



TO: NQF Members and Public  
FR: NQF Staff  
RE: Review of *Composite Performance Measure Evaluation Guidance*  
DA: November 29, 2012

### **Background**

Composite performance measures, which combine information on multiple individual measures into one single measure, are of increasing interest in healthcare performance measurement and public accountability applications. Such measures are complex and require a strong conceptual and methodological foundation with different considerations for testing and analysis.

NQF previously developed guidance to assist steering committees with their evaluation of composite performance measures as part of the NQF endorsement process. Since that time, NQF has gained experience with composite performance measures and has identified some challenges and issues with implementing the prior guidance. Additionally, NQF has updated the standard measure evaluation criteria and guidance for evidence, measure testing, and usability; thus, there also is a need to align the evaluation criteria for composite measures with the updated guidance.

The purpose of the Composite Performance Measure Evaluation Guidance Project was to review and update NQF's criteria and guidance on evaluating composite performance measures for potential NQF endorsement. NQF convened a 12-member Technical Expert Panel (TEP) to achieve the following goals of the project:

- review the existing guidance for evaluating composite performance measures;
- identify any unique considerations for evaluating composite performance within the context of NQF's updated endorsement criteria; and
- modify existing criteria and guidance and/or provide additional recommendations for evaluating composite performance measures.

### **Review and Comment**

The TEP's recommendations are included in the draft document *Composite Performance Measure Evaluation Guidance*. The draft report is posted on the NQF web site for purposes of review and comment only and is not intended to be used for voting purposes.

The recommendations include modifications to the current NQF criteria and guidance for evaluating composite performance measures. Of particular note is the elimination of the requirement that component performance measures must be NQF-endorsed or meet the criteria for endorsement. Instead, the recommended guidance indicates which subcriteria should be met for the component measures (e.g., evidence and performance gap) and which should be met for the composite measure (e.g., reliability and validity). There are two

subcriteria specific to composite performance measures that are incorporated into the standard criteria. Under Importance to Measure and Report, subcriterion 1d requires a clear and logical description of the quality construct, rationale, and how the composite measure construction is consistent with them. Under Scientific Acceptability of Measure Properties, subcriterion 2d requires analyses to demonstrate that the component measures fit the quality construct, the aggregation and weighting rules fit the quality construct and rationale, and extent and handling of missing data.

The TEP recommended that NQF develop examples of types of analyses that are appropriate for various approaches to composite measure construction. NQF staff will continue to work with the TEP to develop Appendix B as a resource for steering committees, staff, and measure developers. Suggestions are welcome.

You may post your comments and view the comments of others on the NQF website using the online submission process.

**All comments must be submitted no later than 6:00 PM ET, December 28, 2012.**

Thank you for your interest in the NQF's work. We look forward to your review and comments.

# Composite Performance Measure Evaluation Guidance

DRAFT REPORT FOR REVIEW

November 29, 2012



NATIONAL  
QUALITY FORUM

# Contents

- Introduction ..... 1
  - Purpose ..... 1
- Background ..... 2
  - Prior Guidance on Evaluating Composite Measures..... 2
  - NQF Experience with Composite Performance Measures ..... 2
- Definition of Composite Performance Measure ..... 4
  - Types of Composite Measures..... 4
  - Key Steps in Developing a Composite Performance Measure..... 4
- Guiding Principles ..... 5
  - Terminology ..... 5
  - Component Performance Measures..... 5
  - Composite Performance Measure ..... 6
- Recommendations for Composite Measure Evaluation ..... 6
  - Importance to Measure and Report ..... 7
  - Scientific Acceptability of Measure Properties..... 8
  - Feasibility ..... 10
  - Usability and Use ..... 10
  - Comparison to Related and Competing Measures ..... 10
- Recommendations for Review Process..... 17
- Notes..... 19
- Appendix A: Glossary ..... 20
- Appendix B: Approaches for Constructing Composite Performance Measures..... 23
- Appendix C: References Consulted ..... 28
- Appendix D: Technical Expert Panel and NQF Staff ..... 31

# 1 Composite Performance Measure Evaluation Guidance

2 DRAFT REPORT

## 3 Introduction

4 Healthcare is a complex and multidimensional activity. While individual performance measures provide  
5 much important information, there is also value in summarizing performance on multiple dimensions.  
6 Composite performance measures, which combine information on multiple individual measures into  
7 one single measure, are of increasing interest in healthcare performance measurement and public  
8 accountability applications. According to the Institute of Medicine,<sup>1</sup> such measures can enhance the  
9 performance measurement enterprise and provide a potentially deeper view of the reliability of the care  
10 system. Further, composite performance measures may be useful for multiple stakeholders, including  
11 consumers, purchasers, and policy makers.

12 Composite performance measures are complex and require a strong conceptual and methodological  
13 foundation with different considerations for testing and analysis. As with individual performance  
14 measures, the methods used to construct composite performance affects the reliability, validity, and  
15 usefulness of the composite measure.

16 Several composite measures are included in NQF's portfolio of endorsed measures, and NQF previously  
17 developed guidance<sup>2</sup> to assist steering committees with their assessment of these measures as part of  
18 the NQF evaluation process. Since that time, however, NQF has updated the standard measure  
19 evaluation criteria and guidance for evidence, measure testing, and usability; thus, there is a need to  
20 align the evaluation criteria for composite measures with the updated guidance.

## 21 Purpose

22 The purpose of the Composite Performance Measure Evaluation Guidance Project was to review and  
23 update NQF's criteria and guidance on evaluating composite performance measures for potential NQF  
24 endorsement. Specifically, the goals of the project were to:

- 25 • review the existing guidance for evaluating composite performance measures;
- 26 • identify any unique considerations for evaluating composite performance within the context of  
27 NQF's updated endorsement criteria;
- 28 • modify existing criteria and guidance and/or provide additional recommendations for evaluating  
29 composite performance measures.

30 To achieve these goals, NQF convened a 12-member Technical Expert Panel (TEP), which was comprised  
31 primarily of methodologists and other experts in the development of composite performance measures.  
32 In addition to participating in several conference calls, the TEP also gathered for a one-day in-person  
33 meeting in Washington, DC on November 2, 2012.

34

## 1 Background

### 2 Prior Guidance on Evaluating Composite Measures

3 In 2008-2009, NQF initiated a project to identify a framework for evaluating composite performance  
4 measures. That developmental work included defining composite performance measures, articulating  
5 principles underlying the evaluation of composite performance measures, and developing an initial set  
6 of specific criteria (to be used in addition to NQF's standard evaluation criteria) with which to evaluate  
7 composite performance measures for potential NQF endorsement.

8 The principles articulated for evaluating composite performance measures reflected the need for a  
9 conceptual construct of quality underlying the composite measure and justification of the methods used  
10 to construct and test the measure for reliability and validity. The criteria emphasized the need for  
11 transparency around the methodology used for composite measure construction and required that both  
12 the components of the composite and the composite measure as a whole meet NQF's measure  
13 evaluation criteria. This work served as the basis for the current project.

### 14 NQF Experience with Composite Performance Measures

15 Since 2007, 28 measures submitted to NQF for potential endorsement have been flagged as composite  
16 measures. Of these, 22 are currently endorsed. The majority of the endorsed composite measures  
17 (n=11) are derived from surveys targeted towards patients or consumers (e.g., the Consumer  
18 Assessment of Healthcare Providers and Systems (CAHPS) surveys). The remainder of the endorsed  
19 composite performance measures are comprised of all-or-none measures (n=3) and composites  
20 constructed using various methods of aggregation and weighting methodologies (n=8). As with NQF-  
21 endorsed individual measures, these composite measures are considered suitable both for performance  
22 improvement and accountability applications.

23 However, the evaluation of composite measures for potential endorsement has not been without  
24 difficulty. The most common issues have revolved around the identification of composite measures,  
25 ambiguity in the guidance when a component measure is not NQF-endorsed, and incomplete  
26 submissions.

#### 27 *Identifying Measures as Composites*

28 Not all composite measures that have met—or potentially have met—the current NQF definition of  
29 composite measures have been flagged by the measure developers as composite measures. These  
30 include all-or-none composites in which the components are assessed at the patient level (i.e., whether  
31 each patient received each required process then aggregated for the healthcare entity); simpler all-or-  
32 none measures that require multiple conditions (e.g., assess vaccination status and administer flu  
33 vaccine); and any-or-all measures that assess whether a patient has exhibited any or all of a list of  
34 complications. For such measures, it is unclear whether the additional analyses indicated for composite  
35 measures (e.g., analysis of components to demonstrate alignment with the conceptual construct and  
36 contribution to the variation in the overall composite score) are applicable, and often these additional  
37 analyses have not been submitted by developers of these measures.

1 *Evaluation of Component Measures*

2 The current guidance indicates that the component measures that make up a composite measure  
3 should be NQF-endorsed or evaluated as meeting the individual measure evaluation criteria as the first  
4 step in evaluating the composite measure. However, the guidance goes on to state that while a  
5 component measure might not be important enough in its own right as an individual measure, it could  
6 be determined to be an important component of a composite. Some developers have interpreted this  
7 guidance to mean that components do not need to meet the Importance to Measure and Report criteria  
8 around evidence, impact, and performance gap. But this interpretation regarding evidence and  
9 performance gap calls into question the basis for including the component measure. Another issue  
10 related to the evidence criterion is whether measures that are distal to desired outcomes could be  
11 included in composite measures. For example, a performance measure of merely obtaining a lab test is  
12 not considered to meet the criteria because it is so distal to the desired outcome and is often based on  
13 expert opinion; however, this type of component has been suggested for inclusion in a composite  
14 measure.

15 It also is not clear whether balancing measures that would not meet the importance criteria should be  
16 included in a composite performance measure. A balancing measure is not the main focus of interest  
17 but is used to identify or monitor potential adverse consequences of measurement. For example, a  
18 performance measure about treating substance use that requires the identification of patients with  
19 substance use problems will not be accurate if most patients are not even screened; therefore, a  
20 screening measure might be considered a balancing measure.

21 Finally, it has been difficult to apply criteria for related and competing measures to composite measures.  
22 The challenges with measure harmonization are amplified with composite measures because, typically,  
23 more measures (involving multiple developers) will be involved in harmonization discussions. While  
24 using previously-endorsed measures as components in a composite measure should ameliorate most  
25 difficulties around harmonization, often the components in submitted composite measures have not  
26 been previously endorsed. In such cases, these components either competed directly with other  
27 endorsed measures or were not harmonized with endorsed measures.

28 *Incomplete Submissions Related to Requirements for Composite Measures*

29 As discussed earlier, if measures are not flagged as composite measures, then the additional information  
30 needed to evaluate them as composite measures may not be submitted by measure developers.  
31 However, non-responsiveness to composite-specific items also has been a problem. For example, the  
32 current criteria state that the purpose/objective of the composite measure and the construct for quality  
33 must be clearly described, yet often little beyond a list of the component measures is provided.

34 Current criteria require testing for reliability and validity of the composite measure (even if the  
35 individual measures have demonstrated reliability and validity), as well as additional analyses to justify  
36 the inclusion of component measures and the specified aggregation and weighting rules. Reliability and  
37 validity testing of the composite measure may not be conducted. Some of the composite questions refer  
38 to correlational analyses, which may not be appropriate for all composite measures. While the current  
39 guidance recognizes this and indicates that the developer could submit other analyses with rationale,  
40 these alternative analyses have not always been submitted (or if submitted, the rationale may not have  
41 been included or may not have been sufficiently explanatory). Analysis of the contribution of individual  
42 components to the composite score often has not been submitted. Without this information, steering

1 committees may be left with little more than face validity as a basis for recommending a composite  
2 performance measure.

### 3 **Definition of Composite Performance Measure**

4 *A composite performance measure is a combination of two or more individual performance measures in*  
5 *a single performance measure that results in a single score.*

6 The TEP reviewed and retained the definition provided in the initial composite report and added explicit  
7 clarification that it refers to composite *performance* measures. Note that the term “composite  
8 measure” also refers to individual-level measures (i.e., instruments and scales used to obtain data from  
9 individuals, such as the CAHPS or PHQ-9). Data from such instruments may be used in performance  
10 measures that aggregate data for all patients served by a healthcare entity, but that does not in itself  
11 make it a composite performance measure. Patient-reported outcomes (PRO) and performance  
12 measurement has been the subject of a recent NQF project (see [PRO report](#)). Throughout this report,  
13 the terms composite measure or component measure refer to a performance measure unless otherwise  
14 indicated.

### 15 **Types of Composite Measures**

16 Composites often are classified according to the empirical and conceptual relationship among the  
17 component measures and between the components and the composite, the rules for combining the  
18 individual components (e.g., all-or-none, opportunity, weighted average), or the type of individual  
19 measures included in the composite (e.g., process, outcome). Regardless of the various approaches, a  
20 coherent quality construct and rationale should guide the composite development and testing and  
21 analysis. The glossary in [Appendix A](#) contains definitions for various approaches to combining the  
22 component measures. [Appendix B](#) provides a description of various conceptual models for the  
23 relationship among the component measures and composite score and identifies relevant analyses. The  
24 TEP suggested that over time, NQF add specific examples of composite performance measures that use  
25 these conceptual models.

26 The TEP decided that a classification of types of composite performance measures would not be  
27 particularly useful and could lead to unnecessary attention to the approach used to construct the  
28 composite. They agreed that the primary concern for NQF endorsement is whether the resulting  
29 composite performance measure is based on sound measurement science, produces a reliable signal,  
30 and is a valid reflection of quality.

### 31 **Key Steps in Developing a Composite Performance Measure**

32 A variety of methods can be used to construct composite performance measures; however, they all  
33 involve the following key steps:<sup>3-9</sup>

- 34 • Describing the quality construct to be measured and rationale for the composite;
- 35 • Selecting the component measures to be combined in the composite measure;
- 36 • Ensuring that the methods used to aggregate and weight the components supports the goal that  
37 is articulated for the measure;
- 38 • Combining the component measure scores, using the specified method; and



- 1       • Testing the composite measure to determine if it is a reliable and valid indicator of quality  
2       healthcare.

### 3   **Guiding Principles**

4   The following key principles were identified by the TEP and guided their recommendations regarding the  
5   evaluation criteria.

### 6   **Terminology**

7   As noted above, the TEP opted for a broad, generic definition of composite performance measure.

- 8       • The term “composite measure” may be applied to many types of measures, including individual-  
9       level instruments as well as aggregate-level performance measures. NQF only endorses  
10      performance measures.  
11      • Approaches to composite measure development and construction are described using a variety  
12      of terms and can vary by discipline. Nonetheless, the construction and evaluation of composite  
13      performance measures should be based on sound measurement science principles. Although  
14      often used in the published literature on composite measures, the TEP wanted to minimize the  
15      use of discipline-specific language and categorizations (e.g., “psychometric” and “clinimetric”) in  
16      the evaluation criteria and guidance.

### 17   **Component Performance Measures**

18   The prior composite evaluation criteria required that each component performance measure be NQF-  
19   endorsed or meet all criteria for NQF endorsement. At times that has been difficult to implement,  
20   particularly for reliability. The TEP noted that individual measures may not be reliable independently  
21   because of rare events or small case volume, but could be used successfully within a composite because  
22   the composite combines multiple measures, which can increase reliability of the composite performance  
23   measure as a whole. Rather than requiring that each component meet all NQF criteria, the TEP focused  
24   on the overall composite and identified those NQF criteria that must be met to justify inclusion of the  
25   individual component measures. The TEP agreed, however, that if an individual component measure is  
26   NQF-endorsed, then those criteria would not need to be demonstrated again.

- 27       • NQF-endorsement of the individual component measures should not be mandatory; however,  
28       NQF endorsement of the component measures could satisfy some requirements for the  
29       component measures included in a composite.  
30       • The individual component measures that are included in a composite performance measure  
31       should be justified based on the clinical evidence (i.e., for process measures, what is being  
32       measured is based on clinical evidence of a link to desired outcomes; for health outcomes, a  
33       rationale that it is influenced by healthcare).  
34       • The individual components in a composite performance measure generally should demonstrate  
35       a gap in performance; however, there may be conceptual or analytical justification for including  
36       components that do not have a gap in performance.  
37       • The individual components may not be sufficiently reliable independently, but could contribute  
38       to the reliability of the composite performance measure.

## 1 Composite Performance Measure

2 The TEP emphasized the need for a coherent quality construct and rationale to guide construction of the  
3 composite as well as to guide evaluation for NQF endorsement. A quality construct is a hypothetical  
4 complex concept of quality. Component measures should be selected based on fit with the quality  
5 construct, and analyses should justify that fit. All composite performance measures share the potential  
6 for simplification when representing one score versus many scores for individual performance  
7 measures. However, that feature alone is not sufficient justification for a composite performance  
8 measure. Each component should fit the construct and be necessary. The composite performance  
9 measure should provide added value over having individual performance measures. Composite  
10 measures are complex with aggregation and weighting rules that are not applicable to the individual  
11 component measures; therefore, reliability and validity of the composite performance measure score  
12 should be demonstrated.

- 13 • A coherent quality construct and rationale for the composite performance measure are essential  
14 for determining:
  - 15 ○ what components are included in a composite performance measure;
  - 16 ○ how the components are aggregated and weighted;
  - 17 ○ what analyses should be used to support the components and demonstrate reliability  
18 and validity; and
  - 19 ○ added value over that of individual measures alone.
- 20 • Reliability and validity of the individual components do not automatically ensure reliability and  
21 validity of the constructed composite performance measure. Reliability and validity of the  
22 constructed composite performance measure should be demonstrated.
- 23 • When evaluating composite performance measures, both the quality construct itself, as well the  
24 empirical evidence for the composite (i.e., supporting the method of construction and methods  
25 of analysis), should be considered.
- 26 • Components of a composite performance measure should be “necessary”—either empirically  
27 (i.e., they contribute to the reliability) or conceptually. A secondary objective is parsimony,  
28 when possible.
- 29 • The individual components in a composite performance measure may or may not be correlated,  
30 depending on the quality construct.
- 31 • Aggregation and weighting rules for constructing composite performance measures should be  
32 consistent with the quality construct and rationale for the composite. A secondary objective is  
33 simplicity, when possible.
- 34 • The standard NQF criteria apply to composite performance measures.

## 35 Recommendations for Composite Measure Evaluation

36 The NQF measure evaluation criteria apply to composite performance measures and their component  
37 performance measures. The goal is to incorporate evaluation of composite performance measures into  
38 the standard NQF criteria and processes to the extent possible. NQF endorsement is not required for  
39 the component measures unless they are intended to be used independently to make judgments about  
40 performance. However, the individual component measures should meet specific subcriteria such as for  
41 clinical evidence and performance gap, although there may be potential exceptions. The TEP agreed that  
42 two additional criteria are needed to evaluate composite performance measures; these are  
43 incorporated into the evaluation criteria in Table 1 (see [1d](#) and [2d](#)).

1 It is important to note the difference between the NQF criteria for evidence and validity. The evidence  
2 subcriterion is included under the Importance to Measure and Report criterion and addresses the  
3 empirical clinical evidence linking processes to desired health outcomes. In contrast, the validity  
4 subcriterion is included under the Scientific Acceptability of Measure Properties criterion and addresses  
5 whether the performance measure *as constructed* is an accurate reflection of quality. The clinical  
6 evidence provides a justification for measurement and a foundation for validity, but the actual  
7 performance measure should be empirically tested to demonstrate validity because how a measure is  
8 constructed can affect whether it is an accurate reflection of quality.

## 9 Importance to Measure and Report

### 10 *Evidence*

11 Each component measure must meet the evidence criterion to justify its inclusion in the composite. As  
12 with individual performance measures, the evidence requirement ensures that efforts for measurement  
13 are devoted to health outcomes or processes of care that will influence desired outcomes. If a  
14 component measure is NQF-endorsed (since the updated evidence requirements were implemented), it  
15 could be considered as meeting the evidence criterion. If all component measures do not meet the  
16 evidence criterion, or do not qualify for the exceptions to the evidence criterion, the composite would  
17 not meet the criterion for Importance to Measure and Report unless those components were removed.  
18 Evidence is required regardless of approach to constructing a composite measure (i.e., all-or-none  
19 scoring or combining scores from individual performance measures).

### 20 *Performance Gap*

21 Each component measure also should meet the criterion of performance gap to justify its inclusion in  
22 the composite. As with individual performance measures, effort for measurement should be directed to  
23 aspects of care where there is variation or overall poor performance. However, the TEP acknowledged  
24 there may be circumstances when a component measure that does not meet the performance gap  
25 criterion could be included in a composite. In such cases, justification for including such a component  
26 would be required (e.g., it contributes to the reliability of the overall composite score or is needed for  
27 face validity). Ideally, the composite performance measure as a whole also should demonstrate a  
28 performance gap.

### 29 *Quality Construct and Rationale of a Composite Performance Measure*

30 A subcriterion specific for composite performance measures is included under Importance to Measure  
31 and Report (see [1d](#) in Table 1). This is consistent with and refines the prior guidance regarding  
32 describing the purpose and quality construct for the composite performance measure. Composite  
33 measures are complex and represent a higher order construct than the individual measures. Justification  
34 for the approach to composite measure construction and analysis stems from the quality construct and  
35 rationale. Therefore, the quality construct should be clearly articulated and logical in order to meet this  
36 subcriterion and the must-pass criterion of Importance to Measure and Report.

37 Although the TEP recommended that the rationale for the composite performance measure be  
38 identified, it acknowledged that NQF endorses performance measures intended for both accountability  
39 and performance improvement and does not endorse measures for a specific accountability application  
40 (e.g., payment vs. public reporting). One TEP member suggested that the rationale should include the  
41 intended decision-making context (e.g., to select a provider for services, select a provider for contracting

1 or payment incentives, to identify or direct resources for improvement). While others noted that it  
2 might be difficult to envision how or why the component measures or the composite construction  
3 methodology should differ despite the decision-making context (given that all the decisions involve  
4 distinguishing good from poor quality), they did agree that measure developers should clearly explain  
5 how the aggregation and weighting of the components are consistent with the stated quality construct  
6 and rationale for the composite measure, including any decision-making context. The TEP  
7 acknowledged that endorsing multiple composite measures for a quality construct (such as quality of  
8 care for patients with congestive heart failure) that are created for different decision-making  
9 motivations could increase confusion and issues with competing measures. The decision-making  
10 context may influence whether a composite measure is useful at all. For example, a composite  
11 performance measure that includes multiple surgical mortality measures may be useful for assessing  
12 overall surgical quality, whereas the individual performance measures are more useful for selecting a  
13 hospital for a specific surgical procedure.

## 14 Scientific Acceptability of Measure Properties

### 15 *Reliability*

16 One cited advantage of composite performance measures is that using multiple indicators (components)  
17 increases reliability (i.e., the ability to detect a provider effect).<sup>5</sup> The purpose of combining individual  
18 measures that assess the quality of care provided to patients by providers or institutions is to determine  
19 whether these measures are useful in detecting a consistent pattern of practice or quality of care across  
20 patients of the provider or institution. That is, does a set of measures that, taken together, are thought  
21 to reflect good quality of care, show a more consistent pattern within a provider's practice or within an  
22 institution, and greater differences between providers or institutions than would be expected by chance  
23 alone? Reliability testing of the composite performance measure should demonstrate that the  
24 composite measure score differentiates signal from noise (i.e., random measurement error). It should be  
25 noted that increased reliability with increased number of indicators does not hold for all-or-none  
26 measures, when multiple indicators are essentially reduced to one data point.<sup>8</sup> Nevertheless, all-or-none  
27 performance measures also should demonstrate signal-to-noise reliability. These examples of measure  
28 reliability are different than a rationale that all-or-none measures are intended to foster a system  
29 perspective of care, sometimes called "system reliability."

30 Although ideal, demonstrated signal-to-noise reliability of the individual component measures is not  
31 essential for having a reliable composite measure. In some cases, an individual performance measure  
32 may not provide a reliable signal because of small volume or rare events. However, that measure could  
33 appropriately be used as a component in a reliable composite performance measure.

### 34 *Validity*

35 Validity testing of the constructed composite performance measure score is more important than  
36 validity testing of the component measures because even if the individual component measures are  
37 valid, the aggregation and weighting rules for constructing the composite could result in a score that is  
38 not an accurate reflection of quality. However, some TEP members thought that requiring validity  
39 testing of the composite as a whole would be difficult to accomplish prior to NQF endorsement,  
40 although others questioned why NQF would endorse a performance measure without empirical  
41 evidence of validity. If validity of the composite performance measure is not demonstrated, then each of

1 the individual component measures must meet the NQF criteria for validity; further, validity testing of  
2 the overall composite measure would be expected by the time of endorsement maintenance.

3 It may be unlikely that another valid measure of the same quality construct (i.e., a criterion measure)  
4 will be available to test the criterion validity of a composite performance measure. Therefore, validity  
5 testing will require understanding and testing of various theoretical relationships. For example, a  
6 composite measure that includes multiple process measures could be tested for its association with a  
7 measure of a desired outcome. Alternatively, a composite measure might be tested for its ability to  
8 predict future outcomes or its ability to differentiate performance between groups known to differ on  
9 the particular quality construct.

### 10 *Additional Testing of the Composite Performance Measure*

11 A subcriterion specific for composite performance measures is included under Scientific Acceptability of  
12 Measure Properties (see [2d](#) in Table 1). This is consistent with and refines the prior guidance regarding  
13 additional analyses to justify the construction of the composite measure (both component selection and  
14 aggregation and weighting rules). The initial criteria for testing were more relevant to composite  
15 measures that are based on correlated components. The modified criterion is neutral in terms of the  
16 analyses required. For example, if the rationale for summarizing the component measures in a  
17 composite is based on their correlation with each other, then analyses based on correlation (e.g., factor  
18 analysis, item-to-total correlation, and inter-item correlation) are appropriate. In such cases, very high  
19 correlations between component measures may suggest that a component is redundant and not  
20 necessary. Conversely, if the rationale for summarizing the measures in a composite is not based on their  
21 correlation with each other, then analyses demonstrating the contribution of each component to the  
22 composite score, or their clinical justification (e.g., correlation of the individual component measures to  
23 a common outcome measure) are indicated.

24 The unit of analysis for which performance measures are calculated is typically the provider or  
25 institution (hospital, clinic, etc.) rather than the individual patient. For such measures, correlational  
26 analysis such as factor analysis or internal consistency reliability should be calculated at the level of the  
27 unit rather than patient, because the unit scores are what will be reported and acted upon. Correlations  
28 at the unit level might be quite different from those at the patient level. For example, in a patient  
29 survey, some respondents might tend to give more positive (or more negative) responses across the  
30 board, creating positive correlations among items that measure entirely distinct aspects of quality.  
31 However, when data are aggregated to the provider level, these patient tendencies average out,  
32 revealing correlations among items related at the provider level. As another example, measures of  
33 cardiac surgery might include complication rates during CABG surgery, during valve repair surgery, and  
34 during valve replacement surgery; since typically any patient undergoes only one of these procedures,  
35 the patient-level correlations of these measures are not defined but correlations at the provider or  
36 hospital level are meaningful and could be examined to assess the validity of a composite surgical  
37 quality measure. However, special statistical methods should be used for estimating such unit-level  
38 correlations, especially when the component measures do not have high unit-level reliability.

39 Composite measures are, by definition, complex; however, secondary objectives for composite measure  
40 construction are parsimony regarding the component selection and simplicity in terms of the  
41 aggregation and weighting. Scientific Acceptability of Measure Properties is a must-pass criterion and

1 measures must meet both reliability and validity. In addition, composite measures must meet this  
 2 additional criterion in order to meet the must-pass criterion of Scientific Acceptability.

3 **Feasibility**

4 The standard feasibility criteria apply to the composite measure as a whole, but must take into account  
 5 all the component measures. That is, feasibility of the composite measure will be influenced by the least  
 6 feasible of the component measures.

7 **Usability and Use**

8 Composite performance measures must meet the updated criteria for Usability and Use. The TEP noted  
 9 that disaggregation of a composite measure is not an absolute requirement because the individual  
 10 component measures need not be independently reliable. However, at a minimum, the components of  
 11 the composite performance measure must be identified. For purposes of improvement, the data must  
 12 be collected so as to facilitate investigation of the individual components.

13 **Comparison to Related and Competing Measures**

14 Composite performance measures are subject to comparison to related and competing measures. If the  
 15 component measures are not NQF-endorsed, they must be harmonized with endorsed measures or  
 16 assessed against competing measures.

17 Table 1. NQF Measure evaluation Criteria and Guidance for Evaluating Composite Performance  
 18 Measures

Measure Evaluation Criteria	Guidance for Composite Performance Measures
<b>Conditions</b>	
<p><b>1. Evidence, Performance Gap, and Priority—Importance to Measure and Report:</b> Extent to which the specific measure focus is evidence-based, important to making significant gains in healthcare quality, and improving health outcomes for a specific high-priority aspect of healthcare where there is variation in or overall less-than-optimal performance. <i>Measures must be judged to meet all subcriteria to pass this criterion and be evaluated against the remaining criteria.</i></p> <p><b>1a. Evidence to Support the Measure Focus</b>            The measure focus is evidence-based, demonstrated as follows:</p> <ul style="list-style-type: none"> <li>• <b>Health outcome:</b> <sup>3</sup> a rationale supports the relationship of the health outcome to processes or structures of care.</li> <li>• <b>Intermediate clinical outcome:</b> a systematic assessment and grading of the quantity, quality, and consistency of the body of evidence <sup>4</sup> that the measured intermediate clinical outcome leads to a desired health outcome.</li> <li>• <b>Process:</b> <sup>5</sup> a systematic assessment and grading of the quantity, quality, and consistency of the body of evidence <sup>4</sup> that the measured process leads to a desired health outcome.</li> </ul>	<p>The evidence criterion (1a) must be met for each component of the composite (unless NQF-endorsed under the current evidence requirements).</p>



Measure Evaluation Criteria	Guidance for Composite Performance Measures
<ul style="list-style-type: none"> <li>• <b>Structure:</b> a systematic assessment and grading of the quantity, quality, and consistency of the body of evidence <sup>4</sup> that the measured structure leads to a desired health outcome.</li> <li>• <b>Experience with care:</b> evidence that the measured aspects of care are those valued by patients and for which the patient is the best and/or only source of information OR that patient experience with care is correlated with desired outcomes.</li> <li>• <b>Efficiency:</b> <sup>6</sup> evidence not required for the resource use component.</li> </ul> <p><b>AND</b></p> <p><b>1b. Performance Gap</b>  Demonstration of quality problems and opportunity for improvement, i.e., data <sup>7</sup> demonstrating</p> <ul style="list-style-type: none"> <li>• considerable variation, or overall less-than-optimal performance, in the quality of care across providers; <b>and/or</b></li> <li>• disparities in care across population groups.</li> </ul> <p><b>AND</b></p> <p><b>1c. High Priority</b>  The performance measure addresses:</p> <ul style="list-style-type: none"> <li>• a specific national health goal/priority identified by DHHS or the National Priorities Partnership convened by NQF;</li> </ul> <p><b>OR</b></p> <ul style="list-style-type: none"> <li>• a demonstrated high-priority aspect of healthcare (e.g., affects large numbers of patients and/or has a substantial impact for a smaller population; leading cause of morbidity/mortality; high resource use (current and/or future); severity of illness; and severity of patient/societal consequences of poor quality).</li> </ul> <p><b>Composite 1d. For composite performance measures, the following must be clearly stated and logical:</b></p> <ul style="list-style-type: none"> <li>• The quality construct, including the representativeness of the component measures and the relationship of the component measures to the composite and to each other; and</li> <li>• The rationale for constructing a composite measure, including how the composite provides a distinctive or additive value over the component measures individually; and</li> <li>• How the aggregation and weighting of the component measures are consistent with and representative of the stated quality construct and rationale.</li> </ul>	<p>The performance gap criterion (1b) must be met for each component (unless NQF-endorsed), and if possible, for the composite performance measure as a whole. If a component measure has little opportunity for improvement, justification for why it should be included in the composite is required (e.g., increase reliability of the composite, clinical evidence).</p> <p>The priority criterion (1c) applies to the composite performance measure as a whole.</p> <p>Subcriterion 1d must be met for a composite performance measure to meet the criterion of Importance to Measure and Report. All three elements must be clearly articulated and represent a logical rationale.</p>
<p><b>2. Reliability and Validity—Scientific Acceptability of Measure Properties:</b> Extent to which the measure, <u>as specified</u>, produces consistent (reliable) and credible (valid) results about the quality of care when implemented. <b>Measures must be judged to meet the subcriteria for both reliability and validity to pass this criterion and be evaluated</b></p>	

Measure Evaluation Criteria	Guidance for Composite Performance Measures
<p><i>against the remaining criteria.</i></p> <p><b>2a. Reliability</b></p> <p><b>2a1.</b> The measure is well defined and precisely specified <sup>8</sup> so it can be implemented consistently within and across organizations and allows for comparability. EHR measure specifications are based on the quality data model (QDM). <sup>9</sup></p> <p><b>Add to Note 8:</b> Composite measure specifications include scoring rules (i.e., how the component scores are combined or aggregated), how missing data are handled, required sample sizes; and when appropriate methods for standardizing scales across component scores and weighting rules (i.e., whether all component scores are given equal or differential weighting when combined into the composite)</p> <p><b>2a2.</b> Reliability testing <sup>10</sup> demonstrates the measure data elements are repeatable, producing the same results a high proportion of the time when assessed in the same population in the same time period and/or that the measure score is precise.</p> <p><b>2b. Validity</b></p> <p><b>2b1.</b> The measure specifications <sup>8</sup> are consistent with the evidence presented to support the focus of measurement under criterion 1c. The measure is specified to capture the most inclusive target population indicated by the evidence, and exclusions are supported by the evidence.</p> <p><b>2b2.</b> Validity testing <sup>11</sup> demonstrates that the measure data elements are correct and/or the measure score correctly reflects the quality of care provided, adequately identifying differences in quality.</p> <p><b>2b3.</b> Exclusions are supported by the clinical evidence; otherwise, they are supported by evidence of sufficient frequency of occurrence so that</p>	<p><b>2a2.</b> For composite performance measures, reliability must be demonstrated for the composite measure score. Reliability of the individual component measures is not sufficient, and in some cases, component measures that are not independently reliable can contribute to reliability of the composite measure. However, if the component measures will be disaggregated in accountability applications, then reliability for the component measures must be demonstrated (unless NQF-endorsed).</p> <p><b>2b2.</b> For composite performance measures, validity should be demonstrated for the composite measure score. If not feasible at the time of initial endorsement, validity of the component measures must meet NQF criteria, and by endorsement maintenance, validity of the composite performance measure must be demonstrated. If the component measures will be disaggregated for accountability applications, then validity for the component measures must be demonstrated (unless NQF-endorsed).</p> <p><b>2b3.</b> Exclusions apply primarily to the component measures and would not</p>



Measure Evaluation Criteria	Guidance for Composite Performance Measures
<p>results are distorted without the exclusion; <sup>12</sup></p> <p><b>AND</b> If patient preference (e.g., informed decisionmaking) is a basis for exclusion, there must be evidence that the exclusion impacts performance on the measure; in such cases, the measure must be specified so that the information about patient preference and the effect on the measure is transparent (e.g., numerator category computed separately, denominator exclusion category computed separately). <sup>13</sup></p> <p><b>2b4.</b> For outcome measures and other measures when indicated (e.g., resource use):</p> <ul style="list-style-type: none"> <li>• an evidence-based risk-adjustment strategy (e.g., risk models, risk stratification) is specified; is based on patient factors that influence the measured outcome (but not factors related to disparities in care or the quality of care) and are present at start of care; <sup>14,15</sup> and has demonstrated adequate discrimination and calibration</li> </ul> <p><b>OR</b></p> <ul style="list-style-type: none"> <li>• rationale/data support no risk adjustment/ stratification.</li> </ul> <p><b>2b5.</b> Data analysis of computed measure scores demonstrates that methods for scoring and analysis of the specified measure allow for identification of statistically significant and practically/clinically meaningful <sup>16</sup> differences in performance;</p> <p><b>OR</b> there is evidence of overall less-than-optimal performance.</p> <p><b>2b6.</b> If multiple data sources/methods are specified, there is demonstration they produce comparable results.</p> <p><b>2c. Disparities</b> If disparities in care have been identified, measure specifications, scoring, and analysis allow for identification of disparities through stratification of results (e.g., by race, ethnicity, socioeconomic status, gender);</p> <p><b>OR</b> rationale/data justifies why stratification is not necessary or not feasible.</p> <p><b>Composite 2d. For composite performance measures, empirical analyses demonstrate:</b></p> <ul style="list-style-type: none"> <li>• the component measures fit the quality construct and are necessary (secondary objective of parsimony to the extent possible);</li> <li>• the aggregation and weighting rules are consistent with the quality construct and rationale (secondary objective of simplicity to the extent possible); and</li> <li>• the extent of missing data and how the specified handling of missing data minimizes bias.</li> </ul>	<p>need to be addressed if validity of the composite performance measure was demonstrated.</p> <p>2b4. This would be required for outcome component measures (unless NQF-endorsed).</p> <p>2b5. Applies to composite performance measures.</p> <p>2b6. Applies to component measures.</p> <p>2c. Applies to composite performance measures.</p> <p>Subcriterion 2d must be met for a composite performance measure to meet the criterion of Scientific Acceptability of Measure Properties.</p>
<p><b>3. Feasibility:</b> Extent to which the required data are readily available or</p>	

Measure Evaluation Criteria	Guidance for Composite Performance Measures
<p>could be captured without undue burden and can be implemented for performance measurement.</p> <p><b>3a.</b> For clinical measures, the required data elements are routinely generated and used during care delivery (e.g., blood pressure, lab test, diagnosis, medication order).</p> <p><b>3b.</b> The required data elements are available in electronic health records or other electronic sources. If the required data are not in electronic health records or existing electronic sources, a credible, near-term path to electronic collection is specified.</p> <p><b>3c.</b> Demonstration that the data collection strategy (e.g., source, timing, frequency, sampling, patient confidentiality, <sup>17</sup> costs associated with fees/licensing of proprietary measures) can be implemented (e.g., already in operational use, or testing demonstrates that it is ready to put into operational use).</p>	<p>3a, 3b, 3c. Apply to composite performance measures as a whole, taking into account all component measures.</p>
<p><b>4. Usability and Use</b> Extent to which potential audiences (e.g., consumers, purchasers, providers, policymakers) are using or could use performance results for both accountability and performance improvement <sup>18</sup> to achieve the goal of high-quality, efficient healthcare for individuals or populations.</p> <p><b>4a. Accountability and Transparency</b> <sup>19</sup> Performance results are used in at least one accountability application <sup>1</sup> within three years after initial endorsement and are publicly reported <sup>19</sup> within six years after initial endorsement (or the data on performance results are available). <sup>20</sup> If not in use at the time of initial endorsement, then a credible plan <sup>21</sup> for implementation within the specified timeframes is provided.</p> <p><b>AND</b></p> <p><b>4b. Improvement</b> <sup>22</sup> Progress toward achieving the goal of high-quality, efficient healthcare for individuals or populations is demonstrated. <sup>22</sup> If not in use for performance improvement at the time of initial endorsement, then a credible rationale describes how the performance results could be used to further the goal of high-quality, efficient healthcare for individuals or populations.</p> <p><b>AND</b></p> <p><b>4c.</b> The benefits of the performance measure in facilitating progress toward achieving high-quality, efficient healthcare for individuals or populations outweigh evidence of unintended negative consequences to individuals or populations (if such evidence exists).</p>	<p>4a. Applies to composite performance measures. To facilitate transparency, at a minimum, the individual components of the composite must be identified; ideally, if the components are NQF-endorsed and/or meet reliability and validity criteria, the component scores would be reported.</p> <p>4b. Applies to composite performance measures, except it should explicitly link back to the quality construct and rationale. To facilitate improvement, data should be collected in such a way as to permit disaggregation.</p> <p>4c. Applies to composite performance measures. If there is evidence of unintended negative consequences for one of the components, the developer should explain how that is handled or justify why that component should remain in the composite.</p>

Measure Evaluation Criteria	Guidance for Composite Performance Measures
<p><b>5. Comparison to Related or Competing Measures</b> If a measure meets the above criteria <u>and</u> there are endorsed or new related measures (either the same measure focus or the same target population) or competing measures (both the same measure focus and the same target population), the measures are compared to address harmonization and/or selection of the best measure.</p> <p><b>5a.</b> The measure specifications are harmonized <sup>23</sup> with related measures; <b>OR</b> the differences in specifications are justified.</p> <p><b>5b.</b> The measure is superior to competing measures (e.g., is a more valid or efficient way to measure); <b>OR</b> multiple measures are justified.</p>	<p>5a and 5b. Applies to composite performance measures as a whole as well as the component measures.</p>

1

2 Table 2. Notes to Measure Evaluation Criteria

Conditions
<p><b>1. Accountability applications</b> are the use of performance results about identifiable, accountable entities to make judgments and decisions as a consequence of performance, such as reward, recognition, punishment, payment, or selection (e.g., public reporting, accreditation, licensure, professional certification, health information technology incentives, performance-based payment, network inclusion/exclusion). <b>Selection</b> is the use of performance results to make or affirm choices regarding providers of healthcare or health plans.</p> <p><b>2.</b> A measure that has not been tested for reliability and validity is only potentially eligible for time-limited endorsement if all of the following conditions are met: 1) the measure topic is not addressed by an endorsed measure; 2) it is relevant to a critical timeline (e.g., legislative mandate) for implementing endorsed measures; 3) the measure is not complex (requiring risk adjustment or a composite); and 4) the measure steward verifies that testing will be completed within 12 months of endorsement.</p>
<p><b>1. Evidence, Performance Gap, and Priority—Importance to Measure and Report</b></p>
<p><b>3.</b> Generally, rare event outcomes do not provide adequate information for improvement or discrimination; however, serious reportable events that are compared to zero are appropriate outcomes for public reporting and quality improvement.</p> <p><b>4.</b> The preferred systems for grading the evidence are the U.S. Preventive Services Task Force (USPSTF) <a href="#">grading definitions</a> and <a href="#">methods</a>, or Grading of Recommendations, Assessment, Development and Evaluation (<a href="#">GRADE guidelines</a>).</p> <p><b>5.</b> Clinical care processes typically include multiple steps: assess → identify problem/potential problem → choose/plan intervention (with patient input) → provide intervention → evaluate impact on health status. If the measure focus is one step in such a multistep process, the step with the strongest evidence for the link to the desired outcome should be selected as the focus of measurement.</p> <p><b>6.</b> Measures of efficiency combine the concepts of resource use <u>and</u> quality (NQF’s <a href="#">Measurement Framework: Evaluating Efficiency Across Episodes of Care</a>; <a href="#">AQA Principles of Efficiency Measures</a>).</p>

7. Examples of data on opportunity for improvement include, but are not limited to: prior studies, epidemiologic data, or data from pilot testing or implementation of the proposed measure. If data are not available, the measure focus is systematically assessed (e.g., expert panel rating) and judged to be a quality problem.

## 2. Reliability and Validity—Scientific Acceptability of Measure Properties

8. Measure specifications include the target population (denominator) to whom the measure applies, identification of those from the target population who achieved the specific measure focus (numerator, target condition, event, outcome), measurement time window, exclusions, risk adjustment/stratification, definitions, data source, code lists with descriptors, sampling, scoring/computation. **Composite measure specifications include scoring rules (i.e., how the component scores are combined or aggregated), how missing data are handled, required sample sizes; and when appropriate methods for standardizing scales across component scores and weighting rules (i.e., whether all component scores are given equal or differential weighting when combined into the composite).**

9. EHR measure specifications include data type from the QDM, code lists, EHR field, measure logic, original source of the data, recorder, and setting.

10. Reliability testing applies to both the data elements and computed measure score. Examples of reliability testing for data elements include, but are not limited to: inter-rater/abstractor or intra-rater/abstractor studies; internal consistency for multi-item scales; test-retest for survey items. Reliability testing of the measure score addresses precision of measurement (e.g., signal-to-noise).

11. Validity testing applies to both the data elements and computed measure score. Validity testing of data elements typically analyzes agreement with another authoritative source of the same information. Examples of validity testing of the measure score include, but are not limited to: testing hypotheses that the measure scores indicate quality of care, e.g., measure scores are different for groups known to have differences in quality assessed by another valid quality measure or method; correlation of measure scores with another valid indicator of quality for the specific topic; or relationship to conceptually related measures (e.g., scores on process measures to scores on outcome measures). Face validity of the measure score as a quality indicator may be adequate if accomplished through a systematic and transparent process, by identified experts, and explicitly addresses whether performance scores resulting from the measure as specified can be used to distinguish good from poor quality.

12. Examples of evidence that an exclusion distorts measure results include, but are not limited to: frequency of occurrence, variability of exclusions across providers, and sensitivity analyses with and without the exclusion.

13. Patient preference is not a clinical exception to eligibility and can be influenced by provider interventions.

14. Risk factors that influence outcomes should not be specified as exclusions.

15. Risk models should not obscure disparities in care for populations by including factors that are associated with differences/inequalities in care, such as race, socioeconomic status, or gender (e.g., poorer treatment outcomes of African American men with prostate cancer or inequalities in treatment for CVD risk factors between men and women). It is preferable to stratify measures by race and socioeconomic status rather than to adjust out the differences.

16. With large enough sample sizes, small differences that are statistically significant may or may not be practically or clinically meaningful. The substantive question may be, for example, whether a statistically significant difference of one percentage point in the percentage of patients who received smoking cessation counseling (e.g., 74 percent v. 75 percent) is clinically meaningful; or whether a statistically significant difference of \$25 in cost for an episode of care (e.g., \$5,000 v. \$5,025) is practically meaningful. Measures with overall less-than-optimal performance may not demonstrate much variability across providers.

## Feasibility

17. All data collection must conform to laws regarding protected health information. Patient confidentiality is of

particular concern with measures based on patient surveys and when there are small numbers of patients.

### Usability and Use

**18.** An important outcome that may not have an identified improvement strategy still can be useful for informing quality improvement by identifying the need for and stimulating new approaches to improvement.

**19. Transparency** is the extent to which performance results about identifiable, accountable entities are *disclosed and available* outside of the organizations or practices whose performance is measured. Maximal transparency is achieved with **public reporting** defined as making comparative performance results about identifiable, accountable entities freely available (or at nominal cost) to the public at large (generally on a public website). *At a minimum, the data on performance results about identifiable, accountable entities are available to the public (e.g., unformatted database).* The capability to verify the performance results adds substantially to transparency.

**20.** This guidance is not intended to be construed as favoring measures developed by organizations that are able to implement their own measures (such as government agencies or accrediting organizations) over equally strong measures developed by organizations that may not be able to do so (such as researchers, consultants, or academics). Accordingly, measure developers may request a longer timeframe with appropriate explanation and justification.

**21. Credible plan** includes the specific program, purpose, intended audience, and timeline for implementing the measure within the specified timeframes. A plan for accountability applications addresses mechanisms for data aggregation and reporting.

**22.** Demonstrated progress toward achieving the goal of high-quality, efficient healthcare includes evidence of improved performance and/or increased numbers of individuals receiving high-quality healthcare. Exceptions may be considered with appropriate explanation and justification.

### Comparison to Related and Competing Measures

**23.** Measure harmonization refers to the standardization of specifications for related measures with the same measure focus (e.g., *influenza immunization* of patients in hospitals or nursing homes); related measures with the same target population (e.g., eye exam and HbA1c for *patients with diabetes*); or definitions applicable to many measures (e.g., age designation for children) so that they are uniform or compatible, unless differences are justified (e.g., dictated by the evidence). The dimensions of harmonization can include numerator, denominator, exclusions, calculation, and data source and collection instructions. The extent of harmonization depends on the relationship of the measures, the evidence for the specific measure focus, and differences in data sources.

1

## 2 Recommendations for Review Process

3 The TEP made several recommendations for the process of evaluating composite performance  
4 measures.

- 5
- 6 • Steering committees should include at least one member who is knowledgeable about composite  
7 measures and/or composite measures should undergo a methodological technical expert  
8 consultation.
  - 9 • If a steering committee recommends the removal of one or more components from the composite  
performance measure—and the developer is agreeable to the revised construction of the

- 1 composite—there should be an opportunity for the developer to respond to the recommendation
- 2 within the project rather than having to completely re-submit the revised measure at a later date.
- 3 • Provide examples of types of analyses for different types of composite performance measures. (See
- 4 Appendix B for a first step.)
- 5

1 **Notes**

- 2 1. Institute of Medicine, *Performance Measurement: Accelerating Improvement*, Washington, DC:  
3 National Academies Press; 2006.
- 4 2. National Quality Forum (NQF), *Composite Measure Evaluation Framework and National Voluntary*  
5 *Consensus Standards for Mortality and Safety-Composite Measures: A Consensus Report*,  
6 Washington, DC: National Quality Forum; 2009.
- 7 3. Booyesen F, An overview and evaluation of composite indices of development, *Social Indicators*  
8 *Research*, 2002;59:115-151.
- 9 4. Fayers PM, Hand DJ, Causal variables, indicator variables and measurement scales: an example from  
10 quality of life, *J R Statist Soc A*, 2002;165 (Part 2):233-261.
- 11 5. Kaplan SH, Normand SL, ., *Conceptual and Analytical Issues in Creating Composite Measures of*  
12 *Ambulatory Care Performance*, Washington, DC: National Quality Forum; 2006.
- 13 6. Nardo M, Saisana M, Saltelli A, et al., *Handbook on Constructing Composite Indicators: Methodology*  
14 *and User Guide. OECD Statistics Working Paper*, Paris, France: OECD Statistics Directorate; 2005.  
15 Report No.: STD/DOC(2005)3.
- 16 7. O'Brien SM, Shahian DM, DeLong ER, et al., Quality measurement in adult cardiac surgery: part 2--  
17 Statistical considerations in composite measure scoring and provider rating, *Ann Thorac Surg*,  
18 2007;83(4 Suppl):S13-S26.
- 19 8. Reeves D, Campbell SM, Adams J, et al., Combining multiple indicators of clinical quality: an  
20 evaluation of different analytic approaches, *Med Care*, 2007;45(6):489-496.
- 21 9. Shahian DM, Edwards FH, Ferraris VA, et al., Quality measurement in adult cardiac surgery: part 1--  
22 Conceptual framework and measure selection, *Ann Thorac Surg*, 2007;83(4 Suppl):S3-12.

1 **Appendix A: Glossary**

2

Term	Definition	Source
<p><b>All-or-None Scoring</b></p> <p><i>Also known as:</i></p> <ul style="list-style-type: none"> <li>• <i>Appropriateness model</i></li> <li>• <i>Conjunctive scoring</i></li> </ul>	<p>A percentage is determined by applying an all-or-none rule at the patient level. The denominator is the number of patients eligible to receive at least one of the identified elements of care, and the numerator is the number of patients who actually received all of the care for which the specific patient was eligible. No partial credit is given.</p>	<p>NQF Composite Guidance Report, 2007</p>
<p><b>Bundle</b></p>	<p>A series of interventions related to a specific condition that, when implemented together, will achieve significantly better outcomes than when implemented individually. This term was developed by faculty at the Institute for Healthcare Improvement. See <a href="http://www.ihl.org/IHI/Topics/CriticalCare/IntensiveCare/ImprovementStories/BundleUpforSafety.htm">www.ihl.org/IHI/Topics/CriticalCare/IntensiveCare/ImprovementStories/BundleUpforSafety.htm</a>.</p>	<p>NQF Composite Guidance Report, 2007</p>
<p><b>Clinimetric approach</b></p>	<p>Approach to developing a scale that relies on the required relationships between the observed items and the attribute for which an index is being defined. The most important attributes to be included in the index are not expected to be homogeneous because they indicate different aspects of a complex clinical phenomenon.</p>	<p>NQF Composite Guidance Report, 2007</p>
<p><b>Component</b></p>	<p>A constituent part or element of a composite measure.</p>	<p>NQF Composite Guidance Report, 2007</p>
<p><b>Composite measure</b></p>	<p>A combination of two or more individual measures into a single measure that results in a single score.</p>	<p>NQF Composite Guidance Report, 2007</p>
<p><b>Construct</b></p>	<p>An abstract phenomenon that is measured indirectly through less abstract indicators.</p>	<p>NQF Composite Guidance Report, 2007</p>
<p><b>Domain</b></p>	<p>A dimension or aspect of a construct.</p>	<p>NQF Composite Guidance</p>



Term	Definition	Source
		Report, 2007
<b>Indicator</b>	Sometimes used interchangeably with measure, but may indicate a more descriptive level than the term “measure,” which indicates the operational definition.	NQF Composite Guidance Report, 2007
<b>Indicator Average</b>	For each indicator, the percentage of times the indicator was met is computed. The scores are averaged across all indicators. This score represents the mean rate at which each audited aspect of care was met.	Reeves, 2007
<b>Item</b>	A single question on a measurement scale or instrument	NQF Composite Guidance Report, 2007
<b>Latent variable</b>	An unobserved trait or characteristic	NQF Composite Guidance Report, 2007
<b>Measure</b>	Numeric quantification of some concept. A quality measure is a numeric quantification of healthcare quality.	NQF Composite Guidance Report, 2007
<b>Opportunity scoring</b>	<p>Scoring used with process measures, determined from the sum of all numerators (achieved the desired process) divided by the sum of all denominators (i.e., number of eligible patients or opportunities, which could vary by measure).</p> <p>If the opportunity score is based on “care events” (patient/provider interactions), the opportunity score is the percentage of all care events that were met. For example, if patient A meets 1 of 1 opportunity and patient B meets 3 of 4 opportunities, then the care event opportunity score =80% [i.e., <math>(1+3)/(1+4)</math>].</p> <p>If the opportunity score is based on patients, the opportunity score is some function (typically the average) of the number of care events that were met for each patient. Using the above example, the patient-based opportunity score =88% [i.e., 100% met for patient A, 75% met for patient B → average over the 2</p>	NQF, Composite Guidance Report, 2007, Aligning Forces, 2010, Reeves, 2007

<b>Term</b>	<b>Definition</b>	<b>Source</b>
	patients= $100+75 / 2$ . (Has also been called “patient average”.)	
<b>Paired measures</b>	Individual measures that should be measured concurrently in the same population; however, the results are not combined into a single score.	NQF Composite Guidance Report, 2007
<b>Percentage Standard</b>	This is a less stringent version of the All-or-None method, where the criterion for success is that some percentage (e.g., 70%) or more of the triggered indicators be met.	Reeves, 2007
<b>Performance measure</b>	Numeric quantification of healthcare quality for a designated accountable healthcare entity, such as hospital, health plan, nursing home, clinician, etc.	PRO Report, 2012
<b>Psychometric approach</b>	Approach to developing a scale that relies on the relationships between the items that have been measured where the multiple component items are all measuring more or less the same single attribute.	NQF Composite Guidance Report, 2007
<b>Quality construct</b>	A hypothetical complex concept of quality.	
<b>Scale</b>	A measure of an attribute composed of a set of related items. A score on the scale represents a point along a continuum representing more or less of the attribute.	NQF Composite Guidance Report, 2007
<b>Subscale</b>	A measure of a dimension of a scale composed of a subset of the items in a scale.	NQF Composite Guidance Report, 2007
<b>Variable</b>	A characteristic or attribute that varies within and among people or the subjects of study.	NQF Composite Guidance Report, 2007

1

## 1 Appendix B: Approaches for Constructing Composite Performance Measures

Quality Construct	Description	Unique Considerations for Testing and Evaluating the Composite
<p>1. The quality construct is seen as causing the component performance measure scores</p> <ul style="list-style-type: none"> <li>• Also known as psychometric, reflective, scale, homogenous scale, dimensional</li> <li>• Example: NQF# 0530: Mortality for Selected Conditions (AHRQ)</li> </ul>	<ul style="list-style-type: none"> <li>• Scores on the component performance measures are considered the effect of quality (or caused by quality)</li> <li>• Component performance measures are considered a random sample of potential indicators of quality and should be interchangeable; therefore, focusing QI only on the component performance measures may not change the composite score</li> <li>• Component performance measures should be correlated with one another because they share common variance; and each component is correlated with the total composite score (omitting the component being assessed)</li> </ul> <p><b>Aggregation:</b></p> <p><b>Combination of multiple individual performance measures</b></p> <p>Various approaches may be used, including:</p> <ul style="list-style-type: none"> <li>▪ Opportunities [sum of all numerators / sum of all denominators]</li> <li>▪ Average/weighted average of component measure scores [score on A + score on B + score on C . . . / # of component performance measures]; or</li> <li>▪ Comparison to some benchmark (e.g., percentage of component performance measures that improved, reached 80%, etc.)</li> </ul>	
<p>2. The quality construct is seen as being caused or defined by the component performance measure scores</p> <ul style="list-style-type: none"> <li>• Also known as clinimetric, formative,</li> </ul>	<ul style="list-style-type: none"> <li>• Component performance measures are considered to cause quality (or define quality)</li> <li>• Component performance measures define the quality construct and should cover the entire scope of the quality construct; therefore, focusing QI on the component performance measures should change the composite score</li> <li>• Component performance measures do not</li> </ul>	

Quality Construct	Description	Unique Considerations for Testing and Evaluating the Composite
<p>index, heterogenous index, categorical</p> <p>Example:</p>	<p>need to be correlated with one another (but could be);?? each component should be correlated with the total composite score (omitting the component being assessed)??</p> <p><b>Aggregation:</b></p> <p><b>Combination of multiple individual performance measures</b></p> <p>Various approaches may be used, including:</p> <ul style="list-style-type: none"> <li>▪ Opportunities [sum of all numerators / sum of all denominators]</li> <li>▪ Average/weighted average of component measure scores [score on A + score on B + score on C . . . / # of component performance measures]; or</li> <li>▪ Comparison to some benchmark (e.g., percentage of component performance measures that improved, reached 80%, etc.)</li> </ul>	
<p>3. The quality construct is viewed as receiving all necessary care</p> <p>• Also known as All-or-None</p> <p>Example: NQF# 0729: Optimal Diabetes Care (MN Community Measurement)</p>	<p>Individual patient scores on component measures are considered to define quality and ALL must be achieved to signal quality</p> <p><b>Aggregation:</b></p> <p><b>Composite numerator</b> - Multiple components specified in the numerator and measured for each patient</p> <p>Percentage of patients who received all necessary components of care [# of patients in the denominator who met all components ( A and B and C and . . . ) / # of patients in target population]</p>	
<p>4. The quality construct is viewed as receiving necessary care, but receiving some is better</p>	<p>Individual patient scores on component measures are considered to define quality and achieving more is a signal of better quality</p>	

Quality Construct	Description	Unique Considerations for Testing and Evaluating the Composite
<p>than none</p> <ul style="list-style-type: none"> <li>• Also known as partial credit, percentage of necessary care</li> </ul> <p>Example:</p>	<p><b>Aggregation:</b></p> <p><b>Composite numerator</b> - Multiple components specified in the numerator and measured for each patient</p> <p>Average percentage of necessary components of care received by patient [Sum of percentage of components met (A, B, C . . .) for each patient in the denominator / # of patients in target population]</p>	
<p>5. The quality construct is viewed as not experiencing any healthcare-acquired adverse event/complication</p> <ul style="list-style-type: none"> <li>• Also known as any-or-all</li> </ul> <p>Example: NQF# 0564: Complications within 30 Days Following Cataract Surgery Requiring Additional Surgical Procedures (PCPI)</p>	<p>Individual patient scores on component measures are considered to define quality and NONE must be achieved to signal quality</p> <p><b>Aggregation:</b></p> <p><b>Composite numerator</b> - Multiple components specified in the numerator and measured for each patient</p> <p>Percentage of patients who experienced any of the component adverse events or complications [# of patient in the denominator who experienced A or B or C or . . . / # of patients in target population]</p>	
<p>6. The quality construct is defined by one concept but uses additional information on average performance to increase precision (reliability)</p> <p>Also known as reliability adjustment, shrinkage estimator</p>	<ul style="list-style-type: none"> <li>• Quality is defined by one performance measure but this measure is considered an unreliable estimate by itself, because of rare events or small case volume</li> </ul> <p><b>Aggregation:</b></p> <ul style="list-style-type: none"> <li>• Combines two rates of the same concept (e.g., a provider’s observed mortality rate and an average mortality rate for a specific category of providers such as quartile by patient case volume)</li> <li>• To-date has been used only with outcome measures</li> </ul>	

Quality Construct	Description	Unique Considerations for Testing and Evaluating the Composite
<p>Example: NQF# 0737: Survival Predictor for Esophagectomy Surgery (Leapfrog)</p>	<ul style="list-style-type: none"> <li>• Uses a provider characteristic (e.g., case volume) to categorize all providers for purposes of creating average rates for different categories</li> </ul> <p>(Weight x observed rate) + (weight x average rate)</p> <p>Weight is based on reliability of the provider observed rate, which is influenced by case volume</p>	

1

**2 The following are examples are not composite performance measures**

Conceptual Model	Description	Unique Considerations for evaluation
<p>7. Multi-item composites to measure individuals regardless of quality construct</p> <p>Example: Model 1-PHQ-9, CAHPS; Model 2-Apgar</p>	<ul style="list-style-type: none"> <li>• Multi-item scale, instrument, index, survey administered to individuals.</li> <li>• Patient data on these scales may be used in an individual performance measure or a composite performance measure; but the scale itself is not a performance measure and not eligible by itself for NQF endorsement.</li> </ul>	<ul style="list-style-type: none"> <li>• Not a composite <i>performance</i> measure</li> <li>• If patient data from such a scale is used in a performance measure, the reliability and validity of the scale also must be demonstrated.</li> <li>• See <a href="#">PRO project</a>.</li> </ul>
<p>8. Multiple aspects of quality are identified, but no quality construct for a composite is provided</p> <p>Example: NQF# 0101 Falls: Screening, Risk-Assessment, and Plan of Care to Prevent Future Falls (NCQA)</p> <p>CAHPS performance measures</p>	<p>The performance measures represent multiple individual aspects of quality</p> <p>There are two variants of this example:</p> <ul style="list-style-type: none"> <li>• Individual performance measures are identified as paired or grouped to be used and reported together to appropriately interpret results</li> <li>• Multiple related individual performance measures are submitted on one submission form, but require computation of individual performance measure scores; may have multiple denominators as well as</li> </ul>	<ul style="list-style-type: none"> <li>• Identifying paired or grouped measures is appropriate in some circumstances; however, individual performance measures need to be evaluated individually against all criteria and should be submitted on separate forms.</li> <li>• Including multiple individual measures in one form could obscure measure-specific information, making evaluation more difficult; it also is more difficult for others to find performance measures when they are catalogued with one</li> </ul>

	numerators. They do not necessarily need to be reported together for appropriate interpretation and if that were the case could be submitted as paired/grouped measures.	measure number.
--	--	-----------------

1

2

## 1 Appendix C: References Consulted

- 2 1. Agency for Healthcare Research and Quality (AHRQ), *Inpatient Quality Indicators Composite*  
3 *Measure Workgroup Final Report*, 2008. Available at  
4 <http://www.qualityindicators.ahrq.gov/Downloads/Modules/PSI/PSI%20Composite%20Development.pdf>.  
5 Last accessed October 2012.
- 6 2. Asch S, Hofer T. Representing overall quality of care: The whole must be more than the sum of the  
7 parts. 2008.
- 8 3. Ashby J, Juarez DT, Berthiaume J, et al., The relationship of hospital quality and cost per case in  
9 Hawaii, *Inquiry*, 2012;49(1):65-74.
- 10 4. Booyesen F, An overview and evaluation of composite indices of development, *Social Indicators*  
11 *Research*, 2002;59:115-151.
- 12 5. Diamantopoulos A, Winklhofer HM, Index construction with formative indicators: An alternative to  
13 scale development, *Journal of Marketing Research*, 2001;38(2):269-277.
- 14 6. Dijkers MP, Psychometrics and clinimetrics in assessing environments. A comment suggested by  
15 Mackenzie et al., 2002, *J Allied Health*, 2003;32(1):38-43.
- 16 7. Dimick JB, Staiger DO, Baser O, et al., Composite measures for predicting surgical mortality in the  
17 hospital, *Health Aff (Millwood)*, 2009;28(4):1189-1198.
- 18 8. Dimick JB, Staiger DO, Osborne NH, et al., Composite measures for rating hospital quality with major  
19 surgery, *Health Serv Res*, 2012;47(5):1861-1879.
- 20 9. Eapen ZJ, Fonarow GC, Dai D, et al., Comparison of composite measure methodologies for rewarding  
21 quality of care: an analysis from the American Heart Association's Get With The Guidelines  
22 program, *Circ Cardiovasc Qual Outcomes*, 2011;4(6):610-618.
- 23 10. Edwards JR, Bagozzi RP, On the nature and direction of relationships between constructs and  
24 measures, *Psychol Methods*, 2000;5(2):155-174.
- 25 11. Fayers PM, Hand DJ, Causal variables, indicator variables and measurement scales: an example from  
26 quality of life, *J R Statist Soc A*, 2002;165 (Part 2):233-261.
- 27 12. Felt-Lisk S, Lavin B, Gold M, Aggregate quality measures for the National Healthcare Quality Report:  
28 Summary of technical advisory panel meetings May/June 2005 (DRAFT), 2005;
- 29 13. Hess BJ, Weng W, Lynn LA, et al., Setting a fair performance standard for physicians' quality of  
30 patient care, *J Gen Intern Med*, 2011;26(5):467-473.
- 31 14. Holmboe ES, Weng W, Arnold GK, et al., The comprehensive care project: measuring physician  
32 performance in ambulatory practice, *Health Serv Res*, 2010;45(6 Pt 2):1912-1933.
- 33 15. Ingenix, Creating quality composite scores: Challenges and issues in physician quality measurement,  
34 2008;October 2012.



- 1 16. Institute of Medicine, *Performance Measurement: Accelerating Improvement*, Washington, DC:  
2 National Academies Press; 2006.
- 3 17. Jacobs R, Goddard M, Smith PC, How robust are hospital ranks based on composite performance  
4 measures?, *Med Care*, 2005;43(12):1177-1184.
- 5 18. Jacobs R, Goddard M, Smith PC, *Public Services: Are Composite Measures a Robust Reflection of*  
6 *Performance in the Public Sector*, 2006. Report No.: CHE Research Paper 16.
- 7 19. Kaplan SH, Griffith JL, Price LL, et al., Improving the reliability of physician performance assessment:  
8 identifying the "physician effect" on quality and creating composite measures, *Med Care*,  
9 2009;47(4):378-387.
- 10 20. Kaplan SH, Normand SL, *Conceptual and Analytical Issues in Creating Composite Measures of*  
11 *Ambulatory Care Performance*, Washington, DC: National Quality Forum; 2006.
- 12 21. Kianifard F, Evaluation of clinimetric scales: Basic principles and methods, *The Statistician*,  
13 1994;43(4):475-482.
- 14 22. Nardo M, Saisana M, Saltelli A, et al., *Handbook on Constructing Composite Indicators: Methodology*  
15 *and User Guide. OECD Statistics Working Paper*, Paris, France: OECD Statistics Directorate; 2005.  
16 Report No.: STD/DOC(2005)3.
- 17 23. National Committee for Quality Assurance (NCQA), MEMO: Summary of Alliance Use of Composite  
18 Measures, 2010;
- 19 24. National Quality Forum (NQF), *Composite Measure Evaluation Framework and National Voluntary*  
20 *Consensus Standards for Mortality and Safety-Composite Measures: A Consensus Report*,  
21 Washington, DC: National Quality Forum; 2009.
- 22 25. Nolan T, Berwick DM, All-or-none measurement raises the bar on performance, *JAMA*,  
23 2006;295(10):1168-1170.
- 24 26. Normand C, Measuring outcomes in palliative care: limitations of QALYs and the road to PALYs, *J*  
25 *Pain Symptom Manage*, 2009;38(1):27-31.
- 26 27. O'Brien SM, Shahian DM, DeLong ER, et al., Quality measurement in adult cardiac surgery: part 2--  
27 Statistical considerations in composite measure scoring and provider rating, *Ann Thorac Surg*,  
28 2007;83(4 Suppl):S13-S26.
- 29 28. Peterson ED, DeLong ER, Masoudi FA, et al., ACCF/AHA 2010 Position Statement on Composite  
30 Measures for Healthcare Performance Assessment: a report of American College of Cardiology  
31 Foundation/American Heart Association Task Force on Performance Measures (Writing Committee  
32 to Develop a Position Statement on Composite Measures), *J Am Coll Cardiol*, 2010;55(16):1755-  
33 1766.
- 34 29. Peterson ED, DeLong ER, Masoudi FA, et al., ACCF/AHA 2010 Position Statement on Composite  
35 Measures for Healthcare Performance Assessment: a report of the American College of Cardiology

- 1 Foundation/American Heart Association Task Force on Performance Measures (Writing Committee  
2 to develop a position statement on composite measures), *Circulation*, 2010;121(15):1780-1791.
- 3 30. Physician Consortium for Performance Improvement (PCPI), *Measures Development, Methodology,*  
4 *and Oversight Advisory Committee: Recommendations to PCPI Work Groups on Composite*  
5 *Measures*, 2010. Available at [http://www.ama-assn.org/resources/doc/cqi/composite-measures-](http://www.ama-assn.org/resources/doc/cqi/composite-measures-framework.pdf)  
6 [framework.pdf](http://www.ama-assn.org/resources/doc/cqi/composite-measures-framework.pdf). Last accessed October 2012.
- 7 31. Reeves D, Campbell SM, Adams J, et al., Combining multiple indicators of clinical quality: an  
8 evaluation of different analytic approaches, *Med Care*, 2007;45(6):489-496.
- 9 32. Ross JS, Correlation of inpatient and outpatient measures of stroke care quality within Veterans  
10 Health Administration hospitals, 2011;42(8):2269-2275.
- 11 33. Scholle SH, Roski J, Adams JL, et al., Benchmarking physician performance: reliability of individual  
12 and composite measures, *Am J Manag Care*, 2008;14(12):833-838.
- 13 34. Scholle SH, Roski J, Adams JL, et al., Benchmarking Physician Performance: Reliability of Individual  
14 and Composite Measures, *American Journal of Managed Care*, 2008;14(12):829-838.
- 15 35. Shahian DM, Edwards FH, Ferraris VA, et al., Quality measurement in adult cardiac surgery: part 1--  
16 Conceptual framework and measure selection, *Ann Thorac Surg*, 2007;83(4 Suppl):S3-12.
- 17 36. Shwartz M, Ren J, Pekoz EA, et al., Estimating a composite measure of hospital quality from the  
18 Hospital Compare database: differences when using a Bayesian hierarchical latent variable model  
19 versus denominator-based weights, *Med Care*, 2008;46(8):778-785.
- 20 37. Staiger DO, Dimick JB, Baser O, et al., Empirically derived composite measures of surgical  
21 performance, *Med Care*, 2009;47(2):226-233.
- 22 38. Streiner DL, Clinimetrics vs. psychometrics: an unnecessary distinction, *J Clin Epidemiol*,  
23 2003;56(12):1142-1145.
- 24 39. Streiner DL, Being inconsistent about consistency: when coefficient alpha does and doesn't matter, *J*  
25 *Pers Assess*, 2003;80(3):217-222.
- 26 40. Timbie JW, Shahian DM, Newhouse JP, et al., Composite measures for hospital quality using quality-  
27 adjusted life years, *Stat Med*, 2009;28(8):1238-1254.
- 28 41. Weifeng W, Hess BJ, Lynn LA, et al., Measuring physicians' performance in clinical practice:  
29 reliability, classification accuracy, and validity, *Eval Health Prof*, 2010;33(3):302-320.

30

1 **Appendix D: Technical Expert Panel and NQF Staff**

2 **TECHNICAL EXPERT PANEL**

3 **Patrick Romano, MD, MPH (Co-Chair)**

4 UC Davis School of Medicine  
5 Sacramento, CA

6 **Elizabeth R. DeLong, PhD (Co-Chair)**

7 Duke University Medical Center  
8 Durham, NC, State

9 **John D. Birkmeyer, MD**

10 University of Michigan  
11 Ann Arbor, MI, State

12 **Dale Bratzler, DO, MPH**

13 Oklahoma University Health Services Center  
14 Oklahoma City, OK, State

15 **James Chase, DO, MPH**

16 Minnesota Community Measurement  
17 Minneapolis, MN, State

18 **Nancy Dunton, PhD, FAAN**

19 University of Kansas Medical Center, School of Nursing  
20 Overland Park, KS, State

21 **Elizabeth Goldstein, PhD**

22 Centers for Medicare and Medicaid Services  
23 Baltimore, MD, State

24 **Sherrie Kaplan, PhD, MPH**

25 The University of California - Irvine  
26 Irvine, CA, State

27 **Lyn Paget, MPH**

28 Informed Medical Decisions Foundation  
29 Boston, MA, State

30 **David Shahian, MD**

31 Massachusetts General Hospital  
32 Boston, MA, State

33 **Steven Wright, PhD**

34 Veteran's Health Administration  
35 Providence, RI, State

1 **Alan Zaslavsky, PhD**  
2 Harvard Medical School  
3 Boston, MA, State  
4

5 **NQF STAFF**

6 **Helen Burstin, MD, MPH**  
7 Senior Vice President  
8 Performance Measures

9 **Heidi Bossley, MSN, MBA**  
10 Vice President  
11 Performance Measures

12 **Karen Pace, PhD, RN**  
13 Senior Director  
14 Performance Measures

15 **Karen Johnson, MS**  
16 Senior Director  
17 Performance Measures

18 **Elisa Munthali, MPH**  
19 Senior Project Manager  
20 Performance Measures