection bias in the context of multiple comparisons refers to a biased estimation resulting  
582 from selecting a specific treatment (e.g. a specific dosage) based on the data that are subsequently  
583 used for estimation.

584 Selection bias is usually lower (but still present) if the selection is performed at an interim analysis.  
585 Selection at an earlier interim analysis leads to a lower bias, although it is less informative. However,  
586 methods are available to reduce selection bias, such as shrinkage estimation or specific model based  
587 analyses. Maximum bias should be gauged in order to account for it in the risk benefit assessment.

### 588 **10.2. Confidence intervals**

589 As can occur with multiple testing, multiple confidence intervals may also increase the chance of false  
590 decisions since the probability that a set of multiple non-adjusted confidence intervals cover correctly  
591 all parameters to be estimated would usually be less than the pre-specified nominal coverage  
592 probability related to the single confidence intervals.

593 Informative confidence regions that correspond to multiplicity procedures may, however, not always be  
594 available or may be difficult to derive. If the confidence regions do not correspond to the hypothesis  
595 testing procedure, different conclusions are possible, e.g. a confidence interval excluding the null  
596 hypothesis combined with a non-significant testing result or *vice versa*. The decision should, however,  
597 be based on the hypothesis test. In that case it is advised to use simple but conservative confidence  
598 interval methods, such as Bonferroni-corrected intervals, ensuring that the uncertainty about the  
599 beneficial effects is properly understood.