

---

# **Patient-Focused Drug Development: Incorporating Clinical Outcome Assessments Into Endpoints For Regulatory Decision-Making**

## **Guidance for Industry, Food and Drug Administration Staff, and Other Stakeholders**

### ***DRAFT GUIDANCE***

**This guidance document is being distributed for comment purposes only.**

Comments and suggestions regarding this draft document should be submitted within \_\_\_ days of publication in the *Federal Register* of the notice announcing the availability of the draft guidance. Submit electronic comments to <https://www.regulations.gov>. Submit written comments to the Dockets Management Staff (HFA-305), Food and Drug Administration, 5630 Fishers Lane, Rm. 1061, Rockville, MD 20852. All comments should be identified with the docket number listed in the notice of availability that publishes in the *Federal Register*.

For questions regarding this draft document, contact (CDER) Office of Communications, Division of Drug Information at [druginfo@fda.hhs.gov](mailto:druginfo@fda.hhs.gov), 855-543-3784, or 301-796-3400 or (CBER) Office of Communication, Outreach and Development, 800-835-4709 or 240-402-8010 or (CDRH) Patient Science and Engagement Program at [CDRH-PRO@fda.hhs.gov](mailto:CDRH-PRO@fda.hhs.gov), 301-796-6715.

**U.S. Department of Health and Human Services  
Food and Drug Administration  
Center for Drug Evaluation and Research (CDER)  
Center for Biologics Evaluation and Research (CBER)  
Center for Devices and Radiological Health (CDRH)**

**April 2023  
Procedural**

---

# **Patient-Focused Drug Development: Incorporating Clinical Outcome Assessments Into Endpoints For Regulatory Decision-Making**

## **Guidance for Industry, Food and Drug Administration Staff, and Other Stakeholders**

*Additional copies are available from:  
Office of Communications, Division of Drug Information  
Center for Drug Evaluation and Research  
Food and Drug Administration  
10001 New Hampshire Ave., Hillandale Bldg., 4<sup>th</sup> Floor  
Silver Spring, MD 20993-0002  
Phone: 855-543-3784 or 301-796-3400; Fax: 301-431-6353  
Email: [druginfo@fda.hhs.gov](mailto:druginfo@fda.hhs.gov)*

<https://www.fda.gov/drugs/guidance-compliance-regulatory-information/guidances-drugs>  
and/or

*Office of Communication, Outreach and Development  
Center for Biologics Evaluation and Research  
Food and Drug Administration  
10903 New Hampshire Ave., Bldg. 71, Room 3128  
Silver Spring, MD 20993-0002  
Phone: 800-835-4709 or 240-402-8010  
Email: [ocod@fda.hhs.gov](mailto:ocod@fda.hhs.gov)*

<https://www.fda.gov/vaccines-blood-biologics/guidance-compliance-regulatory-information-biologics/biologics-guidances>  
and/or

*Office of Policy  
Center for Devices and Radiological Health  
Food and Drug Administration  
10903 New Hampshire Ave., Bldg. 66, Room 5431  
Silver Spring, MD 20993-0002  
Email: [CDRH-Guidance@fda.hhs.gov](mailto:CDRH-Guidance@fda.hhs.gov)*

<https://www.fda.gov/medical-devices/device-advice-comprehensive-regulatory-assistance/guidance-documents-medical-device-and-radiation-emitting-products>

**U.S. Department of Health and Human Services  
Food and Drug Administration  
Center for Drug Evaluation and Research (CDER)  
Center for Biologics Evaluation and Research (CBER)  
Center for Devices and Radiological Health (CDRH)**

**April 2023  
Procedural**

**Contains Nonbinding Recommendations**

*Draft — Not for Implementation*

**TABLE OF CONTENTS**

<b>I.</b>	<b>INTRODUCTION.....</b>	<b>1</b>
<b>A.</b>	<b>Overview of the Series of FDA Guidance Documents on Patient-Focused Drug Development... 1</b>	
<b>B.</b>	<b>Purpose and Scope of PFDD Guidance 4.....</b>	<b>3</b>
<b>II.</b>	<b>COA-BASED ENDPOINT CONSIDERATIONS.....</b>	<b>4</b>
<b>A.</b>	<b>Endpoint of Interest: What Are You Measuring in the Target Study Population?.....</b>	<b>4</b>
1.	<i>Selecting and Justifying Endpoints.....</i>	<i>4</i>
2.	<i>Considerations for Constructing a COA-Based Endpoint.....</i>	<i>6</i>
3.	<i>Clinical Trial Duration and Timing of Assessments for COA-Based Endpoints.....</i>	<i>14</i>
<b>B.</b>	<b>Estimation and Missing Data.....</b>	<b>16</b>
1.	<i>Analysis at a Fixed Time Point.....</i>	<i>16</i>
2.	<i>Analyzing Ordinal Data.....</i>	<i>16</i>
3.	<i>Missing Data.....</i>	<i>17</i>
<b>III.</b>	<b>EVALUATING THE MEANINGFULNESS OF TREATMENT BENEFIT.....</b>	<b>18</b>
<b>A.</b>	<b>Factors Affecting the Interpretability of COA Scores.....</b>	<b>18</b>
1.	<i>How Closely Does the Measured Concept of Interest Correspond to the Patients' Experiences? 19</i>	
2.	<i>How Simple or Familiar is the COA's Metric?.....</i>	<i>19</i>
<b>B.</b>	<b>Approaches for Collecting Evidence to Support Interpretability of COA-Based Endpoints20</b>	
1.	<i>Interpreting in Terms of Meaningful Score Differences.....</i>	<i>20</i>
2.	<i>Interpreting in Terms of Meaningful Score Regions.....</i>	<i>24</i>
3.	<i>Additional Considerations for Justifying Meaningful Differences or Meaningful Score Regions.26</i>	
<b>C.</b>	<b>Applying Information About Meaningful Score Differences or Meaningful Score Regions to Clinical Trial Data.....</b>	<b>27</b>
1.	<i>Interpreting the Meaningfulness of Continuous COA-Based Endpoints.....</i>	<i>28</i>
2.	<i>Interpreting the Meaningfulness of Ordinal and Dichotomous COA-Based Endpoints.....</i>	<i>33</i>
<b>IV.</b>	<b>ADDITIONAL CONSIDERATIONS.....</b>	<b>33</b>
<b>A.</b>	<b>Other Study Design Considerations.....</b>	<b>33</b>
1.	<i>Masking.....</i>	<i>33</i>
2.	<i>Practice Effects.....</i>	<i>33</i>
3.	<i>Use of Assistive Devices.....</i>	<i>35</i>
4.	<i>Considerations When Using a Nonrandomized Design, External Controls, or Nonconcurrent Control.....</i>	<i>36</i>
5.	<i>Analysis of Treatment Effects for Subgroups Based on Post-Baseline Events.....</i>	<i>37</i>
6.	<i>Computerized Adaptive Testing.....</i>	<i>37</i>
7.	<i>Minimizing Participant Burden.....</i>	<i>38</i>
<b>B.</b>	<b>Formatting and Submission Considerations.....</b>	<b>38</b>
	<b>REFERENCES.....</b>	<b>41</b>

*Contains Nonbinding Recommendations*

*Draft — Not for Implementation*

**Patient-Focused Drug Development: Incorporating Clinical Outcome Assessments Into Endpoints for Regulatory Decision-Making Guidance for Industry, Food and Drug Administration Staff, and Other Stakeholders<sup>1</sup>**

This draft guidance, when finalized, will represent the current thinking of the Food and Drug Administration (FDA or Agency) on this topic. It does not establish any rights for any person and is not binding on FDA or the public. You can use an alternative approach if it satisfies the requirements of the applicable statutes and regulations. To discuss an alternative approach, contact the FDA staff responsible for this guidance as listed on the title page.

**I. INTRODUCTION**

**A. Overview of the Series of FDA Guidance Documents on Patient-Focused Drug Development**

This guidance (Guidance 4) is the fourth in a series of four methodological patient-focused drug development (PFDD) guidance documents<sup>2</sup> that describe how stakeholders (patients, caregivers, researchers, medical product developers, and others) can collect and submit patient experience data<sup>3</sup> and other relevant information from patients and caregivers to be used for medical product<sup>4</sup> development and regulatory decision-making. The topics that each guidance document addresses are described below:

<sup>1</sup> This guidance has been prepared by the Center for Drug Evaluation and Research in cooperation with the Center for Biologics Evaluation and Research and the Center for Devices and Radiological Health at the Food and Drug Administration.

<sup>2</sup> The four guidance documents fulfill commitments under section I.J.1 associated with the sixth authorization of the Prescription Drug User Fee Act (PDUFA VI) under Title I of the FDA Reauthorization Act of 2017, as well as requirements under section 3002 of the 21st Century Cures Act (available at <https://www.fda.gov/downloads/forindustry/userfees/prescriptiondruguserfee/ucm563618.pdf>).

<sup>3</sup> “Patient experience data” is defined for purposes of this guidance in Title III, Section 3001 of the 21st Century Cures Act, as amended by section 605 of the Food and Drug Administration Reauthorization Act (FDARA) of 2017, to include data that “(1) are collected by any persons (including patients, family members and caregivers of patients, patient advocacy organizations, disease research foundations, researchers and drug manufacturers); and (2) are intended to provide information about patients’ experiences with a disease or condition, including (A) the impact (including physical and psychosocial impacts) of such disease or condition or a related therapy or clinical investigation; and (B) patient preferences with respect to treatment of the disease or condition.”

<sup>4</sup> For purposes of this guidance a *medical product* refers to a drug (as defined in section 201 of the Federal Food, Drug, and Cosmetic Act (21 U.S.C. 321)) intended for human use, a device (as defined in such section 201) intended for human use, or a biological product (as defined in section 351 of the Public Health Service Act (42 U.S.C. 262)).

## Contains Nonbinding Recommendations

Draft — Not for Implementation

- 28 • Methods to collect patient experience data that are accurate and representative of the  
29 intended patient population (Guidance 1)<sup>5</sup>  
30
- 31 • Approaches to identifying what is most important to patients with respect to their  
32 experience as it relates to burden of disease/condition and burden of treatment  
33 (Guidance 2)<sup>6</sup>  
34
- 35 • Approaches to selecting, modifying, developing, and validating clinical outcome  
36 assessments (COAs) to measure outcomes of importance to patients in clinical trials  
37 (Guidance 3)<sup>7</sup>  
38
- 39 • Methods, standards, and technologies for collecting and analyzing COA data for  
40 regulatory decision-making, including selecting the COA-based endpoint and  
41 determining clinically meaningful change in that endpoint (Guidance 4; current  
42 guidance)  
43

44 Please refer to Guidance 1, Guidance 2, and other FDA guidances<sup>8</sup> for additional information on  
45 collecting patient experience data. When final, the PFDD guidance series will replace the  
46 guidance for industry *Patient-Reported Outcome Measures: Use in Medical Product  
47 Development to Support Labeling Claims* (December 2009).  
48

49 FDA encourages stakeholders to interact early with FDA and obtain feedback from the relevant  
50 FDA review division when considering the collection of patient experience data related to the  
51 burden of disease and the benefits, burdens, and harms of treatment.<sup>9</sup> FDA recommends that  
52 stakeholders engage with patients and other appropriate subject matter experts (e.g., clinical and  
53 disease experts, qualitative researchers, survey methodologists, statisticians, psychometricians,  
54 patient preference researchers) when designing and implementing studies to evaluate the burden  
55 of disease and treatment, and perspectives on treatment benefits and risks.

---

<sup>5</sup> See the FDA guidance for industry, FDA staff, and other stakeholders *Patient-Focused Drug Development: Collecting Comprehensive and Representative Input* (June 2020). We update guidances periodically. For the most recent version of a guidance, check the FDA guidance web page at <https://www.fda.gov/regulatory-information/search-fda-guidance-documents>.

<sup>6</sup> See FDA's guidance for industry, FDA staff, and other stakeholders *Patient-Focused Drug Development Methods to Identify What is Important to Patients* (February 2022).

<sup>7</sup> See the draft FDA guidance for industry, Food and Drug Administration staff, and other stakeholders *Patient-Focused Drug Development: Selecting, Developing, or Modifying Fit-for-Purpose Clinical Outcome Assessments* (June 2022). When final, this guidance will represent the FDA's current thinking on this topic.

<sup>8</sup> See FDA's guidance for industry *Patient Preference Information—Voluntary Submission, Review in Premarket Approval Applications, Humanitarian Device Exemption Applications, and De Novo Requests, and Inclusion in Decision Summaries and Device Labeling* (August 2016) and FDA's guidance for industry, Food and Drug Administration staff, and other stakeholders *Principles for Selecting, Developing, Modifying, and Adapting Patient-Reported Outcome Instruments for Use in Medical Device Evaluation* (January 2022), or subsequent guidances in the PFDD series, when available.

<sup>9</sup> In addition to the general considerations discussed in this guidance, a study may need to meet specific statutory and regulatory standards governing the collection, processing, retention, and submission of data to the FDA to support regulatory decisions regarding a marketed or investigational medical product. This guidance focuses on more general considerations that apply to many types of studies, and you should consult with the review division and applicable guidance regarding any other applicable requirements.

## *Contains Nonbinding Recommendations*

*Draft — Not for Implementation*

56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
68  
69  
70  
71  
72  
73  
74  
75  
76  
77  
78  
79  
80

### **B. Purpose and Scope of PFDD Guidance 4**

This guidance is intended to help sponsors of clinical trials for medical product development, as defined in footnote 4. This guidance focuses on COA issues associated with clinical trial (study) endpoints, design, conduct, and analysis and will be of most relevance for those designing and conducting trials using COAs as well as analyzing and interpreting the trial data.<sup>10</sup> This guidance builds on Guidance 3 by focusing on endpoints constructed from fit-for-purpose<sup>11,12</sup> COAs which are intended to reflect, directly or indirectly, how patients feel, function, or survive. Some COAs provide direct insight on how patients feel or function (e.g., a patient-reported outcome (PRO) instrument measuring pain intensity). Other COAs, however, may provide more indirect information to evaluate clinical benefit (e.g., clinician-reported outcome (ClinRO) instruments measuring extent or activity of disease such as psoriasis area and severity). In these situations, it is important to understand how the COA-based endpoint corresponds to changes relevant to patients (e.g., the type and extent of change that is meaningful to patients).

Section II of this guidance discusses considerations for COA-based endpoints to align the study design, endpoint, and analysis with the clinical study objective to improve study planning and the interpretation of analyses.

Section III of this guidance describes methods to aid in the interpretation of treatment effects on COA-based endpoints in terms of patients' views on the effect of a medical product. This information is important because statistical significance does not, by itself, indicate whether the detected effect corresponds to a clinically meaningful treatment effect.

---

<sup>10</sup> The considerations addressed in this guidance may be relevant to a variety of regulatory decisions that require an assessment of benefit or risk, including but not limited to: drug approval decisions under the standards in section 505(d) of the FD&C Act and regulations in 21 CFR 314; device approval decisions under the standards in sections 513(a)(2) and 515(d) and regulations in 21 CFR part 814; biological product approval decisions under the standards in section 351(a) of the Public Health Service Act and regulations in 21 CFR 601; device classification decisions under the standards in sections 513(a)(2) and 513(f) and regulations in 21 CFR parts 807 and 860; investigational new drug and investigational device exemption applications under sections 21 CFR parts 312 and 812; REMS and PMR requirements under sections 505-1 and 505(o)(3) and device post-approval requirements under 21 CFR part 814 subpart E; labeling decisions under 21 CFR parts 201, 801, and 809. Necessarily, this guidance does not attempt to capture all of the regulatory standards that might apply to a sponsor's intended plan of study; sponsors should consult the relevant review division(s) as necessary to discuss their study plans and are responsible for satisfying applicable requirements.

<sup>11</sup> See the Agency's draft guidance for industry, Food and Drug Administration staff, and other stakeholders *Patient-Focused Drug Development: Selecting, Developing, or Modifying Fit-for-Purpose Clinical Outcome Assessments* (June 2022). When final, this guidance will represent the FDA's current thinking on this topic.

<sup>12</sup> A COA is considered fit-for-purpose when the level of validation is sufficient to support its context of use. Note that having a fit-for-purpose COA is necessary for a strong endpoint rationale, but it is not sufficient. For example, a COA that is considered fit-for-purpose for assessing symptom intensity might be used for an endpoint based on the average symptom intensity score across 7 days. However, if worst intensity were identified as the most relevant patient experience for improvement based on patient input and the product's mechanism of action, the rationale for using an endpoint of average symptom intensity would be very weak—despite being based on a fit-for-purpose COA.

## ***Contains Nonbinding Recommendations***

*Draft — Not for Implementation*

81 Section IV of this guidance includes a list of additional considerations when developing an  
82 endpoint from a COA and formatting and submitting patient experience data from a clinical  
83 study supporting medical product regulatory decision-making.  
84

85 Though the text and examples in this guidance focus mostly on treatment benefit (e.g.,  
86 improvement in disease-related symptoms or impaired functions), COAs also can be used to  
87 assess treatment harms including symptomatic adverse events and other burdens to the patient  
88 associated with the medical product under study. While many of the recommendations in this  
89 guidance will apply to the evaluation of treatment benefit or risk, additional considerations may  
90 be needed when using COAs to inform treatment risks.  
91

92 In general, FDA’s guidance documents do not establish legally enforceable responsibilities.  
93 Instead, guidances describe the Agency’s current thinking on a topic and should be viewed only  
94 as recommendations, unless specific regulatory or statutory requirements are cited. The use of  
95 the word *should* in Agency guidances means that something is suggested or recommended, but  
96 not required.  
97

98

99

## **II. COA-BASED ENDPOINT CONSIDERATIONS**

100

101 This section discusses considerations for selecting COA-based endpoints, including the  
102 development of a well-justified rationale for the endpoints and considerations for statistical  
103 analyses of COA-based endpoints in clinical trials.  
104

105

### **A. Endpoint of Interest: What Are You Measuring in the Target Study Population?**

106

107

108 PFDD Guidance 3 discusses the importance of a fit-for-purpose COA. PFDD Guidance 4  
109 complements PFDD Guidance 3 by focusing on the rationale for the proposed use of COA scores  
110 to construct endpoints that will support inferences about the effect of a medical product on how  
111 patients feel or function. As with the rationale for interpreting COA scores as measures of the  
112 concept of interest, the rationale for the use of COA scores as the basis for an endpoint should be  
113 well-supported by evidence.  
114

114

115

#### *1. Selecting and Justifying Endpoints*

116

117 Generally, endpoints that are based on COAs should (1) reflect an aspect of the patient’s health  
118 that is meaningful; and (2) be capable of supporting an inference of treatment effect within the  
119 context of the planned clinical trial. For a given COA score, there may be multiple options for  
120 constructing a trial endpoint (e.g., mean score at 12 weeks or time to complete symptom  
121 resolution).  
122

122

123 Sponsors should clearly describe the COA-based endpoint, including:  
124

124

- 125 • Type of assessment(s) made (e.g., Patient-Reported Outcome (PRO) measures, Observer-  
126 Reported Outcome (ObsRO) measures, Clinician-Reported Outcome (ClinRO) measures,

## *Contains Nonbinding Recommendations*

*Draft — Not for Implementation*

127 Performance Outcome (PerfO) measures).

128

- 129 • The COA(s) used to measure the concept(s) of interest. Note that it is important for
- 130 endpoints to be assessed using a COA that is fit-for-purpose. For details, see draft PFDD
- 131 Guidance 3.

132

- 133 • Specific score(s) from the COA (e.g., specific subscale score, total score).

134

- 135 • If a multi-component endpoint, the algorithm used to combine scores from two or more
- 136 components into a single endpoint.

137

- 138 • Rules for handling missing item responses or task results when computing COA scores,
- 139 along with justification for the rules.

140

- 141 • Timing of the assessments used to construct the endpoint, the timeframe over which COA
- 142 scores are combined to construct the endpoint, and a detailed description of how COA
- 143 scores collected during the treatment period are combined into an endpoint (e.g., score at
- 144 week 12, average daily scores for 7 days prior to week 12 study visit, maximum value of
- 145 the daily 200 mobile sensor assessments for 7 days prior to the week 36 study visit.).
- 146 Also, if the endpoint is defined in terms of change from baseline to some follow-up
- 147 assessment, then the definition of “baseline” should be clear.

148

149 FDA recognizes that constructing and selecting trial endpoint(s) often involves weighing the

150 strengths and limitations of different approaches. Early in the planning of a clinical trial,

151 sponsors should provide to FDA a well-supported rationale for the selection of the endpoint(s)

152 by explaining why each endpoint is informative for the trial context. The rationale for endpoint

153 selection typically will address the following:

154

- 155 • Concept(s) of interest.

156

- 157 • Clinical trial objective or hypothesis corresponding to the endpoint, ensuring that the
- 158 objective/hypothesis is specific (e.g., “To compare the patient-reported physical
- 159 functioning between arms at 24 weeks” rather than “To compare the patient-reported
- 160 outcomes of product X vs. Y”).

161

- 162 • The role of the endpoint (e.g., primary, secondary, other).

163

- 164 • Intended indication related to the COA-based endpoint.

165

- 166 • Explanation for why the selected COA is fit-for-purpose in the planned trial.

167

- 168 • Support for the importance of the endpoint to patients and/or caregivers from literature



## *Contains Nonbinding Recommendations*

*Draft — Not for Implementation*

169 review and/or primary data collection.<sup>13</sup> In some cases, for endpoints based on a COA  
170 that measures a concept of interest that is indirectly related to some meaningful aspect of  
171 health for the patient (e.g., based on a neurological functioning test that is thought to be  
172 indicative of the patients' cognitive functioning), it might be sufficient to provide support  
173 for the adequacy of the endpoint for measuring this aspect of health. Furthermore, there  
174 are well-established relevant outcomes such as organ failure and death that do not require  
175 additional support. If a multi-component endpoint, justification for the components  
176 included and the algorithm for combining them into the endpoint.

177

178 • Strengths and limitations of the proposed endpoint.

179 An endpoint's use in another trial evaluating a different product may not be adequate support for  
180 the use of the same endpoint for a trial under consideration, because the context of use can vary  
181 in important ways from trial to trial and science and/or policy might have evolved since the  
182 endpoint was last used. When disease-specific FDA guidances exist, sponsors should consult  
183 these for recommendations for suitable endpoints.<sup>14</sup>

184

### 185 2. *Considerations for Constructing a COA-Based Endpoint*

186

187 This section provides guidance on using scores from one or more COAs to construct endpoints  
188 for specific circumstances as well as guidance regarding particular types of endpoints. This is  
189 not a comprehensive review of all possible types of endpoints but rather a discussion of  
190 frequently encountered challenges for COA-based endpoints.

191

#### 192 a. Considerations for baseline administration of COAs relevant to COA- 193 based endpoints

194

195 Prior to discussing the different approaches, several considerations about collecting COAs at  
196 baseline should be noted:

197

198 • Some diseases, conditions, or clinical trial designs may necessitate more than one  
199 baseline assessment or longer/shorter baseline periods.

200

201 • When multiple baseline measurements are taken, the protocol should define how the  
202 baseline value will be calculated from the multiple measurements.

203

204 • A screening visit that includes administration of the COA is often used to ensure that  
205 patients enrolled in the trial have a sufficient level of severity so that improvement could

---

<sup>13</sup> For example, Stone et al. (2021) conducted semi-structured interviews with patients who have chronic pain (as well as clinicians and clinical trial lists) to elicit their understanding of and preferences for seven different endpoints that could be constructed based on intensive longitudinal assessments of pain intensity (e.g., average pain over a week, worst pain intensity over a week, time spent with low or no pain). Patients were asked to rank the different endpoints in the order of what they were "most hoping for as a result of treatment."

<sup>14</sup> Please see the FDA guidance web page <https://www.fda.gov/regulatory-information/search-fda-guidance-documents>

## *Contains Nonbinding Recommendations*

*Draft — Not for Implementation*

206 be observed. To avoid regression to the mean and other potential sources of bias,<sup>15</sup> the  
207 COA score obtained at screening should not be used as the patient’s baseline value.  
208 Rather, a separate, later pre-randomization assessment should be used as the patient’s  
209 baseline value.

210  
211 • If the trial includes a run-in period during which the patient’s score from the COA might  
212 be expected to change (e.g., medication washout, patient behavior modification), then this  
213 should be considered when planning the timing of assessments.

214  
215 b. Endpoints based on COA scores at a fixed time point or a summary of  
216 COA scores over time

217  
218 In most situations in which a COA produces ordinal or continuous (interval or ratio scale) scores,  
219 the best and recommended endpoint will be the COA score at a predefined assessment point or  
220 summarized over some predefined post-baseline assessment period, and the most straightforward  
221 analysis will be a comparison of randomized groups with respect to the follow-up score(s) after  
222 adjusting for the baseline value (e.g., with a linear model to compare average follow-up scores).

223  
224 When the endpoint is based on COA scores at a predefined assessment point, sponsors should  
225 justify the use of, and time at which, an analysis at a fixed time point (e.g., 12 weeks) is to be  
226 performed. For example, an analysis at a fixed time point might be justified if the COA score is  
227 not highly variable over time and the chosen time point (e.g., end of study) would be useful for  
228 reflecting the durability of the treatment effect. Justification of the fixed time point should also  
229 take the recall period of the COA (where applicable) into consideration.

230  
231 When considering an endpoint based on summarizing COA scores over some predefined post-  
232 baseline assessment period, different summaries may be appropriate depending on the research  
233 questions. Common types of summaries include the patient’s mean score over a fixed time  
234 period, the maximum (or minimum) score during some period (e.g., worst pain recorded during a  
235 7-day period). For some types of summaries, an alternative approach is to use repeated measures  
236 modeling of all observed COA scores and derive summary estimates from the model.

237 Regardless of the approach taken, sponsors who wish to construct an endpoint based on  
238 summaries of patients’ COA scores over time should consider the robustness of the summary (or  
239 model) and any modeling assumptions, handling of missing COA scores, statistical power, and  
240 interpretability.

241  
242 c. Endpoints constructed by dichotomizing COA scores

243  
244 COA scores are often ordinal or continuous (interval or ratio scales) in nature. When this is the  
245 case, defining the endpoint using the ordinal or continuous COA score, rather than making the  
246 endpoint dichotomous, uses all the information and therefore usually maximizes statistical  
247 power. In some cases, dichotomized endpoints (e.g., “responder” status) are well-established  
248 and can be reasonable choices when it is important to evaluate the effect of treatment on the

---

<sup>15</sup> Shaw PA, Johnson LL, Proschan MA, 2018, Intermediate Topics in Biostatistics, In JI Gallin, FP Ognibene, LL Johnson (Eds) *Principles and Practice of Clinical Research* (4<sup>th</sup> ed), London: Academic Press, pp. 384-409.

## Contains Nonbinding Recommendations

Draft — Not for Implementation

249 probability of achieving clearly defined and important health states. Examples of such health  
250 states might be complete patient-reported symptom resolution or investigator’s global assessment  
251 of acne lesions as “clear” or “almost clear” (see the May 2018 guidance for industry *Acne*  
252 *Vulgaris: Establishing Effectiveness of Drugs Intended for Treatment*). If a sponsor wishes to  
253 use an endpoint based on dichotomization from either ordinal or continuous data, the sponsor  
254 should prespecify a single score threshold and provide evidence to justify the dichotomization in  
255 the endpoint rationale. For example, FDA recommends that the rationale include evidence that  
256 patients and/or their caregivers view health states above the threshold to be meaningfully  
257 different from health states below the threshold. This recommendation also applies to the use of  
258 ordinal or continuous COA data to define an event for a time-to-event endpoint. Of note, data  
259 used to derive a score threshold(s) should be different than that used to demonstrate effectiveness  
260 (e.g., data from registration trial(s)). In addition to prespecifying a single score threshold,  
261 sponsors should also conduct analyses to explore treatment effects over a range of thresholds.

262  
263 Sometimes the motivation for dichotomizing an ordinal or continuous COA-based score is to  
264 make the endpoint more interpretable for patients, caregivers, and/or clinicians. This is typically  
265 possible without creating a dichotomized endpoint for the primary analysis of treatment effect.  
266 (See Section III, *Evaluating the Meaningfulness of Treatment Benefit*).

267  
268 d. Endpoints constructed by computing change from baseline or percent  
269 change from baseline COA scores

270  
271 As discussed in Section II.A.2.a, in comparative trials, the preferred method for adjusting for  
272 baseline status is to do so in the context of a statistical model. Using the COA score’s change-  
273 from-baseline as an endpoint is another option, but it has some important considerations:

- 274
- 275 • COA scores that are ordinal are challenging to interpret in terms of change from baseline  
276 because the difference between two ordinal scores cannot be assumed to have the same  
277 meaning across scores (e.g., for an ordinal score with 5 levels—when interpreting level 3  
278 relative to level 1 and level 5 relative to level 3—both differ by two levels but might not  
279 correspond to the same degree of change in the underlying health state). Put another  
280 way, there might not be a linear relationship between the ordinal values and the true level  
281 of symptom severity or functioning being measured.
  - 282 • If it aids interpretation to express treatment effects in terms of change-from-baseline, this  
283 can be done in the context of most models used to compare treatment groups on follow-  
284 up scores adjusting for baseline. For example, an ANCOVA model could be used to  
285 derive the predicted follow-up score on treatment for patients with a given baseline score,  
286 and these two values could be used to compute a predicted change-from-baseline score.
  - 287 • For situations in which it is not possible to conduct a randomized, controlled trial and a  
288 single arm trial is done instead (e.g., to evaluate some devices), a change-from-baseline  
289 endpoint might be the best available option.

290  
291  
292  
293 A similar endpoint that could be considered is the percent change-from-baseline. An advantage  
294 of this approach might be easier interpretability, but in addition to the considerations presented

## *Contains Nonbinding Recommendations*

*Draft — Not for Implementation*

295 for change-from-baseline endpoints, several important challenges are worth noting about percent  
296 change-from-baseline:

- 297
- 298 • Interpretation can be complicated by the fact that percent change-from-baseline is  
299 asymmetric; that is, it treats the baseline and follow-up COA scores differently (Berry  
300 and Ayers 2006). For example, consider two patients who are randomized to receive a  
301 new medical product. The first patient's COA score improves from 5 to 10 (change =  
302 +5) and the second patient's score decreases from 10 to 5 (change = -5). In both cases,  
303 the absolute change is 5, but the percent change is very different: +100% and -50%. This  
304 has important implications, including the fact that the average change on the original  
305 scale (0) indicates no overall change, whereas the average percent change ( $[+100 - 50]/2$   
306 = +25%) suggests an overall improvement.
  - 307
  - 308 • Percent change-from-baseline is undefined if the baseline score on a COA is zero, and  
309 some kind of imputation is required to include the observation in the analysis.
  - 310
  - 311 • Compared to follow-up scores or change-from-baseline scores, percent change-from-  
312 baseline scores may have highly non-normal distributions that can be challenging to  
313 model.
- 314

315 If the reason for considering percent change-from-baseline is that the treatment effect is expected  
316 to be multiplicative rather than additive (e.g., treatment improves a patient's symptom severity  
317 by 20% of the patient's severity level without treatment), then a logarithmic or similar  
318 transformation could be applied to continuously distributed COA scores prior to comparing  
319 groups (Senn 2007).

320

321 e. Endpoint strategies when a disease affects multiple aspects of feeling and  
322 functioning

323

324 A disease might manifest in multiple ways, in which case it is important to consider how or  
325 whether a medical product affects different aspects of health. Some aspects of health might be  
326 relevant for almost all patients with a given condition (e.g., pain associated with migraine).  
327 Other affected aspects of health might differ between patients and within patients over time with  
328 certain conditions (e.g., lupus, sarcoidosis, primary mitochondrial diseases, schizophrenia, and  
329 many rare diseases). In these situations, it may be challenging to identify one specific aspect of  
330 the disease for evaluating treatment benefit. It may be necessary to consider several different  
331 aspects to adequately assess benefit. FDA recognizes that selection of the endpoint(s) in these  
332 situations is likely to involve weighing the strengths and limitations of various approaches.  
333 When possible, sponsors can evaluate multiple endpoints in earlier phase trials to inform the  
334 selection of endpoints for later trials.

335

336 This section reviews three general strategies for constructing endpoints when multiple  
337 aspects of health might be of interest: (1) separate endpoints for each aspect of health, (2) a  
338 multi-component endpoint, and (3) a personalized endpoint.

339

340 *Construct Separate Endpoints for Each Aspect of Health*

## *Contains Nonbinding Recommendations*

*Draft — Not for Implementation*

- 341  
342 As described in the guidance for industry *Multiple Endpoints in Clinical Trials* (October  
343 2022), if a separate endpoint will be constructed for each aspect of health, their role should  
344 be described, with the main options as follows:  
345
- 346 • *One primary endpoint and multiple secondary endpoints.* This option might be useful  
347 when there is one core or cardinal manifestation of a disease (primary endpoint) that most  
348 patients can be expected to experience and that is regarded by patients and/or caregivers  
349 as important. Secondary endpoints can be created for aspects of health that might not be  
350 experienced by all patients and/or are viewed as relatively less critical, but still important,  
351 to patients and/or caregivers.  
352
  - 353 • *Multiple primary endpoints.* This option might be useful when an improvement in at  
354 least one aspect of health would be regarded as evidence of treatment benefit.  
355
  - 356 • *Co-primary endpoints.* This option may be appropriate when there are multiple aspects  
357 of health that are critically important to the disease being studied, such that a treatment  
358 benefit can only be concluded if the medical product has an effect on each of the  
359 designated endpoints.

360  
361 By creating a separate endpoint for each relevant aspect of health, there is clarity about which  
362 aspect of health has or has not been affected by the medical product, because each endpoint  
363 corresponds to only one aspect of health. But there are several issues with this approach that also  
364 should be considered. First, for diseases with many possible manifestations, the approach may  
365 be challenging to use if it is not known ahead of time which aspects of health are most likely to  
366 improve as a result of using the medical product under study. Second, depending upon the roles  
367 of the multiple endpoints, multiplicity adjustments might be needed, necessitating a larger  
368 sample size to ensure sufficient statistical power. Finally, if patients differ from one another in  
369 their symptoms or functional impacts due to the disease, then the treatment effect estimated for  
370 any one endpoint will be diluted by the patients for whom the endpoint is not relevant (e.g.,  
371 patients who never had a given symptom cannot improve with treatment). Consult the guidance  
372 for industry *Multiple Endpoints in Clinical Trials* (October 2022) for additional information on  
373 constructing and analyzing multiple endpoints in a single trial.  
374

### *Construct a Multi-Component Endpoint*

375  
376  
377 A multi-component endpoint is based on a within-patient combination of two or more  
378 components, each reflecting a different aspect of health. Constructing the endpoint for an  
379 individual patient requires observation of all the specified components for that patient. Then  
380 a single overall rating or status on the endpoint is determined according to a prespecified  
381 algorithm.  
382

## *Contains Nonbinding Recommendations*

*Draft — Not for Implementation*

383 A COA-based multi-component endpoint may take many forms. The individual components  
384 could be (a) scores from different COAs, (b) scores from multiple subscales of a single COA,  
385 or (c) responses to individual items or tasks that make up a single COA.<sup>16</sup>  
386

387 Some COA-based multi-component endpoints are constructed by combining the patient’s  
388 scores—in their original metric or transformed (e.g., dichotomized)—from two or more  
389 components according to an algorithm. Some examples include:  
390

- 391 • An overall symptom index score created by using a well-justified weighted combination  
392 of responses to separate items that each assess a different type of symptom.  
393
- 394 • Patients’ endpoint values (“improved” versus “not improved”) are assigned based on a  
395 more complex algorithm, for example, an algorithm requiring some minimum change-  
396 from-baseline for one COA and some minimum change on at least two of four other  
397 COAs.  
398

399 Other multi-component endpoints are constructed with the objective of demonstrating the  
400 absence of all symptoms. Examples include:  
401

- 402 • Achievement of complete resolution of all symptoms
- 403 • Total time without any symptoms during some predefined post-baseline period
- 404 • Time until complete resolution of all symptoms
- 405 • Time to sustained clinical recovery assessed over an appropriate duration  
406

407 There are several advantages to using a multi-component endpoint, including:  
408

- 409 • A multi-component endpoint has the potential to evaluate the entire range of important  
410 disease manifestations. Because patients may experience some aspects of a disease more  
411 than others—and some aspects, not at all—a multi-component endpoint lends itself to  
412 capturing a treatment effect more so than an endpoint that evaluates a narrower aspect of  
413 the disease.  
414
- 415 • No multiplicity adjustment is needed to control the chance of erroneous conclusions (e.g.,  
416 Type 1 error) for a multi-component endpoint compared to the use of multiple separate  
417 endpoints.  
418
- 419 • The use of within-patient multi-component endpoints can be efficient if the treatment  
420 effects on the different components are generally concordant.  
421

---

<sup>16</sup> Responses to individual items or tasks that make up a single COA could be treated as individual components of a multicomponent endpoint only when the COA is based on a composite indicator measurement model. In a composite measurement model, responses to the items or tasks are not assumed to be reflective of or caused by a single underlying aspect of health (as they would be for a reflective measurement model). Instead, each item or task addresses a separate health concept and, when combined, responses to all the items or tasks define the overall concept of interest. See the draft guidance for industry *Patient-Focused Drug Development: Selecting, Developing, or Modifying Fit-for-Purpose Clinical Outcome Assessments* (June 2022), Section IV.E.

## *Contains Nonbinding Recommendations*

*Draft — Not for Implementation*

422 These advantages should be weighed against important concerns and limitations with  
423 constructing certain types of multi-component endpoints, including:

- 424
- 425 • For endpoints that are based on complete resolution of all symptoms, it might be difficult  
426 to achieve complete resolution with a medical product in the context of a clinical trial.  
427 Furthermore, some patient populations might not require complete resolution of all  
428 symptoms to feel they have benefitted from treatment. Other endpoints may be advisable  
429 to assess treatment-related improvements in individual symptom intensity or frequency.  
430
  - 431 • For multi-component endpoints that sum or average over scores from multiple  
432 components, a clinically meaningful improvement in one COA becomes increasingly  
433 diluted as more COAs are included in the construction of the endpoint. For example, if a  
434 patient’s only manifestation of a disease is symptom *A*, then the patient might appear to  
435 show little improvement if the multi-component endpoint averages the status on symptom  
436 *A* with symptoms *B*, *C*, *D*, and *E*. Therefore, sponsors considering this type of multi-  
437 component endpoint should balance the ability to observe improvements in any of several  
438 aspects of health with the chance that improvements in one aspect will be diluted by  
439 aspects that were never a problem for the patient. Sponsors might also consider the use  
440 of a personalized endpoint in such situations (see *Construct a Personalized Endpoint*  
441 below).  
442
  - 443 • All multi-component endpoints are based on some implicit or explicit weighting scheme.  
444 This includes multi-component endpoints that imply that all components have reasonably  
445 similar clinical importance, such as when taking the average across multiple COAs or  
446 assigning the status of “improved” to a patient who shows improvement in scores for any  
447 1 of 5 COAs. Sponsors should be explicit about how each component is weighted in  
448 constructing the endpoint and provide justification for the weights.  
449
  - 450 • When a treatment effect is found using a multi-component endpoint, it may be helpful to  
451 examine the treatment effect for individual components. For more detail about when and  
452 how to examine individual components, see the guidance for industry *Multiple Endpoints*  
453 *in Clinical Trials* (October 2022).  
454
  - 455 • There are several challenges for endpoints that rely on categorizing meaningful changes  
456 in one or more COAs.  
457
    - 458 – Endpoint values are strongly dependent on the thresholds selected for meaningful  
459 improvement and/or worsening and choosing such thresholds can be challenging.  
460 Thresholds for each COA should be predefined and justified. Sponsors should also  
461 conduct sensitivity analyses that explore treatment effects over a range of thresholds.  
462
    - 463 – There is the potential for bias when those completing or administering the COA are  
464 aware of the thresholds for being considered a meaningful improvement (or  
465 worsening). It is important when possible that clinicians (for ClinRO measures),  
466 caregivers (for ObsRO measures), and/or any research staff (for PerfO measures)

## Contains Nonbinding Recommendations

Draft — Not for Implementation

467 involved in assessment are not made aware of the threshold definitions and are  
468 masked<sup>17</sup> to treatment assignment.

469  
470 - Endpoints that assign values of worsened = -1, no change = 0, and improved = +1  
471 assumes that the patients view the degree of improvement and deterioration in a  
472 concept of interest as symmetric, which may not be the case.

473

### 474 *Construct a Personalized Endpoint*

475

476 Personalized endpoints are sometimes proposed to reflect what is important to each  
477 individual patient enrolled in a clinical trial, especially for diseases with variable clinical  
478 manifestations that impact patients differently. Several examples include the following:

479

480 • The “most bothersome symptom” approach in which patients identify at baseline the one  
481 disease-related symptom that is most bothersome to them. The patient’s status on that  
482 symptom post-randomization then becomes the outcome to be analyzed (Duke Margolis  
483 Center for Health Policy 2017). A similar approach is based on patients identifying at  
484 baseline the symptom that is “most severe” for them (which may or may not be the  
485 symptom that is most bothersome for them).

486

487 • Goal Attainment Scaling (GAS; Krasny-Pacini et al. 2016) in which each patient  
488 identifies a prespecified number of personal goals (e.g., being able to work in the garden)  
489 at baseline. At one or more post-randomization assessments, the patient records their  
490 status with respect to each goal using a standardized response scale and the responses are  
491 summarized across the patient’s goals. Whereas the “most bothersome” and “most  
492 severe” symptom approaches are based on assessments of symptoms, GAS usually is  
493 based on assessment of functioning.

494

495 Personalized endpoints have several advantages, including:

496

497 • They are very patient focused in their attempt to reflect how each patient feels or  
498 functions in terms of what is most important to them at baseline.

499

500 • Because each patient’s endpoint value is based only on what was identified as an issue  
501 for them at baseline, there is no dilution of treatment effect due to mixing affected and  
502 unaffected patients (i.e., when treating each aspect of health as its own endpoint) or  
503 mixing affected and unaffected aspects of health within a patient (i.e., when constructing  
504 some multi-component endpoints).

505

506 • Depending upon the context of use, a personalized endpoint could be considered along  
507 with another endpoint to inform decisions about the effect of a medical product. For  
508 example, the FDA guidance for industry *Migraine: Developing Drugs for Acute*

---

<sup>17</sup> Keeping study group assignment hidden from those involved in a study or trial is commonly referred to as “blinding” or “masking.” Those who do not know the assignment are referred to as “blinded” or “masked.” The term “masked” is used in this guidance.



## *Contains Nonbinding Recommendations*

*Draft — Not for Implementation*

509 *Treatment* (February 2018) describes using two co-primary endpoints: (1) having no  
510 headache pain at 2 hours after dosing; and (2) a demonstrated improvement on the  
511 patient’s most bothersome migraine-related symptom at 2 hours after dosing. (Note that  
512 this approach is specific to the context of use and might not be appropriate in other  
513 contexts of use.)  
514

515 These advantages should be weighed against several concerns, including:  
516

- 517 • For personalized endpoints that rely on patients choosing a single “most bothersome” or  
518 “most severe” symptom, it might be difficult for patients to select a single symptom.  
519
- 520 • Changes might occur over the duration of a clinical trial in what patients regard as their  
521 “most bothersome” symptom, “most severe” symptom, or their most important personal  
522 goals.  
523
- 524 • It is possible that patients might choose symptoms or areas of functioning (for GAS) at  
525 baseline that are not targeted by the product being evaluated or that might not be realistic  
526 to achieve for patients in the target population.  
527
- 528 • The outcomes chosen by patients might not reflect new or worsening symptoms and/or  
529 functional limitation(s) that occur during the trial duration. For this reason, the same set  
530 of outcome assessments should be assessed for all patients regardless of their own  
531 personalized endpoint.  
532
- 533 • The processes for eliciting personalized endpoints have the potential for inconsistency.  
534 Therefore, the process to construct a personalized endpoint should be standardized and  
535 the criteria for selecting the outcome assessments should be consistent across sites and  
536 patients.  
537
- 538 • As with multi-component endpoints, it is challenging to describe the specific effect of the  
539 treatment on a personalized endpoint. For this reason, it is important to measure all  
540 relevant symptoms and areas of functioning in addition to those identified as most  
541 important to the individual patients. This will make it possible to conduct prespecified  
542 treatment comparisons for individual symptoms and types of functioning.  
543

### 544 3. *Clinical Trial Duration and Timing of Assessments for COA-Based Endpoints* 545

546 Generally, COA data should be collected over the duration of the clinical trial, as indicated for  
547 other measures of effectiveness or safety in the clinical trial protocol.  
548

549 The timing of assessments plays a vital role in gaining reliable and meaningful information on  
550 the concept(s) of interest reflected in the COA-based endpoint and should be selected carefully  
551 and be scientifically justified. Clinical trials using COAs should include a schedule of COA  
552 administration as part of the overall study assessment schedule in the protocol. The COA  
553 schedule should consider the natural course of the disease or condition (i.e., acute, chronic, or  
554 episodic), the research questions to be addressed, the trial duration, patient burden, the disease

## *Contains Nonbinding Recommendations*

*Draft — Not for Implementation*

555 stage of the target patient population, the expected time frame when the investigational product  
556 is likely to affect the COA-based endpoint, and timing of collection of COAs if temporary study  
557 interruptions or discontinuation of study interventions are anticipated to occur.  
558

559 In general, COA assessment frequencies or the rules governing when the COA is measured  
560 should be the same for all treatment arms (see event-triggered data collection below). In many  
561 instances, such as when a COA is planned to be frequently measured (e.g., event-triggered data  
562 collection) or when the COA is complex and potentially burdensome, sponsors might consider  
563 seeking input from members of the patient community to ensure that the planned length of the  
564 trial and timing of COA assessments is feasible and as convenient as possible for the patients  
565 and/or caregivers. This input may help to reduce missed assessments and study dropout.  
566 Sponsors can further reduce patient burden by including only those assessments that are well  
567 justified within the context of the study objectives. See Section IV.A.7 (*Minimizing Participant*  
568 *Burden*) for more discussions.  
569

570 Other important considerations for determining the most appropriate timing of assessments for  
571 COA-based endpoints include, but are not limited to, the following:  
572

- 573 • *Event-triggered data collection:* In some studies, COA administration may be triggered  
574 to occur during or following events such as urination or an asthma exacerbation. For this  
575 type of data collection, consider the windows for data collection around an event and  
576 whether it would be appropriate to prompt to ensure that all events were collected (i.e., at  
577 the end of the diary day). For example, for a trial evaluating a treatment for a disorder  
578 that results in difficulty or excessive frequency of urination, a participant could be asked  
579 to record each urination episode and complete a short assessment immediately following  
580 the event (e.g., pain or burning during urination, post-micturition dribble). Then, at the  
581 end of the diary day, the patient could be shown a list of reported urination episodes and  
582 asked if they had any other urination episodes that needed to be reported and assessed.  
583
- 584 • *Anticipated rate of change in the underlying concept of interest to be measured:* The  
585 timing of assessments should align with the anticipated nature and rate of change in the  
586 underlying concept of interest to be measured. For example, if the concept of interest to  
587 be measured is expected to change rapidly over the course of the study period, then  
588 assessments should be placed closer together. If the concept of interest is expected to  
589 change slowly, then assessments can be placed further apart.  
590
- 591 • *Ability to assess time-to-event endpoints:* If the trial endpoint is based on time to achieve  
592 an outcome of interest (e.g., time to complete symptom resolution), the frequency of  
593 assessment should be sufficient to assess clinically meaningful differences in the time to  
594 the outcome of interest. If assessments are made too infrequently, important differences  
595 between trial arms may not be detected.  
596
- 597 • It will typically be of interest to understand treatment effects regardless of adherence to  
598 treatment, such that the protocol should include plans to continue to follow patients and  
599 administer the COA after discontinuation of treatment.  
600

## Contains Nonbinding Recommendations

Draft—Not for Implementation

### 601           **B.       Estimation and Missing Data**

602  
603       The statistical analysis considerations for COA-based endpoints are similar to the statistical  
604       considerations for any other endpoint used in medical product development. This section briefly  
605       discusses several considerations that commonly arise when estimating COA-based estimands,<sup>18</sup>  
606       including missing data.

#### 607 608           1.       *Analysis at a Fixed Time Point*

609  
610       For evaluating a treatment effect on COA scores at a fixed time point, the statistical power of the  
611       treatment group comparison is generally better when the comparison is statistically adjusted for  
612       patients' baseline scores<sup>19</sup> on the COA (see the draft guidance for industry *Adjusting for*  
613       *Covariates in Randomized Clinical Trials for Drugs and Biological Products* (May 2021)).<sup>20</sup>  
614       This recommendation also applies when the endpoint is the change in COA score from baseline  
615       to a predefined time point.

616  
617       If a COA-based endpoint is collected repeatedly, data from intermediate time points (i.e.,  
618       measurements taken prior to the fixed time point) can still be included in a longitudinal (e.g.,  
619       mixed-effects or generalized estimating equations) model in which a treatment contrast is made  
620       for a prespecified fixed time point.

#### 621 622           2.       *Analyzing Ordinal Data*

623  
624       Sometimes COA scores are used to construct an endpoint that results in an ordinal metric.  
625       Several analytic options exist for ordinally scaled endpoints. The choice of analytic approach  
626       might depend on the type of ordinal endpoint. For COA-based endpoints, there are generally  
627       two situations that generate an ordinal scale:

- 628  
629       •   *An ordinal endpoint based on a COA measuring a single aspect of health.* For example,  
630       a group comparison at a fixed time point might be made using a single item COA  
631       measuring the intensity of musculoskeletal pain might have response options of *none*,  
632       *mild*, *moderate*, and *severe*, which are scored as 0, 1, 2, and 3. The steps between  
633       successive levels might not reflect equal increments in pain, and so it might be  
634       challenging in some cases to interpret an estimate of treatment effect in terms of mean  
635       differences (e.g., as generated by an ANCOVA). On the other hand, an approach that  
636       tries to simplify the endpoint for analytic purposes by dichotomizing (e.g., [0 or 1] vs [2  
637       or 3]) risks ignoring important information about patients' relative standing on the

---

<sup>18</sup> An *estimand* is defined as a precise description of the treatment effect reflecting the clinical question posed by the trial objective. It summarizes at a population-level what the outcomes would be in the same patients under different treatment conditions being compared (see the ICH guidance for industry *E9(R1) Statistical Principles for Clinical Trials: Addendum: Estimands and Sensitivity Analysis in Clinical Trials* (May 2021) (ICH E9(R1)).

<sup>19</sup> Patient or clinician global impressions of severity, when used as anchor variables (see Section III), should be assessed at baseline. Note that patient or clinician global impressions of change used as anchor variables are not administered at baseline. Also, the concept of baseline or baseline symptoms may be complicated in certain study designs such as prophylaxis trials. Finally, some endpoints defined using event-triggered assessments might not be possible to assess at baseline.

<sup>20</sup> When final, this guidance will represent the FDA's current thinking on this topic.

## *Contains Nonbinding Recommendations*

*Draft — Not for Implementation*

638 concept of interest. An ordinal modeling approach (e.g., cumulative logistic regression;  
639 Agresti, 2013; Harrell, 2015) has different assumptions than a general linear model and  
640 may incorporate more information in the endpoint than the dichotomization approach.  
641 The key point when choosing an analytic approach is that the results are interpretable and  
642 address the appropriate clinical question. Regardless of the approach taken, sponsors  
643 should explore the potential impact of violation of assumptions.  
644

- 645 • *A multi-component endpoint constructed by assigning ordinal values based on scores*  
646 *reflecting multiple aspects of health.* This type of multi-component ordinal endpoint  
647 might mix distinct aspects of a disease, such as symptom levels, hospitalization, and  
648 death. The ordinal values are assigned by an algorithm to reflect increasingly severe  
649 disease states. While the same analytic approaches could be considered for this type of  
650 ordinal endpoint, greater caution is required in interpreting the findings. There could be a  
651 situation where ordinal multi-component endpoints that mix distinct aspects of a disease  
652 in which treatments are beneficial in terms of one aspect of health (e.g., severity of  
653 symptoms) but are harmful in terms of another aspect (e.g., mortality). It is possible in  
654 these situations that estimates of treatment effect from common analytic methods such as  
655 ANCOVA and cumulative logistic regression may show overall treatment benefit but  
656 could obscure harmful effects. Sponsors should consult FDA when developing analytic  
657 plans for such ordinal, multi-component endpoints.  
658

### 3. *Missing Data*

659  
660  
661 Missing data are problematic because they may lead to reduced power and potential bias in the  
662 estimated treatment effect when missingness is related to treatment effectiveness or to adverse  
663 events from the treatment. Two types of missingness can occur for COA-based endpoints: (1)  
664 missing responses to items or tasks that make up a COA; and (2) missing an entire COA at a  
665 given time point.<sup>21</sup>  
666

667 Every effort should be made to avoid missing COA data. This begins with collecting only those  
668 COAs necessary to assess the endpoint (e.g., for efficacy, safety, tolerability) and designing a  
669 data collection plan that is least burdensome and as easy as possible for patients and/or  
670 caregivers. This includes counseling patients on the importance of completing the COA and  
671 providing reminders when the patient needs to complete the COA. When a person does not  
672 complete a COA at a given time point, the site should be notified so that research staff can  
673 contact the appropriate person (patient, caregiver, study, or site staff) to obtain the needed  
674 assessment. It is important to collect reasons for missing data to inform suitable sensitivity  
675 analyses of the study endpoints considering different approaches to account for the missing data.  
676 The ability of the COA-based endpoint to address the clinical question of interest will depend on  
677 the amount of and reasons for missing data and how plausible the missing data assumptions are  
678 for the study.  
679

---

<sup>21</sup> Missing data should be distinguished from intercurrent events (e.g., death). Within the estimand framework, intercurrent events are things that happen after randomization that might affect the ability to observe or the interpretation of an endpoint. Potential intercurrent events and methods to handle intercurrent events should be addressed in the statistical analysis plan. For additional details, please see ICH E9(R1).

## *Contains Nonbinding Recommendations*

*Draft — Not for Implementation*

680 Missing item-level COA data should be handled based on the scoring algorithm for the  
681 instrument. In cases where patient-level COA data are missing for the entire domain(s) or the  
682 entire measurement(s), sponsors should propose statistical methods that properly account for  
683 missing data with respect to a particular estimand.

684  
685 Methods to handle the missing data for a COA-based endpoint should be aligned with the  
686 estimand of interest and addressed in the statistical analysis plan.

687  
688

### 689 **III. EVALUATING THE MEANINGFULNESS OF TREATMENT BENEFIT<sup>22</sup>**

690  
691 In regulatory decision-making, FDA evaluates how well results of a COA-based endpoint  
692 correspond to a treatment benefit that is meaningful to patients. For endpoints based on COAs  
693 intended to reflect how patients feel or function (see Section I.B), sponsors should provide  
694 supporting evidence to justify the meaningfulness of an observed treatment benefit. Section III  
695 discusses what supporting evidence is recommended, how it could be collected, and how it can  
696 be applied to help interpret the trial results. FDA strongly recommends that sponsors seek FDA  
697 input as early as possible regarding the evaluation of meaningful treatment benefit.

698  
699

#### 699 **A. Factors Affecting the Interpretability of COA Scores**

700  
701 To determine whether a medical product has a positive, meaningful effect on how a patient feels  
702 or functions (i.e., a treatment benefit<sup>23</sup>), FDA recommends that sponsors measure how a patient's  
703 status on a COA-based endpoint corresponds to the way they feel and/or function in their daily  
704 life. For example, if a treatment is shown to reduce scores on a performance outcome measure  
705 by an average of 2 points on a 15-point scale, it would be helpful to know whether a 2-point  
706 difference corresponds to something that patients would notice as important in their daily lives.  
707 Or, if a treatment is expected to increase a patient's score on a measure of functioning from 12 to  
708 18, it would be helpful to know what kinds of things the patient could do (or do more easily)  
709 corresponding to a score of 18 versus 12. Knowing how COA scores relate to patients'  
710 experiences is central to interpreting the meaningfulness of a COA-based endpoint result(s).  
711 This is true whether the endpoint is based on scores generated from a single COA or multiple  
712 COAs (as in a multi-component endpoint).

713  
714 Some COAs might produce scores that are easier to interpret than other COAs in terms of  
715 patients' experiences. How easily one can interpret a COA score depends on at least two  
716 factors:

717

---

<sup>22</sup> Most of the methods described in this section for interpreting trial results can apply to treatment impacts other than those described as "benefit." These could include treatment tolerability or harm in terms of how the patients feel, function, or survive. However, for brevity this section will refer only to treatment benefit.

<sup>23</sup> Treatment benefit is also demonstrated by a favorable effect on how patients survive, but this is not relevant for the discussion of COA-based endpoints.

## *Contains Nonbinding Recommendations*

*Draft — Not for Implementation*

718  
719 1. *How Closely Does the Measured Concept of Interest Correspond to the Patients’*  
720 *Experiences?*  
721

722 Some COAs measure a concept of interest that is a directly interpretable reflection of the  
723 patients’ health-related experiences, such as a PRO measure of current pain intensity. For such  
724 measures, it may be relatively easy to infer how different scores on the measure correspond to  
725 different experiences the patients might have. Other COAs might measure a concept of interest  
726 that is more indirectly related to the patient’s health-related experiences, such as an ObsRO  
727 measure of the patient’s pain behavior (which is indirectly related to the patient’s actual pain) or  
728 a PerfO measure of leg strength (which is indirectly related to activities that require lower limb  
729 function such as walking or climbing stairs).<sup>24</sup> For these types of measures, it may be more  
730 challenging to infer how different scores on the measure correspond to different experiences the  
731 patients might have; this means that additional empirical support is needed to translate scores on  
732 the measures to corresponding patient experiences in their daily lives.  
733

734 2. *How Simple or Familiar is the COA’s Metric?*  
735

736 In addition to how closely the concept of interest corresponds to the patient’s direct experience,  
737 the metric that is used to express the COA scores can also be more or less easy to interpret.  
738 Some COAs produce scores that are easier to interpret on their own because they use a metric  
739 that is relatively simple and/or familiar. For example, a daily diary that records the number of  
740 times per night that a patient woke up to urinate would generate a directly interpretable metric  
741 (i.e., number of times per night). Another example might be a simple ordinal rating of pain  
742 severity (e.g., none, mild, moderate, severe) that generates a score that most patients have little  
743 trouble interpreting in terms of noticeable gradations between patients’ experiences. Cognitive  
744 interview data might confirm that patients are comfortable evaluating their symptom severity  
745 with this scale and that patients view each category as corresponding to a meaningfully distinct  
746 experience. In this case, the scores themselves are directly interpretable in terms of patients’  
747 experiences, and therefore, additional supporting evidence may not be necessary for  
748 interpretation.  
749

750 Other COAs produce scores that are more difficult to interpret on their own because they use a  
751 metric that is unfamiliar and/or abstract, such as a COA measure that produces transformed  
752 scores (e.g., linear transformation of a 0-4 raw score scale to a 0-100 score scale). There might  
753 be very good reasons to generate a score on such a metric, but it increases the complexity of  
754 describing the endpoint in labeling. In this case, FDA recommends additional evidence to justify  
755 how scores relate to meaningful patient experiences.  
756

---

<sup>24</sup> Indirect measures of patients’ experiences could be recommended for many reasons, including the patients being incapable of self-reporting (e.g., too young, suffering from cognitive impairments) or a concern that heterogeneity in environments will create undesirable noise in self-reports of functioning (which may suggest the use of a PerfO measure).

## Contains Nonbinding Recommendations

Draft—Not for Implementation

757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800

### **B. Approaches for Collecting Evidence to Support Interpretability of COA-Based Endpoints**

Sponsors should first review any existing evidence in support of the interpretability of the COA scores used to construct the endpoints. If the body of evidence supporting the interpretability of COA scores (e.g., from existing literature) is not sufficient, FDA recommends conducting empirical studies to support interpretability of COA scores prior to conducting a registration trial. When feasible, it is advantageous to use multiple methods to inform interpretations of scores. It is expected that empirical approaches will generate a range of plausible estimates reflecting the inherent uncertainty in interpreting scores. Based on such empirical studies, sponsors should prespecify the range of estimates that will be used to interpret the treatment effect(s) in a registration trial. The following sections describe two general approaches for conducting empirical studies to support the interpretability of COA scores—interpreting in terms of meaningful score differences (III.B.1) and in terms of meaningful score regions (III.B.2).

#### *1. Interpreting in Terms of Meaningful Score Differences*

This first approach identifies what size difference between any two COA scores would be viewed as meaningful for patients. This will be referred to as the *meaningful score difference (MSD)*. Often, *MSD* is determined based on what patients would regard as a clinically meaningful within-patient change (i.e., improvement or deterioration from the patient’s perspective), but other approaches might also be appropriate (e.g., those based on the patient’s perception of the differences between hypothetical vignettes representing different degrees of symptom severity or functioning). Note that patients differ in their views of what might count as *MSD*, but for purposes of evaluating the results of clinical trials, a range of *MSD* should be selected that reflects most patients.

Regardless of the approach used to determine the *MSD*, the *MSD* can be used in at least two ways: (1) to evaluate the expected treatment effect for the average patient in some target population; or (2) to use as a threshold in descriptive analyses that identify individual patients who might have changed by a meaningful amount. Both of these applications will be discussed (see III.C) following a review of approaches for selecting a value or range of values for *MSD*.

Key assumptions should be identified and evaluated before *MSD* can be used to interpret the meaningfulness of a treatment effect in a clinical trial. Two common assumptions that should be evaluated are the following:

- The value of *MSD* is the same regardless of the baseline COA score (Crosby et al. 2003). For example, if *MSD* is specified as 4 points, then score differences of 5-1, 10-6, and 15-11 should all be regarded as meaningful differences by patients. If this assumption is not true, it is possible to use different values for *MSD* depending on the patient’s baseline status.

## Contains Nonbinding Recommendations

Draft—Not for Implementation

- 801       • The value of *MSD* is the same for improvement and deterioration (Crosby et al. 2003). If  
802 this assumption is not true, then it is possible to use different values for *MSD* depending  
803 on the direction of change.  
804

805 Sponsors can consider the use of anchor-based methods for identifying *MSD*. An anchor is some  
806 external variable, not derived from the COA whose scores require interpretation, for which  
807 meaningful differences are directly interpretable or already known.<sup>25</sup> Meaningful differences on  
808 the anchor can then be mapped onto differences in terms of the COA scores. For example, a  
809 patients' categorizations of their change in symptom severity (much better, a little better, no  
810 change, a little worse, much worse) could be used to find the range of changes in a multi-item  
811 COA that correspond to patients endorsing their change in symptom severity as "much better."  
812 (Considerations for the use of anchors are discussed in the next two sections.) Distribution-  
813 based methods (e.g., effect sizes, certain proportions of the standard deviation and/or standard  
814 error of measurement) do not directly consider the patient voice, and as such, are insufficient to  
815 serve as the sole basis for identifying an *MSD*. Distribution-based methods can provide helpful  
816 information about measurement variability. FDA is open to discussion about other well-justified  
817 methods developed for determining thresholds for *MSD* (e.g., Idio Scale Judgment; Cook et al.  
818 2017).

819  
820           a.       Choice of anchor variables

821  
822 FDA recommends that sponsors use multiple anchor measures to inform decisions about a  
823 plausible range of *MSD* values. Several factors should be considered when choosing anchor  
824 measures and, in the case of multiple anchor variables, when deciding how much weight to give  
825 an anchor when specifying *MSD* values:  
826

- 827       • Ideally, the concept assessed by an anchor variable should match or be inclusive of the  
828 concept of interest being assessed by the COA-based endpoint. For example, a sponsor  
829 might propose a single item assessing the patient's global impression of severity for a  
830 symptom to use as an anchor variable to help interpret scores on a multi-item patient-  
831 reported outcome measure of severity for the same symptom. Sometimes it may not be  
832 possible to find an anchor that is a direct reflection of the patients' experiences related to  
833 the concept of interest measured by the COA-based endpoint. In such cases, sponsors  
834 can consider using multiple, less directly related anchors to aid in the interpretation of a  
835 meaningful difference in scores.  
836
- 837       • An anchor should be plainly understood by respondents in the context of use. FDA  
838 recommends testing the proposed anchor item(s), including their response categories, in  
839 cognitive interviews.  
840
- 841       • An anchor should have a well-justified definition for meaningful change or for  
842 meaningful increments. For example, consider the case of a single-item ordinal anchor to

---

<sup>25</sup> While it might be similar to the COA, an anchor variable would typically not be useful as the basis for the trial endpoint because it may be less sensitive than the COA and/or address a concept of interest that is broader or more specific than the concept of interest measured by the COA.



## Contains Nonbinding Recommendations

Draft — Not for Implementation

843 measure patients' perceptions of their symptom severity (e.g., with response options of  
844 none, mild, moderate, severe). Such an anchor might be used, for example, to help  
845 interpret scores from a multi-item COA intended to measure a symptom's severity.  
846 Qualitative data collected as part of cognitive interviews with patients could help to  
847 establish whether patients believe that the anchor variable's response options adequately  
848 represent meaningfully different experiences in their daily lives.

- 849
- 850 • Differences in COA scores should be related to differences documented by one or more  
851 anchors.<sup>26</sup> The stronger the relationship, the more confidence in translating differences in  
852 the anchor to differences in COA scores.  
853
- 854 • Selected anchors should be assessed at comparable time points to the target COA.  
855 Sponsors should also ensure that, where applicable, the recall period of the anchor  
856 measure is consistent with the period covered by the COA-based endpoint.  
857
- 858 • Sometimes sponsors wish to use a Global Impression of Change as an anchor, for  
859 example, a Patient Global Impression of Change (PGIC), in which patients report the  
860 direction and extent of change they have undergone between baseline and a follow-up  
861 time point using an ordinal categorical response scale. There should be evidence that the  
862 Global Impression of Change reflects the patient's/observer's/clinician's perception of  
863 the change they experienced (in the case of the patient) or observed (in the case of an  
864 observer or clinician). The usefulness of the Global Impression of Change as an anchor  
865 is reduced when there is excessive recall error and/or present state bias (i.e., the  
866 impression of change is influenced by the patient's status at follow-up more than the  
867 patient's actual change).  
868
- 869 • Sometimes sponsors wish to use a Global Impression of Severity as an anchor, for  
870 example Patient Global Impression of Severity (PGIS), in which  
871 patients/observers/clinicians report the current or recent status of the severity or  
872 observation of symptoms or degree of functioning using a single ordinal response  
873 scale. Note that PGIS can be used to support either an *MSD* approach (by relating  
874 changes in the PGIS to changes in COA scores) or, as will be discussed in Section  
875 III.B.2, a meaningful score regions (MSRs) approach (by relating COA scores to their  
876 most likely PGIS response category).  
877

878 In some situations, an acceptable anchor variable will not exist. When a suitable anchor cannot  
879 be found, sponsors can consider other methods to inform the choice of *MSD*, such as Idio Scale  
880 Judgment (Cook et al. 2017).  
881

---

<sup>26</sup> Note that "differences in COA scores" is used here as a general term that includes differences that occur over time within a patient, i.e., changes in COA scores.

## Contains Nonbinding Recommendations

Draft — Not for Implementation

882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
  
916  
917

b. Analyses of anchors to inform choice of meaningful score difference

There are several options for relating differences in COA scores to anchor measures to arrive at *MSDs* (Coon and Cook 2018). Regardless of the analytic approach used, the following principles apply:

- Examine the distribution of the anchor scores or changes in anchor scores to ensure there is adequate variability for purposes of analysis. When changes in anchor scores are of interest, changes in the anchor scores should also be examined by baseline anchor score.
- Clearly describe the relationship between the COA score differences and the anchor (e.g., PGIC) or change in the anchor score (e.g., PGIS).
- Represent the distribution of COA difference scores corresponding to each response level of the anchor (e.g., PGIC) or each level of change in the anchor (e.g., PGIS). This presentation helps to inform a reasonable range of *MSD* estimates based on the heterogeneity among the patients studied.
- For ordinaly-scaled anchors measured at two time points (e.g., PGIS), sponsors should first determine, based on evidence, what size changes in the anchor are regarded as meaningful (e.g., 1-category, 2-category). For each level of potentially meaningful change in the anchor (e.g., 1-category), sponsors should examine the distribution of COA difference scores separately by baseline anchor response. See Table 1 for an example table shell that could be used to determine for patients who experienced a 1-category improvement in the PGIS whether the COA change scores are distributed differently depending upon the patient’s baseline PGIS category.

In Table 1, the lowest PGIS category of “None” is not shown because it is impossible for a patient with no severity to experience improvement in their PGIS.

**Table 1. Sample Table Shell To Display the Distribution of COA Change-From-Baseline Scores for Patients With a 1-Category Improvement in Patient Global Impression of Severity.**

PGIS at Baseline	N (%)	Change in COA Score from Baseline to End of Study				
		10 <sup>th</sup> Percentile	25 <sup>th</sup> Percentile	50 <sup>th</sup> Percentile	75 <sup>th</sup> Percentile	90 <sup>th</sup> Percentile
<i>Mild</i>						
<i>Moderate</i>						
<i>Severe</i>						
<i>Very Severe</i>						

## *Contains Nonbinding Recommendations*

*Draft — Not for Implementation*

- 918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955
- To select a range of thresholds to define *MSD*, sponsors should consider the following:<sup>27</sup>
    - Any choice of threshold *MSD* that attempts to distinguish between meaningful and non-meaningful differences will not correspond to some patients' experiences. That is, a difference below *MSD*, as measured, could be experienced as meaningful by some patients or a difference above *MSD*, as measured, could be experienced as not meaningful by some patients. Sponsors should consider and seek FDA input on how best to balance these two types of errors in the context of use. Note that this issue applies to any method used to derive thresholds, including anchor-based methods.
    - Generally, a wider range of thresholds should be selected when there is greater uncertainty about what patients would regard as an impactful difference. (Note that subsequent use of a wider range of thresholds to interpret a treatment effect will translate into correspondingly greater uncertainty about whether an obtained treatment effect is considered meaningful to patients.) A wider range of thresholds should be considered when any of the following are true:
      - There is a lower association between the COA difference scores and the anchor values, resulting in substantial overlap in the distributions of COA difference scores corresponding to different levels of the anchor scores (or differences between anchor scores). The greater the overlap, the less certainty there is that a given difference in COA score corresponds to a noticeable difference as indicated by the anchor. (See Coon and Cook 2018 for analytic approaches to examining overlap in distributions.)
      - Analyses of multiple anchor variables have generated different estimates of *MSD*. Note that in considering the range of *MSD*, threshold estimates from some anchors can be weighted more heavily than those estimates from other anchors based on the quality of the anchor (see III.B.1.a).
      - Analyses of the same anchor variable across multiple studies have generated different estimates of *MSD*.
      - There are several important prespecified patient subgroups, and analyses of the same anchor variable might generate different findings for different patient subtypes.

### 2. *Interpreting in Terms of Meaningful Score Regions*

956  
957  
958  
959  
960

Another approach for interpreting the meaningfulness of treatment effect is to specify the meaning of individual COA scores so that it is easier to judge whether two or more scores (e.g., treatment group means at a prespecified time point) correspond to distinct health-related

---

<sup>27</sup> For a discussion of different methods for determining a threshold of meaningful score differences, see Coon and Cook 2018.

## Contains Nonbinding Recommendations

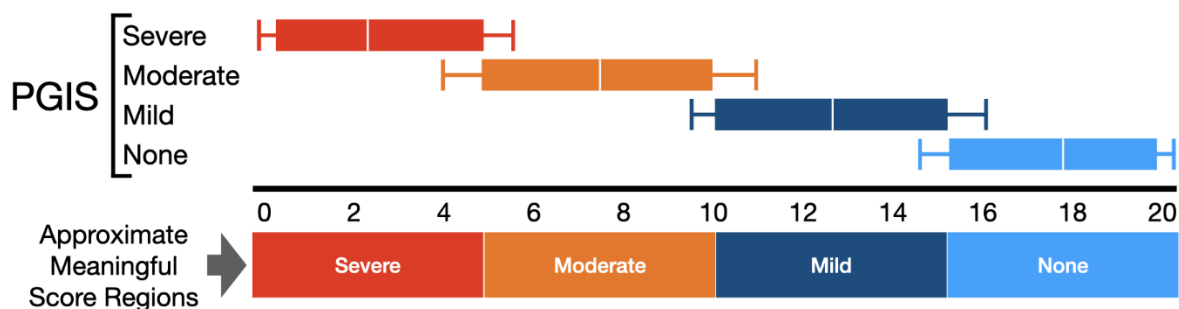
Draft—Not for Implementation

961 experiences of patients. For example, consider a measure of functioning that can generate scores  
962 from 0 to 20. Based on a study conducted with an independent sample of patients using the  
963 PGIS as an anchor, a figure can be constructed (Figure 1) to illustrate how different scores  
964 correspond to patients' global judgments of their functional impairment (none, mild, moderate,  
965 or severe). Assuming that the criteria for a strong anchor have been met (see III.B.1.a), the  
966 distributions of COA scores by PGIS response category could be examined to inform an  
967 approximate division of the COA score range into meaningful score regions (MSRs), as shown at  
968 the bottom of Figure 1. (Note that the figure shows an example in which the MSRs have equal  
969 widths; in other cases, the widths might differ.) In a later section (III.C), it is shown how MSRs  
970 could be used to help interpret a treatment effect on a COA-based endpoint.

971  
972 In Figure 1, Box-and-whisker plots display the 25<sup>th</sup> (left edge of box), 50<sup>th</sup> (white line inside the  
973 box), and 75<sup>th</sup> (right edge of box) percentiles of the COA score distributions corresponding to  
974 each PGIS level. Whiskers indicate scores  $\pm 1.5$  interquartile range. *Approximate meaningful  
975 score regions* denote groups of scores that are thought to be similar to one another and different  
976 from other groups of scores in terms of the patient's experience of the symptom(s) measured by  
977 the COA.

978

979 **Figure 1. Example of Approach for Interpreting COA Scores in Terms of Meaningful**  
980 **Score Regions Corresponding to Patient Global Impression of Severity (PGIS).**



981  
982 Different approaches to translate COA scores into their corresponding patient experiences may  
983 be appropriate if the approach is well justified within the context of use. Such approaches might  
984 include the following:

- 985
- 986 • Bookmarking or similar methods in which patients, caregivers, and/or clinicians make  
987 judgments to sort patient experiences into a small number of ordinal categories (e.g.,  
988 none, mild, moderate, or severe) (Cook et al. 2019). By determining the COA scores  
989 corresponding to those patient experiences, it is possible to identify the COA score ranges  
990 or zones that correspond to the different ordinal levels.
  - 991 • For COAs containing multiple items that are all thought to reflect the same underlying  
992 concept of interest, such as lower limb mobility, another way to facilitate interpretation of  
993 COA scores is to use one or more illustrative items from the COA measure to help  
994

## *Contains Nonbinding Recommendations*

*Draft — Not for Implementation*

995 identify *MSRs*.<sup>28</sup> Essentially, this approach uses one or more of the COA’s own items to  
996 serve as a kind of internal anchor variable.<sup>29</sup> If the illustrative item’s response categories  
997 are easy to interpret in terms of patients’ experiences, then this can be done by showing  
998 the predicted illustrative item responses for two or more COA scores. This allows a  
999 comparison of COA scores in terms of different ways the patient might feel or function as  
1000 described by the illustrative item. For example, imagine a multi-item PRO measure of  
1001 lower limb mobility with scores that range from 0 (poor mobility) to 100 (excellent  
1002 mobility). Assume that the sponsor predefined *MSRs* based on data collected prior to the  
1003 clinical trial by examining the relationship between scores on the PRO measure and  
1004 responses to an individual item from the same measure that asks about difficulty walking  
1005 up a flight of stairs. In this case, the response options for the individual item serve as  
1006 approximate *MSRs* to guide interpretation of the expected scores in each treatment group.  
1007 Suppose the mean scores for the randomized groups at a predefined follow-up time were  
1008 40 and 60. The *MSR* corresponding best to a score of 40 is “much difficulty” walking up  
1009 a flight of stairs, compared to “little difficulty” for people whose score is 60. Items  
1010 selected to serve as illustrative items should have item responses that are easily  
1011 interpretable and are strongly associated with the COA score.  
1012

- 1013 • For measures developed using Item Response Theory (IRT) (Chang and Reeve 2005), the  
1014 meaning of different scores can be enhanced by using IRT item parameters to locate  
1015 different items onto the measure’s metric. For example, if a sponsor were using an IRT-  
1016 based measure whose items assessed the level of assistance a patient needs to do different  
1017 activities, the sponsor could show the activities that patients would be predicted to do  
1018 “with no assistance” for different scores.

### 1020 3. *Additional Considerations for Justifying Meaningful Differences or Meaningful* 1021 *Score Regions*

- 1023 • FDA recommends that sponsors seek FDA input early regarding plans for determining  
1024 *MSDs* (III.B.1) or *MSRs* (III.B.2). Ideally sponsors should evaluate and provide  
1025 estimates of meaningful differences or scores prior to the start of the registration trial(s).  
1026
- 1027 • When justifying a meaningful difference using transformed data, the sponsor should  
1028 provide the threshold on the transformed and raw scales to aid in interpretation. For  
1029 multi-item measures using a transformed scale, it is critical that the threshold *MSD* be at  
1030 least equal to or greater than a one-category change for at least one item on the raw  
1031 (untransformed) scale.
- 1032
- 1033 • For situations in which it is not feasible to obtain information to inform meaningful  
1034 differences or scores before a registration trial (e.g., rare disease trials), sponsors can  
1035 consider using exit interviews or surveys (refer to PFDD Guidance 2). Patients or their  
1036 caregivers could be asked questions such as whether the patient experienced a change in

---

<sup>28</sup> This approach is known more generally as *content-based interpretation* (section 11.1.4 in Cappelleri et al. 2014).

<sup>29</sup> It is “internal” in the sense that the item is part of the COA and is used along with other items to generate a score for the COA.

## Contains Nonbinding Recommendations

Draft — Not for Implementation

1037 their symptoms from baseline, whether the change was an improvement or worsening,  
1038 and whether they believe the change in symptoms was meaningful (e.g., they can now  
1039 walk around their house without assistance). The interviews should be conducted after  
1040 the patients complete the main portion of the study to avoid any potential compromise to  
1041 trial integrity. Note that this approach is susceptible to greater bias than other approaches  
1042 and generally should only be used in trials in which patients and/or caregivers are  
1043 unaware of their study group assignment. Sponsors who are considering conducting exit  
1044 interviews or surveys should submit a study protocol and interview guide to FDA for  
1045 review as early as possible, ideally prior to beginning the registration trial.<sup>30</sup>  
1046

- 1047 • If sponsors wish to use data cited in the literature to propose *MSDs* or *MSRs*, sponsors  
1048 should explain why it is reasonable to generalize the *MSDs* and *MSRs* from the literature  
1049 to aid in interpreting the results of their registration trial. It is important to evaluate the  
1050 comparability of context between the literature and the registration trial under  
1051 consideration in terms of relevant factors such as disease, patient population, background  
1052 standard of care, location, calendar time,<sup>31</sup> COA version,<sup>32</sup> endpoints, and length of  
1053 follow-up.  
1054

### 1055 C. Applying Information About Meaningful Score Differences or Meaningful 1056 Score Regions to Clinical Trial Data

1057 Information about meaningful differences or scores can be used to help interpret the  
1058 meaningfulness of treatment effects within a clinical trial. Determining whether a medical  
1059 product produces an effect that is meaningful to patients involves careful consideration of  
1060 multiple sources of information. This could include findings from multiple endpoints (e.g.,  
1061 primary and secondary endpoints), multiple anchors that inform a range of *MSDs* or *MSRs*,  
1062 prespecified sensitivity analyses to supplement the main trial analysis of the COA-based  
1063 endpoint, analyses to examine heterogeneity of treatment effect, and graphical and/or exploratory  
1064 analyses to examine analytic assumptions or illustrate findings in alternative ways. Stakeholders  
1065 should consider the strength of evidence to support decision making and the general  
1066 considerations described in this section when creating justifications to support  
1067 the interpretation of clinical trial data. In the broader picture of marketing authorization  
1068 decisions, there are many factors to weigh simultaneously when making a decision about  
1069 meaningfulness.  
1070

1071 Sponsors should prespecify the method(s) used to interpret COA-based treatment effects and to  
1072 convey the uncertainty around guides for score interpretation (e.g., estimates of *MSD* or *MSRs*)  
1073 through describing a range of likely values, confidence intervals, or other representations of the  
1074 uncertainty. The specific method of applying *MSDs* or *MSRs* will depend on the type of COA-  
1075 based endpoint and the approaches taken to analyze the trial outcomes. The considerations and  
1076

---

<sup>30</sup> For a review of emerging qualitative methods for informing estimates of meaningful differences, see Staunton et al. 2019.

<sup>31</sup> Consider any changes relevant to the estimation of *MSDs* and *MSRs* that might have occurred since the time the study or studies in the literature were conducted.

<sup>32</sup> Note that a COA refers to any instructions, administration materials, content, formatting, and scoring rules associated with a COA.

## Contains Nonbinding Recommendations

Draft—Not for Implementation

1077 examples in this section are meant to provide general suggestions for how to approach the  
1078 interpretation of COA-based treatment effects.

1079  
1080 Note that the roles of *MSD* or *MSRs* differ depending upon the type of endpoint. For endpoints  
1081 based on continuous COA scores, the *MSD* or *MSRs* help to interpret the treatment effect. For  
1082 this application, the sponsor can prespecify a range of *MSD* or *MSRs* that will be used to aid  
1083 interpretation. For endpoints based on categorizing COA scores (e.g., a “responder” endpoint),  
1084 the *MSD* or *MSRs* define the endpoint. In that case, the sponsor should prespecify a single  
1085 threshold (for *MSD*) or set of thresholds (for *MSRs*) that will be used to define the endpoints.

1086  
1087 *1. Interpreting the Meaningfulness of Continuous COA-Based Endpoints*

1088  
1089 Different approaches can be used for interpreting treatment effects in terms of continuous COA-  
1090 based endpoints depending upon whether *MSDs* or *MSRs* are used to aid in interpretation.

1091  
1092 *a. Meaningful score difference approach*

1093  
1094 An important consideration when applying *MSDs* to interpret a continuous COA-based endpoint  
1095 is whether the estimates of *MSD* are relatively the same regardless of the patients’ baseline COA  
1096 scores. Sponsors who plan to interpret trial results in terms of *MSDs* should have already  
1097 collected or cited evidence to evaluate this possibility.<sup>33</sup>

1098  
1099 • *If there is evidence that MSD is relatively consistent over all baseline scores:* In this case,  
1100 the difference between study arms may be compared to the value(s) of *MSD* to  
1101 understand the meaningfulness of the treatment effect. For example, in a hypothetical  
1102 clinical trial comparing a new product A to a current product B, scores (0-20) on a PRO  
1103 measure of functioning were analyzed using an ANCOVA with baseline PRO  
1104 functioning scores as the covariate. The primary prespecified group comparison was  
1105 conducted at 12 weeks post-randomization. Figure 2 displays the treatment effect and  
1106 95% confidence interval.<sup>34</sup> Based on three different anchor-based analyses conducted  
1107 using an independent sample of patients, the sponsor prespecified a range of *MSD* for the  
1108 PRO functioning measure of 3 to 5 points. (The sponsor also conducted analyses to show  
1109 that the value of *MSD* did not vary substantially by baseline COA score.) Because the x-  
1110 axis reflects possible differences between scores on the PRO functioning measure, one  
1111 can graph both the expected difference in scores between products A and B (i.e., the  
1112 average treatment effect) and the range of *MSDs* thought to correspond to meaningfully  
1113 different patient experiences. Figure 2 shows that values of the treatment effect that are

---

<sup>33</sup> Caution is needed when evaluating the potential baseline dependency of the *MSD*, because simple stratification on the baseline COA scores may lead to an erroneous finding of baseline dependency. There are other approaches that can be used (see Terluin B, Roos EM, Terwee CB, Thorlund JB, and LH Ingelsrud, 2021, Assessing Baseline Dependency of Anchor-Based Minimal Important Change (MIC): Don’t Stratify on the Baseline Score! *QualLife Res*, 30(10):2773-2782, doi:10.1007/s11136-021-02886-2).

<sup>34</sup> This treatment effect can be interpreted as a conditional treatment effect—that is, the treatment effect is assumed to be approximately constant across subgroups defined by the baseline PRO score in the ANCOVA model. In other words, this treatment effect is the difference in PRO score we would expect for the average patient. See FDA’s draft guidance for industry *Adjusting for Covariates in Randomized Clinical Trials for Drugs and Biological Products* (May 2021). When final, this guidance will represent the FDA’s current thinking on this topic.

## Contains Nonbinding Recommendations

Draft — Not for Implementation

1114 consistent with the observed data (reflected by the 95% confidence interval) are above  
1115 the maximum estimate of the threshold for *MSD*. This strongly suggests that the average  
1116 treatment effect corresponds to a difference in experience that most patients would  
1117 consider meaningful. In contrast, Figure 3 displays a scenario that does not clearly  
1118 correspond to a meaningful overall difference due to treatment using the predefined  
1119 *MSDs*, although a small portion of patients might experience a treatment effect that they  
1120 regard as meaningful.

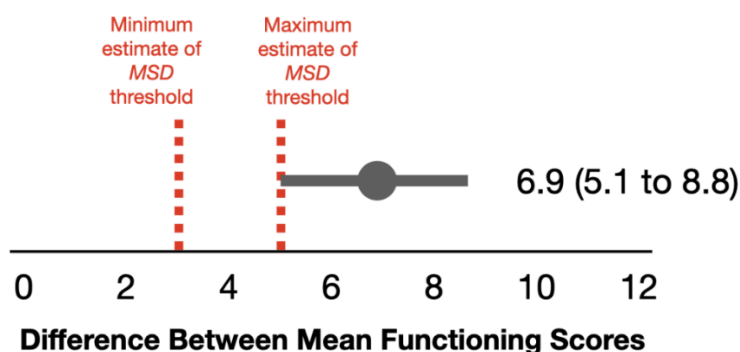
1121

1122 In Figure 2, dotted red lines indicate the minimum and maximum estimates of meaningful  
1123 difference thresholds (D thresholds) obtained from anchor-based studies conducted  
1124 independently of the registration trial. Differences greater than a threshold estimate are  
1125 considered noticeably different by patients.

1126

1127 **Figure 2. Estimated Difference in Adjusted Means (With 95% Confidence Interval)**  
1128 **Between Products A and B on Functioning Measure Scores at Follow-Up Time Point**  
1129 **Relative to Thresholds for Meaningful Score Differences**

1130



1131

1132 In Figure 3, dotted red lines indicate the minimum and maximum estimates of meaningful  
1133 difference thresholds (D thresholds) obtained from anchor-based studies conducted  
1134 independently of the registration trial. Differences greater than a threshold estimate are  
1135 considered noticeably different by patients.

1136

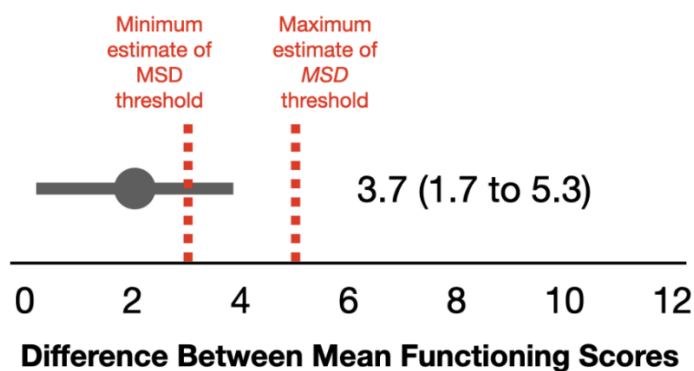
1137



## Contains Nonbinding Recommendations

Draft — Not for Implementation

1138  
1139 **Figure 3. Estimated Difference in Adjusted Means (With 95% Confidence Interval)**  
1140 **Between Products A and B on Functioning Measure Scores at Follow-up Time Point**  
1141 **Relative to Thresholds for Meaningful Differences**



- 1142
- 1143
- 1144
- 1145
- 1146
- 1147
- 1148
- 1149
- 1150
- 1151
- 1152
- 1153
- 1154
- *If there is evidence that MSD varies substantially depending upon the patients' baseline scores:* This might occur if, for example, estimates of *MSD* ranged from 2 to 8 with larger values of *MSD* found for patients whose baseline COA scores reflected lower severity. In this case, if the treatment effect is larger than the largest estimate of *MSD* (e.g., *MSD* = 8 for patients who are least severe at baseline), this suggests that the treatment effect corresponds to a meaningful difference in patients' experiences. If the treatment effect is smaller than the largest estimate of *MSD*, it means that the treatment effect might be meaningful for only some or even none of the patients depending upon their baseline COA scores. To explore this, the sponsor could compare the treatment effect estimate to the estimates of *MSD* corresponding to each level of baseline COA score to better understand the meaningfulness of the treatment effect in patients across the range of baseline severity.

1155

1156 In addition to directly interpreting the estimate of treatment effect as described above, other

1157 analyses and displays may aid interpretation. If within-patient changes from baseline in the

1158 COA-based endpoint can be meaningfully estimated and interpreted from the trial data, sponsors

1159 can also plot the empirical probability density function (ePDF) or empirical cumulative

1160 distribution function (eCDF) of changes from baseline for each trial arm. The graphs should be

1161 annotated with a range of *MSD* values and the proportion of patients in each trial arm whose

1162 change-from-baseline exceeds one or more values of *MSD*. At times, other descriptive statistics

1163 by trial arm, such as the median and other quantiles of the change-from-baseline distributions,

1164 can provide additional relevant information.

1165

1166 These and other supplementary analyses should be interpreted in the context of the estimates of

1167 treatment effect overall and, if applicable, by prespecified patient subgroups. A judgement about

1168 the overall meaningfulness of the treatment effect could be made based on all the different

1169 analyses described in the example, along with data from complementary endpoints, any other

1170 clinical trials, and other factors that define the context of use.

1171

## Contains Nonbinding Recommendations

Draft — Not for Implementation

### b. Meaningful score regions approach

1172  
1173  
1174 Figure 1 (presented earlier) illustrated how a study conducted with an independent sample of  
1175 patients using the PGIS as an anchor informed a decision about approximate *MSRs*. These  
1176 regions corresponded to patients' experiences of their health state as *none, mild, moderate*, or  
1177 *severe*.

1178  
1179 When examining the treatment effect in terms of *MSRs*, sponsors should predefine whether a  
1180 difference of 1, 2, or more regions is required for patients to view the treatment effect as  
1181 meaningful. The discussion that follows uses a 1 region difference, which would need to be  
1182 supported by patient and/or caregiver input and which might not apply to other COAs and  
1183 contexts of use.

1184  
1185 An important consideration when applying the *MSRs* approach to interpret a continuous COA-  
1186 based endpoint is whether the widths of the *MSRs* are relatively similar. For example, the widths  
1187 of the regions in Figure 1 are all approximately 5 points. The following are general  
1188 considerations regarding the width of the *MSRs* and the size of the treatment effect:

- 1189
- 1190 • *If there is evidence that the widths of the MSRs are relatively similar:* In this case, if  
1191 the treatment effect is larger than the width of each of the *MSRs*, this suggests the  
1192 treatment effect could be considered meaningful (i.e., because no matter where along  
1193 the score range the treatment effect occurs, the average treatment effect will always  
1194 correspond to a difference in score regions). This is illustrated in which the overall  
1195 treatment effect is shown in terms of the adjusted means at the predefined follow-up  
1196 time generated from an ANCOVA. However, if the average treatment effect is  
1197 smaller than the common width of the *MSRs*, then additional analyses may be  
1198 necessary to understand the nature of the treatment effect, such as exploring predicted  
1199 COA scores at follow-up for each study arm over a range of baseline COA scores.  
1200 This analysis may help identify which, if any, COA values at baseline are associated  
1201 with a treatment effect that crosses two or more *MSRs*.

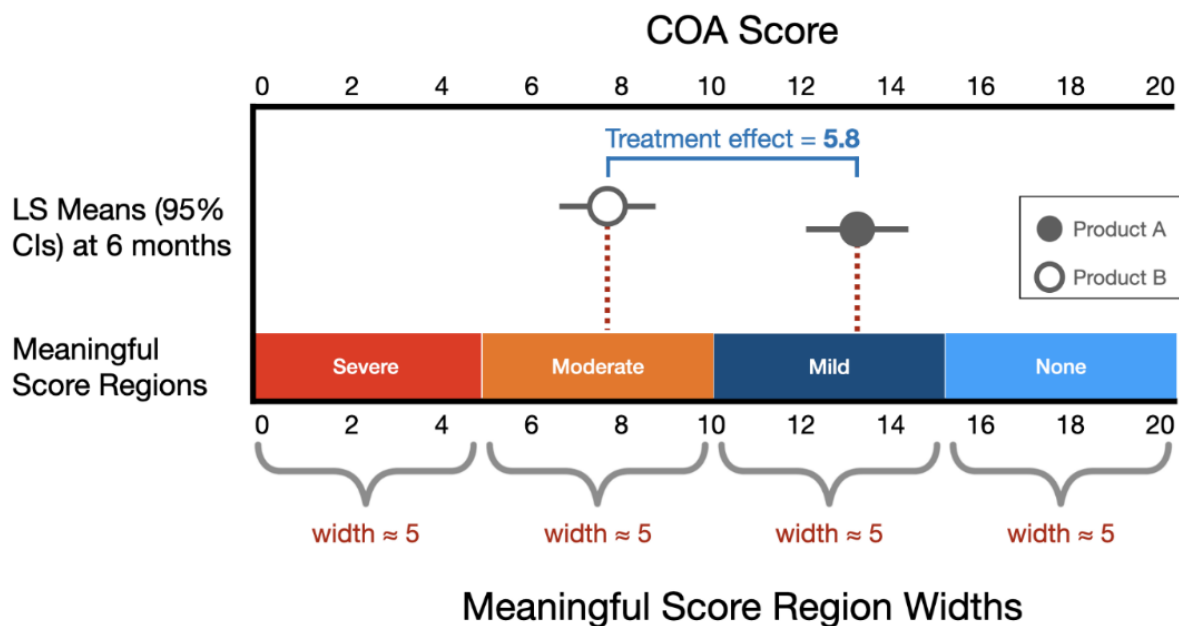
1202  
1203 In Figure 4, dotted red lines are drawn to illustrate how adjusted means are mapped onto  
1204 meaningful score regions derived using PGIS data.  
1205

## Contains Nonbinding Recommendations

Draft — Not for Implementation

1206  
1207  
1208  
1209

**Figure 4. Least Squares (LS) Means Scores (With 95% Confidence Interval) on Functioning Measure Scores at Follow-up Time Point for Products A and B Relative to Meaningful Regions of Scores Based on Patient Global Impression of Severity.**



1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229

- *If there is evidence that the widths of the score regions are relatively different:* In this case, if the treatment effect is larger than the width of the widest score region, this suggests that the treatment effect reflects a meaningful difference to patients and/or caregivers. If the treatment effect is smaller than the widest score region, then the meaningfulness of the treatment effect may be different for different patients, depending upon their baseline status. This possibility could be explored as described in the prior bullet, by examining predicted COA scores at follow-up for each study arm over a range of baseline COA scores.

In addition to directly interpreting summaries of COA scores by treatment group, sponsors may also plot the ePDF or eCDF of COA scores separately by treatment group, annotating the graph with a guide for *MSRs* (e.g., as shown in the X axis of Figure 4). Such graphs might help to assess whether, for example, a small average treatment difference is driven by a small location shift in the entire curve or by a bigger shift in a small part of the curve. Sponsors may also compute, separately by treatment group, the proportion of patients with scores at follow-up that are greater (or less) than a specific score corresponding to the border between two *MSRs* (e.g., in the example used for Figure 4, scores less than 10 would reflect moderate to severe problems with functioning).

## Contains Nonbinding Recommendations

Draft — Not for Implementation

1230  
1231           2.     *Interpreting the Meaningfulness of Ordinal and Dichotomous COA-Based*  
1232                    *Endpoints*  
1233  
1234     When a COA-based endpoint is on an ordinal scale, interpreting effects in terms of  
1235     meaningfulness to patients will depend upon the COA. Some measures produce an ordinal score  
1236     consisting of a small number of categories that may have already been shown through cognitive  
1237     interviews to be well understood and to reflect meaningfully distinct experiences of the patients  
1238     (e.g., pain intensity rating of none, mild, moderate, severe). For these types of ordinal scales, no  
1239     additional work may be needed to interpret the meaningfulness of the score, though additional  
1240     analyses might need to be done to understand the nature of the treatment effect. In contrast,  
1241     some measures might produce an ordinal score with many levels (e.g., 0 – 7) that may have been  
1242     shown through cognitive interviews to be less interpretable in terms of patients’ experiences.  
1243     Additional work is recommended using the *MSRs* approach to understand which score ranges  
1244     correspond to distinct experiences of patients.

1245  
1246     Some endpoints are based on defining a state or status with respect to a COA score (see  
1247     II.A.2.b). The status could be defined based on an *MSD* approach by classifying patients’  
1248     changes from baseline (e.g., as “observed improvement,” “observed worsening,” “no change”).  
1249     The endpoint could also be defined using a *MSRs* approach (e.g., patients scoring below some  
1250     thresholds are classified as “symptoms resolved” and those scoring at or above the threshold are  
1251     classified as “symptomatic”). For these situations, the sponsor should prespecify the threshold  
1252     (in the case of *MSD*) or set of thresholds (in the case of *MSRs*) that will be used to define the  
1253     endpoint.

1254  
1255

## 1256     **IV.     ADDITIONAL CONSIDERATIONS**

1257

### 1258           **A.     Other Study Design Considerations**

1259

#### 1260                    1.     *Masking*<sup>35</sup>

1261

1262     Patients’, clinicians’, and/or caregivers’ knowledge of treatment assignment (e.g., in single arm  
1263     trials, open label trials, open-label treatment extension periods) is likely to influence how they  
1264     report information on a PRO, ClinRO, or ObsRO measure, or how they engage with PerFO tasks  
1265     (e.g., amount of encouragement given to patients when measuring walking distance), which will  
1266     bias estimates of treatment effect. The protocol should specify to what extent masking will be  
1267     maintained among the investigators, evaluators/raters, and reporters (e.g., clinicians, patients,  
1268     caregivers).

1269

#### 1270                    2.     *Practice Effects*

1271

1272     A practice effect (sometimes also called a learning effect) is any change that results from  
1273     practice or repetition of particular tasks or activities including repeated exposure to an

---

<sup>35</sup> See footnote 20.

## *Contains Nonbinding Recommendations*

*Draft — Not for Implementation*

1274 instrument. A simple example is taking a math test, which measures math ability. After  
1275 completing the same test three times, a person’s speed (and accuracy in answering) likely will  
1276 improve because they recognize the questions and have ‘learned’ the test. While potentially an  
1277 issue for any COA, practice effects may be of particular concern in studies utilizing PerfOs with  
1278 within-subject designs in which repeated measurements are taken over time, (i.e., over the course  
1279 of the study period; American Psychological Association 2018; Shadish et al. 2002).

1280  
1281 Practice effects may be problematic for studies conducted to support a medical product  
1282 regulatory application. If severe enough, practice effects could lead to improvements in the  
1283 score of the assessment that might change the effective range of an assessment (e.g., if it creates  
1284 a ceiling effect), potentially limiting the size of the observed treatment effect, which might  
1285 impact the study’s statistical power. Aside from this possibility, in a randomized, double-  
1286 masked<sup>36</sup> trial, practice effects are unlikely to bias the difference of the outcomes between arms.  
1287 For randomized trials that are not masked, differences might arise between trial arms in practice  
1288 effects (e.g., due to differences in patient motivation or in how research staff interact with  
1289 patients) and could impact group differences in the endpoint in a way that is not due to the  
1290 treatment effect. For non-randomized trials, especially trials using external controls whose COA  
1291 assessment schedule differs from treated patients, an apparent difference (or lack of difference)  
1292 between trial arms may be due to practice effects and not due to any difference in the medical  
1293 products.

1294  
1295 Currently, approaches exist for attenuating, but not eliminating, practice effects (Jones 2015). In  
1296 addition, no consensus on best practices for attenuating practice effects has yet been reached.  
1297 Some general strategies for mitigating practice effects are summarized below. These strategies  
1298 can be used in isolation but may be more effective when used in combination.

- 1299
- 1300 • **Consider available evidence on practice effects when identifying an instrument:**  
1301 Some instruments may be more robust to practice effects than others. When selecting an  
1302 instrument, one may wish to consider available evidence of the candidate instrument’s  
1303 robustness (or vulnerability) to practice effects. Such evidence can be obtained through,  
1304 for example, a review of the literature and/or consulting the instrument’s user manual or  
1305 developer. If no evidence exists for a candidate measure, sponsors can conduct their own  
1306 empirical study of potential practice effects.  
1307
  - 1308 • **Increase length of time (spacing) between assessments:** In general—and all else being  
1309 equal—the magnitude of practice effects is expected to decrease as time between  
1310 assessments increases (Shadish et al. 2002). Decisions regarding the length of time  
1311 (spacing) to place between assessments should take into consideration how rapidly (or  
1312 slowly) change in the underlying construct is expected to occur.  
1313
  - 1314 • **Increase the length and number of assessments for the run-in period:** In general, the  
1315 magnitude of practice effects is largest at the beginning of a study and gradually levels  
1316 off or decreases as the number of assessments increases. Having a long run-in period  
1317 allows large practice effects to occur for the first few assessments until its magnitude

---

<sup>36</sup> Also referred to as “double-blind.”

## *Contains Nonbinding Recommendations*

*Draft — Not for Implementation*

1318 does not significantly increase such that the baseline and post-baseline scores are  
1319 minimally affected by practice effects. Note that this strategy would not reduce any  
1320 ceiling or floor effects caused by practice.  
1321

- 1322 • **Use alternative forms (sometimes also referred to as parallel forms or equivalent**  
1323 **forms):** Alternative forms are different versions of an instrument “that are considered  
1324 interchangeable, in that they measure the same constructs in the same ways, are built to  
1325 the same content and statistical specifications, and are administered under the same  
1326 conditions using the same directions” (American Educational Research Association  
1327 2014).  
1328

### 1329 3. *Use of Assistive Devices*

1330  
1331 If a patient starts to use an assistive device after beginning the clinical trial, the interpretation of  
1332 COA-based endpoints can be affected. Use of assistive devices may particularly impact PerFO  
1333 assessment of mobility and can impact other types of COAs (e.g., use of a walker may impact  
1334 both PerFO and PRO measures assessing physical functioning). For diseases where patients’  
1335 underlying disease status is expected to change during the trial, with corresponding changes in  
1336 the use and the type of assistive device, sponsors should consider the following:  
1337

- 1338 • Some COAs address the use of assistive devices in the instructions or administration  
1339 manual, detailing how the conduct of the assessment and scoring should occur when a  
1340 patient is using an assistive device. If this is the case, sponsors should follow the  
1341 directions for administering and scoring the chosen COA.  
1342
- 1343 • When the COA does not explicitly address how to incorporate assistive devices into the  
1344 assessment, then the sponsor should consider one of the following two strategies:<sup>37</sup>  
1345
  - 1346 – If the use of the assistive device could be influenced by treatment and altering the  
1347 need for the assistive device is one of the primary goals of treatment, then incorporate  
1348 the information on the use of assistive device into the COA-based endpoint  
1349 construction, as the use of an assistive device may reflect either an improvement or a  
1350 deterioration in the patient’s disease status.  
1351
  - 1352 – If the use of the assistive device could be influenced by treatment and altering the  
1353 need for the assistive device is not a primary goal of treatment, construct a supportive  
1354 endpoint based on whether an assistive device is used.  
1355
- 1356 • Case report forms for data collection should include information on whether an assistive  
1357 device (and what type) was used during the test.  
1358

---

<sup>37</sup> These strategies are based on the estimand framework—namely the ways to address intercurrent events (i.e., things that happen after randomization that might affect the ability to observe or the interpretation of an endpoint). For additional details, see ICH E9(R1).

## Contains Nonbinding Recommendations

Draft — Not for Implementation

### 1359 4. Considerations When Using a Nonrandomized Design, External Controls, or 1360 Nonconcurrent Control

1361  
1362 Whenever possible, COA-based endpoints should be assessed in the context of randomized,  
1363 controlled clinical trial designs. Sponsors considering COA-based endpoints in nonrandomized,  
1364 external control, or nonconcurrent control (randomized groups but at different calendar times)  
1365 trial designs should be aware of the significant potential for bias in estimating treatment effects:  
1366

- 1367 • Depending on the study, the inability to effectively mask treatment assignment could  
1368 cause group differences due to expectations of outcome held by patients, caregivers,  
1369 clinicians, or research staff. To mitigate this risk, sponsors using these designs may  
1370 consider assessing concepts of interest that require less subjective judgments (e.g.,  
1371 ability to do certain activities instead of perceived difficulty in doing activities).  
1372 Though there might still be effects of patient expectation, sponsors could also use  
1373 PerfO measures for which the patient's performance is rated by study personnel who  
1374 are masked to treatment assignment or rated automatically by some device or  
1375 computer.  
1376
- 1377 • There might be differences in the measures used to assess the concept(s) of interest,  
1378 method of COA administration, and/or the COA assessment frequency/schedule that  
1379 could lead to differences between the groups that is unrelated to the effect of  
1380 treatment. It is important to establish comparability of the COAs across the groups,  
1381 to use well-defined and reliable COA-based endpoints in conjunction with  
1382 standardized rater training and instructions for administration within each comparator  
1383 arm and across comparator arms. Every effort should be made to ensure  
1384 comparability in the assessment methods and timing of COA administration, together  
1385 with the use of standardized data collection methods (e.g., standardized modes of  
1386 administration).  
1387
- 1388 • There might be preexisting differences between the groups that affect the estimate of  
1389 treatment effect. (This potential source of bias is not unique to COA-based  
1390 endpoints.)  
1391

1392 These considerations apply to clinical trials as well as natural history studies,<sup>38</sup> disease registries,  
1393 baseline-controlled trials, and trials with a more complicated sequential on-off-on (medical  
1394 product-control-medical product) designs. Considerations for the various types of control groups  
1395 are discussed at length in the International Council for Harmonisation of Technical Requirements  
1396 for Pharmaceuticals for Human Use (ICH) guidance for industry *E10 Choice of Control Group  
1397 and Related Issues in Clinical Trials* (May 2001).<sup>39</sup>  
1398

---

<sup>38</sup> See the draft guidance for industry *Rare Diseases: Natural History Studies for Drug Development* (March 2019), *Rare Diseases: Common Issues in Drug Development* (January 2019) and final guidance for industry *Use of Real-World Evidence to Support Regulatory Decision-Making for Medical Devices* (August 2017).

<sup>39</sup> Available at the FDA guidance web page.

## *Contains Nonbinding Recommendations*

*Draft — Not for Implementation*

1399

1400

### 5. *Analysis of Treatment Effects for Subgroups Based on Post-Baseline Events*

1401

1402

1403

1404

1405

1406

If subgroups of a trial population are defined based on post-baseline events (e.g., patients who are alive and on treatment), interpretation of direct comparisons between treatment arms are likely to be misleading. By no longer reflecting the randomization intended to support a strong inference, the treatment arms will likely no longer be comparable due to differences in patient characteristics based on post-baseline events.

1407

1408

### 6. *Computerized Adaptive Testing*

1409

1410

1411

1412

1413

1414

1415

1416

1417

1418

1419

1420

1421

1422

One option for collecting scores from patients in clinical trials is to use computerized adaptive testing (CAT). This involves the use of an algorithm to iteratively select and administer items from a bank of items based on previous responses of the person being assessed. With each item that is answered, an updated estimate of the person's status on the concept of interest (e.g., symptom severity) is generated. That updated estimate is used by the CAT algorithm to select items that best match the current estimated severity and provide the most information for further estimation. The general goal of CAT is to provide individualized testing on a large scale by automatically selecting the most appropriate items for a person. However, generally the item selection is based on the likelihood that an item will be helpful in improving the estimate of the person's score, not on the relevance of the item content. (Note that special CATs can be constructed to ensure that items reflecting particular content are administered.) Thus, FDA recommends special considerations to assess whether CAT is appropriate for a given concept of interest and context of use.

1423

1424

1425

1426

1427

Because a CAT is based on IRT modeling, sponsors who wish to use CAT should demonstrate that (1) the underlying IRT parameters are statistically sound and come from the population of interest; (2) the assumptions of the IRT model and CAT are tenable; and (3) the adaptive and scoring algorithms were correctly implemented.

1428

1429

1430

1431

1432

1433

1434

Sponsors should consider the concept of interest and if the specific items have sufficient content coverage when using CAT. Hybrid CAT, where a small number of static items (i.e., those seen by all respondents) are administered along with the administration of items using the CAT algorithm, may be useful when CAT administration of items serves to supplement the static short form. When thoughtfully implemented, CAT or hybrid CAT may present advantages over static administrations, such as short forms.

1435

1436

1437

1438

1439

1440

1441

1442

1443

1444

In general, sponsors should consider whether administering items from an item bank via CAT will be more advantageous than administering a short form consisting of the same set of items to every patient in the trial. In some cases, CAT administration can bring statistical efficiency and help lower patient burden. It allows for tighter control of score reliability, while often reducing the number of items administered. However, depending on the concept of interest being measured and the range of severity in the target population, CAT may or may not provide a significant advantage over a short form in terms of precision, number of items recommended, and/or ceiling/floor effects. Research has shown that in some cases, CAT only provides benefits to measurement precision on the very high and low levels of severity when the sample is



## ***Contains Nonbinding Recommendations***

*Draft — Not for Implementation*

1445 representative of the full spectrum of severity (Choi et al. 2010; Rothrock et al. 2019; Amtmann  
1446 et al. 2018). When weighing CAT versus short form in clinical study settings, sponsors should  
1447 consider the make-up of their target population throughout the study, including at baseline, peak  
1448 effect, and end of study. For specific populations with a limited range of severity, a short form  
1449 can be created from the same item bank to target precise measurement over the range of severity  
1450 expected in the study.

1451  
1452 Sponsors should carefully consider the potential benefits and drawbacks to employing CAT in a  
1453 clinical study. Discussion and alignment with the appropriate review division are strongly  
1454 encouraged.

### 1455 1456 **7. *Minimizing Participant Burden***

1457  
1458 To demonstrate respect for the patients and/or caregivers who participate and maximize the  
1459 quality and completeness of information collected in a clinical trial, sponsors should consider  
1460 ways to minimize the burden of participation and increase the convenience and value of  
1461 participation to patients and/or caregivers. Early engagement with patient communities (see  
1462 PFDD Guidance 1) and the involvement of patient representatives in the development of a  
1463 clinical trial can improve the patient-centeredness of trial procedures and assessments. With  
1464 respect to COA-based endpoints, patient communities can provide input on the relevance, type,  
1465 length, and frequency of COAs. Pilot testing of procedures for recruitment and assessment can  
1466 also help minimize patient burden. A failure to evaluate and address potential issues with burden  
1467 or fatigue can result in a trial with greater missing data, poorer quality data (e.g., when overly  
1468 burdened participants quickly respond and select the first response to every item rather than  
1469 carefully reading and considering their answer), and/or more dropout.

### 1470 1471 **B. *Formatting and Submission Considerations***

1472  
1473 Regardless of how patient experience data is collected in a given study, patient experience data  
1474 collected and submitted to FDA to support a regulatory medical product application are subject  
1475 to statutory and regulatory submission requirements that apply to the study data and submission  
1476 type. Guidance documents that address data formatting and submission include, but are not  
1477 limited to, the following:

- 1478  
1479 • ICH guidance for industry *M8 Electronic Common Technical Document (eCTD) v4.0*  
1480 *DRAFT Implementation Guide v2.0*; and *eCTD Implementation Package DRAFT*  
1481 *Specification for Submission Formats v2.0* (April 2015)
- 1482  
1483 • Code of Federal Regulations, (CFR) Title 21, Chapter 1 (21 CFR Chapter 1)—with  
1484 particular attention given to Parts 11, 21, 312.57, 312.62(b) and (c), and 812.140
- 1485  
1486 • FDA draft guidance for industry *Use of Electronic Records and Electronic Signatures in*  
1487 *Clinical Investigations Under 21 CFR Part 11—Questions and Answers* (June 2017)<sup>40</sup>  
1488

---

<sup>40</sup> When final, this guidance will represent the FDA’s current thinking on this topic.

## Contains Nonbinding Recommendations

Draft — Not for Implementation

- 1489
- 1490
- 1491
- 1492
- 1493
- 1494
- 1495
- 1496
- 1497
- 1498
- 1499
- 1500
- 1501
- 1502
- 1503
- 1504
- 1505
- 1506
- 1507
- 1508
- 1509
- 1510
- 1511
- 1512
- 1513
- 1514
- 1515
- 1516
- 1517
- 1518
- 1519
- 1520
- 1521
- 1522
- 1523
- 1524
- 1525
- 1526
- 1527
- FDA guidance for industry *Computerized Systems Used in Clinical Investigations* (May 2007)
  - FDA guidance for industry *Electronic Source Data in Clinical Investigations* (September 2013)
  - FDA guidance for industry *Providing Regulatory Submissions in Electronic Format—Standardized Study Data* (June 2021)
  - FDA guidance for industry *Providing Regulatory Submissions in Electronic Format—Submissions Under Section 745A(a) of the Federal Food, Drug, and Cosmetic Act* (December 2014)
  - FDA guidance for industry *Providing Regulatory Submissions in Electronic Format—Certain Human Pharmaceutical Product Applications and Related Submissions Using the eCTD Specifications* (February 2020)
  - FDA Study Data Standards Resources<sup>41</sup> which includes links to FDA technical specifications documents and guidances for CDER, CBER, and CDRH including the *Study Data Technical Conformance Guide* and the *eCTD Technical Conformance Guide*. This resource and its documents are frequently updated.
- Electronic devices used to administer COAs in studies conducted to support a regulatory medical product application can present special development, testing, and deployment considerations common to digital health technologies. For example, usability studies may be needed to assess study participants’ ability to enter timely and accurate data. The following FDA guidances have more information about these considerations:
- FDA draft guidance for industry, FDA staff, and other stakeholders *Patient-Focused Drug Development: Selecting, Developing, or Modifying Fit-for-Purpose Clinical Outcome Assessments* (June 2022).<sup>42</sup>
  - FDA draft guidance for industry and FDA staff *Contents of a Complete Submission for Threshold Analyses and Human Factors Submissions to Drug and Biologic Applications* (September 2018)<sup>43</sup>
  - FDA draft guidance for industry *Comparative Analyses and Related Comparative Use Human Factors Studies for a Drug-Device Combination Product Submitted in an ANDA* (January 2017)<sup>44</sup>

---

<sup>41</sup> Available at <https://www.fda.gov/industry/fda-resources-data-standards/study-data-standards-resources>.

<sup>42</sup> When final, this guidance will represent the FDA’s current thinking on this topic.

<sup>43</sup> Ibid.

<sup>44</sup> Ibid.

## ***Contains Nonbinding Recommendations***

*Draft — Not for Implementation*

- 1528 • FDA draft guidance for industry, investigators, and other stakeholders *Digital Health*  
1529 *Technologies for Remote Data Acquisition in Clinical Investigations* (December 2021)<sup>45</sup>  
1530
- 1531 • FDA guidance for industry and FDA staff *Applying Human Factors and Usability*  
1532 *Engineering to Medical Devices* (February 2016)  
1533
- 1534 • FDA draft guidance for industry and FDA staff *Human Factors Studies and Related*  
1535 *Clinical Study Considerations in Combination Product Design and Development*  
1536 (February 2016)<sup>46</sup>  
1537
- 1538 • FDA guidances with digital health content<sup>47</sup>  
1539

1540 Sponsors may also consult SPIRIT (Calvert et al. 2018, 2021) and CONSORT (Calvert et al.  
1541 2013), consensus documents that include an extensive, detailed discussion of PRO information  
1542 that can be included in trial protocols and manuscripts to improve the completeness and clarity of  
1543 reporting. Much of the discussion in SPIRIT and the CONSORT PRO extension is applicable to  
1544 other types of COAs as well.  
1545

---

<sup>45</sup> Ibid.

<sup>46</sup> Ibid

<sup>47</sup> Available at <https://www.fda.gov/medical-devices/digital-health-center-excellence/guidances-digital-health-content>.

## Contains Nonbinding Recommendations

Draft — Not for Implementation

### REFERENCES

- 1546  
1547  
1548 *Please note that the citation of a scientific reference in this guidance does not constitute FDA's*  
1549 *endorsement of approaches or methods presented in that reference for any particular study.*  
1550 *Study designs are evaluated on a case-by-case basis under applicable legal standards.*  
1551  
1552 Agresti, A, 2013, *Categorical Data Analysis (3<sup>rd</sup> Ed)*, Hoboken, NJ: Wiley-Interscience.  
1553  
1554 American Educational Research Association, American Psychological Association, and The  
1555 National Council on Measurement in Education, 2014, *The Standards for Educational*  
1556 *and Psychological Testing*, Washington (DC): American Educational Research  
1557 Association Publications.  
1558  
1559 American Psychological Association, 2018, *APA Dictionary of Psychology*, accessed July 16,  
1560 2019, <https://dictionary.apa.org>.  
1561  
1562 Amtmann, D, A Bamer, J Kim, F Bocell, H Chung, R Park, R Salem, and BJ Hafner, 2018, A  
1563 Comparison of Computerized Adaptive Testing and Fixed-Length Short Forms for the  
1564 Prosthetic Limb Users Survey of Mobility (PLUS-M™), *Prosthet Orthot Int*, 42(5):476-  
1565 482.  
1566  
1567 Berry, DA and GD Ayers, 2006, Symmetrized Percent Change for Treatment Comparisons, *The*  
1568 *American Statistician*, 60(1):27-31.  
1569  
1570 Calvert, M, J Blazeby, DG Altman, DA Revicki, D Moher, MD Brundage, and CONSORT PRO  
1571 Group, 2013, Reporting of Patient-Reported Outcomes in Randomized Trials: the  
1572 CONSORT PRO Extension, *JAMA*, 309(8):814-822.  
1573  
1574 Calvert, M, M King, R Mercieca-Bebber, O Aiyegbusi, D Kyte, A Slade, A Chan, E Basch, J  
1575 Bell, A Bennett, V Bhatnagar, J Blazeby, A Bottomley, J Brown, M Brundage, L  
1576 Campbell, JC Cappelleri, H Draper, AC Dueck, C Ells, L Frank, RM Golub, I Griebisch,  
1577 K Haywood, A Hunn, B King-Kallimanis, L Martin, S Mitchell, T Morel, L Nelson, J  
1578 Norquist, D O'Connor, M Palmer, D Patrick, G Price, A Regnault, A Retzer, D Revicki, J  
1579 Scott, R Stephens, G Turner, A Valakas, G Velikova, M von Hildebrand, A Walker, and  
1580 L Wenzel, 2021, SPIRIT-PRO Extension Explanation and Elaboration: Guidelines for  
1581 Inclusion of Patient-Reported Outcomes in Protocols of Clinical Trials, *BMJ Open*  
1582 11(6):1-35.  
1583  
1584 Calvert, M, D Kyte, R Mercieca-Bebber, A Slade, A Chan, MT King, the SPIRIT-PRO Group, A  
1585 Hunn, A Bottomley, A Regnault, A Chan, C Ells, D O'Connor, D Revicki, D Patrick, D  
1586 Altman, E Basch, G Velikova, G Price, H Draper, J Blazeby, J Scott, J Coast, J Norquist,  
1587 J Brown, K Haywood, LL Johnson, L Campbell, L Frank, M von Hildebrand, M  
1588 Brundage, M Palmer, P Kluetz, R Stephens, RM Golub, S Mitchell, and T Groves, 2018,  
1589 Guidelines for Inclusion of Patient-Reported Outcomes in Clinical Trial Protocols: The  
1590 SPIRIT-PRO Extension, *JAMA*, 319(5):483-494.  
1591

## *Contains Nonbinding Recommendations*

*Draft — Not for Implementation*

- 1592 Cappelleri, JC, KH Zou, AG Bushmakin, JMJ Alvir, D Alemayehu, and T Symonds, 2014,  
1593 Patient-Reported Outcomes: Measurement, Implementation, and Interpretation, Boca  
1594 Raton: CRC Press.  
1595
- 1596 Chang, CH and BB Reeve, 2005, Item Response Theory and its Application to Patient-Reported  
1597 Outcomes Measurement, *Eval Health Prof*, 28(3):264-282.  
1598
- 1599 Choi, SW, SP Reise, PA Pilkonis, RD Hays, and D Cella, 2010, Efficiency of Static and  
1600 Computer Adaptive Short Forms Compared to Full-Length Measures of Depressive  
1601 Symptoms, *Qual Life Res*, 19(1):125-136.  
1602
- 1603 Coon, CD and KF Cook, 2018, Moving From Significance to Real-World Meaning: Methods for  
1604 Interpreting Change In Clinical Outcome Assessment Scores, *Qual Life Res*, 27(1):33-40.  
1605
- 1606 Cook, KF, D Cella, and BB Reeve, 2019, PRO-Bookmarking to Estimate Clinical Thresholds for  
1607 Patient-Reported Symptoms and Function, *Medical Care*, 57(5 Suppl 1):S13-S17.  
1608
- 1609 Cook, KF, MA Kallen, CD Coon, D Victorson, and DM Miller, 2017, Idio Scale Judgment:  
1610 Evaluation of a New Method for Estimating Responder Thresholds, *Qual Life Res*,  
1611 26(11):2961-2971.  
1612
- 1613 Crosby, RD, RL Kolotkin, and GR Williams, 2003, Defining Clinically Meaningful Change in  
1614 Health-Related Quality of Life, *J Clin Epidemiol*, 56(5):395-407.  
1615
- 1616 Duke Margolis Center for Health Policy, 2017, Developing Personalized Clinical Outcome  
1617 Assessments, accessed August 15, 2021,  
1618 [https://healthpolicy.duke.edu/sites/default/files/2020-03/discussion\\_guide\\_4\\_5\\_17.pdf](https://healthpolicy.duke.edu/sites/default/files/2020-03/discussion_guide_4_5_17.pdf).  
1619
- 1620 FDA Draft Guidance for Industry, Benefit-Risk Assessment for New Drug and Biological  
1621 Products, September 2021.  
1622
- 1623 FDA Draft Guidance for Industry, Comparative Analyses and Related Comparative Use Human  
1624 Factors Studies for a Drug-Device Combination Product Submitted in an ANDA, January  
1625 2017.  
1626
- 1627 FDA Draft Guidance for Industry and FDA Staff, Contents of a Complete Submission for  
1628 Threshold Analyses and Human Factors Submissions to Drug and Biologic Applications,  
1629 September 2018.  
1630
- 1631 FDA Draft Guidance for Industry, Investigators, and Other Stakeholders, Digital Health  
1632 Technologies for Remote Data Acquisition in Clinical Investigations, December 2021.  
1633
- 1634 FDA Draft Guidance for Industry and FDA Staff, Human Factors Studies and Related Clinical  
1635 Study Considerations in Combination Product Design and Development, February 2016.  
1636
- 1637 FDA Guidance for Industry, Multiple Endpoints in Clinical Trials, October 2022.

***Contains Nonbinding Recommendations***

*Draft — Not for Implementation*

- 1638  
1639 FDA Draft Guidance for Industry, Rare Diseases: Common Issues in Drug Development,  
1640 January 2019.  
1641  
1642 FDA Draft Guidance for Industry, Rare Diseases: Natural History Studies for Drug  
1643 Development, March 2019.  
1644  
1645 FDA Draft Guidance for Industry, Use of Electronic Records and Electronic Signatures in  
1646 Clinical Investigations Under 21 CFR Part 11—Questions and Answers, June 2017.  
1647  
1648 FDA Guidance for Industry and Food and Drug Administration Staff, Applying Human Factors  
1649 and Usability Engineering to Medical Devices, February 2016.  
1650  
1651 FDA Guidance for Industry, Computerized Systems Used in Clinical Investigations, May 2007.  
1652  
1653 FDA Guidance for Industry, Electronic Source Data in Clinical Investigations, September 2013.  
1654  
1655 FDA Guidance for Industry, Food and Drug Administration Staff, and other Stakeholders,  
1656 Patient-Focused Drug Development: Collecting Comprehensive and Representative  
1657 Input, June 2020.  
1658  
1659 FDA Guidance for Industry, Food and Drug Administration Staff, and Other Stakeholders,  
1660 Patient-Focused Drug Development: Methods to Identify What Is Important to Patients,  
1661 February 2022.  
1662  
1663 FDA Guidance for Industry, Food and Drug Administration Staff, and Other Stakeholders,  
1664 Patient Preference Information—Voluntary Submission, Review in Premarket Approval  
1665 Applications, Humanitarian Device Exemption Applications, Humanitarian Device  
1666 Exemption Applications and De Novo Requests, and Inclusion in Decision Summaries  
1667 and Device Labeling, August 2016.  
1668  
1669 FDA Guidance for Industry, Patient-Reported Outcome Measures: Use in Medical Product  
1670 Development to Support Labeling Claims, December 2009.  
1671  
1672 FDA Guidance for Industry, Food and Drug Administration Staff and Other Stakeholders,  
1673 Principles for Selecting, Developing, Modifying and Adapting Patient-Reported Outcome  
1674 Instruments for Use in Medical Device Evaluation, January 2022.  
1675  
1676 FDA Guidance for Industry, Providing Regulatory Submissions in Electronic Format—Certain  
1677 Human Pharmaceutical Product Applications and Related Submissions Using the eCTD  
1678 Specifications, February 2020.  
1679  
1680 FDA Guidance for Industry, Providing Regulatory Submissions in Electronic Format—  
1681 Standardized Study Data, June 2021.  
1682

***Contains Nonbinding Recommendations***

*Draft—Not for Implementation*

- 1683 FDA Guidance for Industry, Providing Regulatory Submissions in Electronic Format—  
1684 Submissions Under Section 745A(a) of the Federal Food, Drug, and Cosmetic Act,  
1685 December 2014.  
1686
- 1687 FDA Guidance for Industry and Food and Drug Administration Staff, Use of Real-World  
1688 Evidence to Support Regulatory Decision-Making for Medical Devices, August 2017.  
1689
- 1690 Gawlicki, MC, SM McKown, MJ Talbert, and BA Brandt, 2014, Application of Bother in Patient  
1691 Reported Instruments Across Cultures, *Health Qual Life Outcomes*, 12(18):1-7.  
1692
- 1693 Goldberg, TE, PD Harvey, KA Wesnes, PJ Snyder, and LS Schneider, 2015, Practice Effects  
1694 Due to Serial Cognitive Assessment: Implications for Preclinical Alzheimer's Disease  
1695 Randomized Controlled Trials, *Alzheimer's Dement (Amst)* 1(1):103-111.  
1696
- 1697 Guyatt, GH, GR Norman, EF Juniper, and LE Griffith, 2002, A Critical Look at Transition  
1698 Ratings, *J Clin Epidemiol*, 55(9):900-908.  
1699
- 1700 Harrell, F, 2015, *Regression Modeling Strategies (2<sup>nd</sup> Ed.)*, New York: Springer.  
1701
- 1702 ICH Guidance for Industry, E9(R1) Statistical Principles for Clinical Trials: Addendum on  
1703 Estimands and Sensitivity Analysis in Clinical Trials, May 2021.  
1704
- 1705 ICH Guidance for Industry, E10 Choice of Control Group and Related Issues in Clinical Trials,  
1706 May 2001.  
1707
- 1708 ICH Guidance for Industry, M8 Electronic Common Technical Document (eCTD) v4.0 DRAFT  
1709 Implementation Guide v2.0; and eCTD Implementation Package DRAFT Specification  
1710 for Submission Formats v2.0, April 2015.  
1711
- 1712 Jones RN, 2015, Practice and Retest Effects in Longitudinal Studies of Cognitive Functioning,  
1713 *Alzheimer's Dement (Amst)*, 1(1):101-102.  
1714
- 1715 Krasny-Pacini, A, J Evans, M Moore Sohlberg, and M Chevignard, 2016, Proposed Criteria for  
1716 Appraising Goal Attainment Scales Used as Outcome Measures in Rehabilitation  
1717 Research, *Arch Phys Med Rehab*, 97(1):157-170.  
1718
- 1719 Lachin JM, 1999, Worst-Rank Score Analysis With Informatively Missing Observations in  
1720 Clinical Trials, *Control Clin Trials*, 20(5):408-422.  
1721
- 1722 Norman, GR, P Stratford, and G Regehr, 1997, Methodological Problems in the Retrospective  
1723 Computation of Responsiveness to Change: The Lesson of Cronbach, *J Clin Epidemiol*,  
1724 50(8):869-879.  
1725
- 1726 Rothrock, NE, AJ Kaat, MS Vrahas, RV O'Toole, SK Buono, S Morrison, and RC Gershon,  
1727 2019, Validation of PROMIS Physical Function Instruments in Patients With an  
1728 Orthopaedic Trauma to a Lower Extremity, *J Orthop Trauma*, 33(8):377-383.

## ***Contains Nonbinding Recommendations***

*Draft — Not for Implementation*

- 1729  
1730 Senn, S, 2007, *Statistical Issues in Drug Development* (2<sup>nd</sup> Edition), Chichester (England): John  
1731 Wiley & Sons.  
1732
- 1733 Shadish, WR, TD Cook, and DT Campbell DT, 2002, *Experimental and Quasi-Experimental*  
1734 *Designs for Generalized Causal Inference*, Accessed August 2, 2021,  
1735 <https://www.alnap.org/system/files/content/resource/files/main/147.pdf>.  
1736
- 1737 Shiffman, S, AA Stone, and MR Hufford, 2008, *Ecological Momentary Assessment*, *Annu Rev*  
1738 *Clin Psychol*, 4(1):1-32.  
1739
- 1740 Song, MK, and SE Ward, 2015, *Assessment Effects in Educational and Psychosocial*  
1741 *Intervention Trials: An Important But Often-Overlooked Problem*, *Res Nurs Health*,  
1742 38(3):241-247.  
1743
- 1744 Staunton, H, T Willgoss, L Nelsen, C Burbridge, K Sully, D Rofail, and R Arbuckle, 2019, *An*  
1745 *Overview of Using Qualitative Techniques to Explore and Define Estimates of Clinically*  
1746 *Important Change on Clinical Outcome Assessments*, *J Patient Rep Outcomes*, 3(1):16.  
1747
- 1748 Stone, AA, JE Broderick JE, Goldman RE, Junghaenel DU, A Bolton, M May, and S Schneider,  
1749 2021, *I. Indices of Pain Intensity Derived from Ecological Momentary Assessments:*  
1750 *Rationale and Stakeholder Preferences*, *J Pain*, 22(4):359-370.  
1751
- 1752 Terluin B, Roos EM, Terwee CB, Thorlund JB, Ingelsrud LH, 2021, *Assessing Baseline*  
1753 *Dependency of Anchor-Based Minimal Important Change (MIC): Don't Stratify on the*  
1754 *Baseline Score!* *Qual Life Res*, 30(10):2773-2782, doi:10.1007/s11136-021-02886-2.  
1755
- 1756 Uryniak, T, ISF Chan, VV Fedorov, Q Jiang, L Oppenheimer, SM Snappin, CH Teng, and J  
1757 Zhang, 2011, *Responder Analyses—A PhRMA Position Paper*, *Statistics in*  
1758 *Biopharmaceutical Research*, 3(3):476-487.  
1759
- 1760 Wilson, IB, and PD Cleary, 1995, *Linking Clinical Variables With Health-Related Quality of*  
1761 *Life*, *JAMA*, 273(1):59-65.  
1762
- 1763 Vickers AJ, 2001, *The Use of Percentage Change From Baseline as an Outcome in a Controlled*  
1764 *Trial is Statistically Inefficient: A Simulation Study*, *BMC Med Res Methodol*, 1(1):6-6.  
1765  
1766