

Review and Update of Guidance for Evaluating Evidence and Measure Testing

DRAFT TECHNICAL REPORT FOR
REVIEW

~~August 8~~October 2, 2013

To: CSAC
From: Karen Pace
Date: October 2, 2013

This document includes red-line changes as a result of comments and discussion with CSAC/HITAC on September 10. There were two outstanding issues that we need to discuss.

CSAC Action:

- Approve updated guidance reflected in this document (or suggest changes).
- Make decision regarding face validity (see p. 19-20).
- Make decision regarding Trial Measures for Implementation and Testing (see p. 23-25).



NATIONAL
QUALITY FORUM

Contents

Background	4
Purpose	5
Evidence	5
Health Outcomes and Patient-Reported Outcomes (PRO).....	6
Quantity, Quality, Consistency of the Body of Evidence and Exceptions	7
Proposed Guidance for Evaluating the Clinical Evidence – Algorithm 1.....	7
Algorithm 1. Guidance for Evaluating the Clinical Evidence	9
Table 1: Evaluation of Quantity, Quality, and Consistency of Body of Evidence for Structure, Process, and Intermediate Outcome Measures.....	9
Measure Testing.....	12
Testing Data Elements vs. Performance Measure Score	13
More explicit Guidance on Minimum Thresholds and Types of Testing	14
Proposed Guidance on Evaluating Reliability and Validity – Algorithms 2 and 3	14
Algorithm 2. Guidance for Evaluating Reliability	16
Algorithm 3. Guidance for Evaluating Validity	17
Applying NQF Criteria for Endorsement to eMeasures	18
Clarification of Requirements for Endorsing eMeasures.....	20
Proposed Approval for Trial Use for eMeasures.....	22
Appendix A: Current Criteria and Guidance related to Clinical Evidence	25
1. Evidence, Performance Gap, and Priority (Impact)—Importance to Measure and Report	25
Table A-1: Evaluation of Quantity, Quality, and Consistency of Body of Evidence for Structure, Process, and Intermediate Outcome Measures	27
Table A-2: Evaluation of Subcriterion 1a Based on the Quantity, Quality, and Consistency of the Body of Evidence.....	28
Appendix B: Current Criteria and Guidance related to Reliability and Validity Testing	29
2. Reliability and Validity—Scientific Acceptability of Measure Properties	29
Table B-1: Evaluation Ratings for Reliability and Validity	32
Table B-2: Evaluation of Reliability and Validity of eMeasures	34
Table B-3: Evaluation of Scientific Acceptability of Measure Properties Based on Reliability and Validity Ratings.....	35
Table B-4: Scope of Testing Required at the Time of Review for Endorsement Maintenance	36
Appendix C: Current Criteria and Guidance related to Feasibility.....	37
3. Feasibility:	37
NATIONAL QUALITY FORUM	2

Guidance on Evaluating Feasibility	37
Table C-1: Generic Scale for Rating Feasibility Subcriteria	37
Table C-2. Data Element Feasibility Scorecard	38
Appendix D: Project Steering Committee and NQF Staff.....	40
Consensus Standards Approval Committee Member Roster	40
Health Information Technology Advisory Committee Roster.....	42
NQF Staff	45

1 Review and Update of Guidance for Evaluating Evidence 2 and Measure Testing

3 DRAFT TECHNICAL REPORT

4 Background

5 NQF endorses performance measures that are suitable for both accountability applications (e.g., public
6 reporting, accreditation, performance-based payment, network inclusion/exclusion, etc.) as well as
7 internal quality improvement efforts. NQF's measure evaluation criteria and subcriteria are used to
8 determine the suitability of measures for use in these activities. Because endorsement initiates
9 processes and infrastructure to collect data, compute performance results, report performance results,
10 and improve and sustain performance, NQF endorsement is intended to identify those performance
11 measures that are most likely to facilitate achievement of high quality and efficient healthcare for
12 patients. The criteria and subcriteria also relate to the concept of "fit for purpose". For example, the
13 clinical evidence should support use of a measure with a specific target patient population (e.g., foot
14 care for patients with diabetes) and testing of the measure as specified indicates under what
15 circumstances reliable and valid results may be obtained (i.e., using the measure with a specified data
16 source and level of analysis or for the accountable entity whose performance is being measured).

17 Throughout the various iterations of the NQF [measure evaluation criteria](#), the basic criteria and
18 concepts have remained largely unchanged. However, the measure evaluation guidance—which focuses
19 on the specificity and rigor with which the criteria are applied—has become more comprehensive and
20 more specific over time. The guidance on measure evaluation is intended first for steering committees
21 that evaluate performance measures and make recommendations for NQF endorsement, as well as the
22 staff who assist them. Second, the guidance informs measure developers about how to demonstrate
23 that a measure meets the criteria. Third, the guidance informs NQF members and the public about how
24 measures are evaluated and informs those who use NQF-endorsed performance measures about what
25 endorsement means.

26 In 2010, the NQF convened two task forces to help provide guidance for evaluating the clinical evidence
27 and the measure testing results for reliability and validity that is submitted in support of a measure. The
28 approved recommendations were implemented in 2011. Testing of eMeasures also was addressed in
29 the 2011 guidance and in some subsequent draft policy statements.

30 Purpose

31 The purpose of this project was to review the implementation of the 2011 guidance on evaluating
32 evidence and measure testing (including eMeasure testing requirements) and to propose modifications
33 to address any major challenges. Modifications that would potentially increase consistency and clarity
34 in the evaluation of performance measures for potential NQF endorsement also were considered.

35 The specific goals of the project included:

- 36 • promote consistency in evaluation across measures and projects;
- 37 • clarify common misunderstandings about the criteria and guidance;
- 38 • remain consistent with the criteria and principles from the 2011 guidance (i.e., do not change
39 the “bar” for endorsement or the information requested for a measure submission); and
- 40 • address the current challenges with eMeasure testing.

41
42 This project was not intended to suggest changes to the basic measure evaluation criteria or to the
43 consensus development process (CDP). Other related concerns, such as levels of endorsement,
44 endorsement for specific applications, endorsing measures intended only for quality improvement, and
45 definitions of multistakeholder consensus are being addressed through the Board strategic planning
46 process, to be followed by additional work as indicated.

47 The Consensus Standards Approval Committee (CSAC) reviewed and discussed the measure evaluation
48 criteria and guidance at its in-person meetings in March and July 2013, as well as in their monthly calls in
49 May and June. A smaller subcommittee of the CSAC, formed to more thoroughly consider the issues and
50 offer suggestions for modifications than was possible for the full CSAC, met via conference calls in June
51 and July. The Health Information Technology Advisory Committee (HITAC) discussed eMeasure testing
52 requirements via conference call in May 2013 and at its in-person meeting in July 2013. A
53 subcommittee of the HITAC also was formed to offer specific recommendations regarding eMeasure
54 testing; this subcommittee met via conference call in August 2013.

55 This report presents some potential modifications to the 2011 guidance for evaluating evidence and
56 measure testing (including eMeasure testing) for public and NQF member review and comment. Also
57 included in this report is a proposal for another pathway to endorsement for eMeasures. The associated
58 criteria and prior guidance are provided in the appendices.

59 Although simplicity is desired when possible, the evaluation of evidence, reliability, and validity is
60 complex, requiring both objective information such as the clinical evidence and testing results and
61 steering committee judgment to review and reach a conclusion regarding what is sufficient to
62 recommend a performance measure for NQF endorsement.

63 Evidence

64 The most common issues and challenges related to implementing the 2011 guidance on evaluating the
65 clinical evidence (Appendix A) included:

- 66 • Measures were submitted without a summary of the quantity, quality, and consistency of the
67 evidence from a systematic review of a body of evidence. The reasons varied across measures

68 and developers, but the end result was that the rating scale could not be applied consistently.
69 Therefore, the steering committees either rated this subcriterion as insufficient evidence or
70 relied upon their own knowledge and memory of the evidence. This resulted in inconsistency
71 across measures and/or projects.

- 72 • Inconsistent handling of exceptions to the evidence requirement for measures that were not
73 directly evidence-based or focused on distal process steps (e.g., document a diagnosis, order a
74 lab test) with either indirect evidence or no empirical evidence.
- 75 • Submitted evidence was about something other than what was being measured, or provided
76 only indirect evidence.
- 77 • A common misunderstanding was that the guidance on evidence required randomized
78 controlled trials (RCT).

79 In addition, the patient-reported outcomes (PROs) project raised the question of whether NQF should
80 apply the same evidence requirements for PROs and health outcomes.

81 The CSAC and its subcommittee addressed three key questions.

- 82 1. Should NQF require a systematic review of the evidence that health outcomes and PROs are
83 influenced by healthcare processes or structures?
- 84 2. Should NQF's current guidance requiring evidence that is based on a systematic review of the
85 body evidence to support intermediate clinical outcomes, processes, and structures be less
86 stringent?
- 87 3. When should an exception to the evidence requirement be considered?

88 Health Outcomes and Patient-Reported Outcomes (PRO)

89 NQF has stated a hierarchical preference for performance measures of health outcomes. Current
90 criteria require a rationale that such outcomes are influenced by healthcare processes or structures but
91 do not require a review of the quantity, quality, and consistency of evidence. The approved
92 recommendations from the project on [PROs in Performance Measurement](#) established that PROs should
93 be treated the same as other health outcomes and that the CSAC should review the question of
94 evidence requirements. PROs include health-related quality of life/functional status, symptom and
95 symptom burden, experience with care, and health-related behaviors.

96 Outcomes such as improved function, survival, or relief from symptoms are the reasons patients seek
97 care and providers deliver care; they also are of interest to purchasers and policymakers. Outcomes are
98 integrative, reflecting the result of all care provided over a particular time period (e.g., an episode of
99 care). Measuring performance on outcomes encourages a "systems approach" to providing and
100 improving care. Measuring outcomes also encourages innovation in identifying ways to impact or
101 improve outcomes that might have previously been considered not modifiable (e.g., rate of central line
102 infection). Due to differences in severity of illness and comorbidities, not all patients are expected to
103 have the same probability of achieving an outcome; therefore, performance measures of health
104 outcomes and PROs are subject to the additional criterion of risk adjustment under validity.

105 The CSAC reaffirmed the prior guidance for health outcomes (now also applied to PROs) that requires
106 only a rationale that the measured outcome is influenced by at least one healthcare process, service
107 intervention, treatment, or structure.

108 **Quantity, Quality, Consistency of the Body of Evidence and Exceptions**

109 The CSAC also reaffirmed the criteria and guidance that calls for an assessment of the strength of the
110 evidence from a systematic review of the body of evidence for performance measures of intermediate
111 clinical outcomes, processes, or structures. This is consistent with the standards established by the
112 Institute of Medicine (IOM) for [systematic reviews](#) and [guidelines](#). The evidence should demonstrate
113 that the intermediate outcome, process, or structure influences desired outcomes. Evidence refers to
114 empirical studies, but is not limited to RCTs. Because endorsement sets in motion an infrastructure to
115 address the performance measure, the intent of the evidence subcriterion is to ensure that endorsed
116 measures focus on those aspects of care known to influence patient outcomes.

117 The CSAC and subcommittee also reaffirmed the need for exceptions to the evidence subcriterion. Not
118 all healthcare is evidence-based and systematic reviews as called for by the IOM may not be currently
119 available or the details readily accessible to obtain information on the quantity, quality, and consistency
120 of the evidence. However, exceptions should not be considered routine and more specific guidance is
121 needed to promote greater consistency.

122 **Proposed Guidance for Evaluating the Clinical Evidence – Algorithm 1**

123 Algorithm 1 presents a modified approach to guide steering committee evaluation of the evidence
124 submitted with a performance measure. It is consistent with the prior guidance (Appendix A) but is
125 intended to clarify and promote greater consistency and transparency.

126 The key features of this proposed guidance include:

- 127 • Preserves current requirement for a rationale for measures of health outcomes and PROs.
- 128 • Preserves the basic principles of transparency and evaluating the quantity, quality, and
129 consistency of the evidence.
- 130 • Accommodates the fact that some evidence reviews for guidelines may not be up to IOM
131 standards or the information on quantity, quality, and consistency of the body of evidence may
132 not be available. If evidence was graded but the submission did not include a summary of
133 quantity, quality, and consistency, it could potentially receive a moderate rating.
- 134 • Explicitly addresses what to do if a summary of quantity, quality, and consistency of the body of
135 evidence from a systematic review is not provided in the submission form– i.e., moderate is the
136 highest potential rating (see boxes 4 and 6).
- 137 • Preserves flexibility for exceptions to the evidence, but identifies specific questions for
138 considering the exception (boxes 7-9).
- 139 • Explicitly identifies how to handle measures that are based on expert opinion, indirect evidence,
140 or distal process steps (box 3 and exceptions) and therefore need to be explicitly addressed as a
141 potential exception.
- 142 • Uses specific examples of grades from [JSPSTF](#) and [GRADE](#) in addition to the NQF rating scale
143 (Table 1).
- 144 • The final ratings (other than for health outcomes and PROs) are high, moderate, low, and
145 insufficient evidence and are consistent with the prior guidance where high and moderate
146 ratings would be acceptable for endorsement. The ratings would indicate different levels of
147 strength/certainty of the evidence, magnitude of net benefit, as well as transparency, which
148 may be useful to implementers.
- 149 • The guidance still requires judgment of the steering committee.

Formatted: Space After: 11 pt, Line spacing: single

150 The comments reflected some of the differences in perspectives about evidence requirements for
151 measures of health outcomes and PROs. Most commenters did not mention evidence for outcome
152 measures; two commenters specifically agreed with the current approach; and three commenters
153 suggested that measures of outcomes have the same evidence requirements or additional empirical
154 analysis. In July 2013, the CSAC had reaffirmed the current criteria and guidance that requires a
155 rationale that supports that the outcome is influenced by at least one healthcare structure, process,
156 intervention, or service. This updated guidance was intended to reflect the current criteria, not change
157 or “raise the bar.” The CSAC again reaffirmed the current criteria and guidance related to health
158 outcomes and PROs.

Formatted: Indent: Left: 0"

159
160 The CSAC is particularly interested in receiving requested comments on when exceptions to the
161 evidence criterion should be considered. Most commenters agreed that the guidance would help
162 promote greater consistency. In response to the comments, some revisions to the algorithm were made
163 as follows.

- 164 • In Box 2 “plausible” rationale was replaced with a question that mirrors the criterion and the
165 information provided in the measure submission.
- 166 • A section was added to address when there is empirical evidence that has not yet been
167 systematically reviewed and graded (boxes 6-8).
- 168 • Some clarification provided regarding exceptions (boxes 9-11).

Formatted: List Paragraph, Bulleted + Level: 1
+ Aligned at: 0.25" + Indent at: 0.5"

Formatted: Font:

169
170

<p>1. Does the measure assess performance on a health outcome (e.g., mortality, function, health status, complication) or PRO (e.g., HRQoL/function, symptom, experience, health-related behavior)?</p>	<p>→ Yes → → → → → → → →</p>	<p>2. Does the SC agree that the relationship between the measured health outcome/PRO and at least one healthcare action (structure, process, intervention, or service) is identified (stated or diagrammed) and supported by the stated rationale?</p>	<p>→ Yes</p>	<p>Pass</p>
		<p>→ No → → → → → → → → → → → →</p>	<p>→</p>	<p>No pass</p>

↓ No

<p>3. For measures that assess performance on an intermediate clinical outcome, process, or structure - is it based on a systematic review (SR) and grading of the BODY of empirical evidence where the specific focus of the evidence matches what is being measured? (Evidence means empirical studies of any kind, the body of evidence could be one study; SR may be associated with a guideline)</p> <p>Answer NO if any: *Evidence is about something other than what is measured *Empirical evidence submitted but not systematically reviewed *Based on expert opinion *No evidence because it won't be studied (e.g., document a diagnosis) *Distal process step is not the specific focus of the evidence (e.g., monitor BP each visit, when evidence is about treatment of hypertension or relationship to mortality)</p>	<p>4. Is a summary of the quantity, quality, and consistency (QQC) of the body of evidence from a SR provided in the submission form?</p> <p>A SR is a scientific investigation that focuses on a specific question and uses explicit, prespecified scientific methods to identify, select, assess, and summarize the findings of similar but separate studies. It may include a quantitative synthesis (meta-analysis), depending on the available data. (IOM)</p> <p>Answer NO if: *Specific information on QQC not provided (general statements/conclusions, lists/descriptions of individual studies is not sufficient)</p>	<p>5a. Does the SR conclude: *Quantity:Mod/High; Quality:High; Consistency:High (Table 1) *High certainty that the net benefit is substantial (e.g., USPSTF-A) *High quality evidence that benefits clearly outweigh undesirable effects (e.g., GRADE-Strong) *If measuring Inapprop. care, Mod/Hi certainty of no net benefit or harm outweighs benefit (USPSTF-D)</p> <p>5b. Does the SR conclude: *Quantity:Low-High; Quality:Mod; Consistency:Mod/High (Table 1) *Moderate certainty that the net benefit is substantial OR moderate-high certainty the net benefit is moderate (e.g., USPSTF-B)</p> <p>5c. Does the SR conclude: *Consistency:Low (Table 1); *Moderate/high certainty that the net benefit is small (e.g., USPSTF C); OR no net benefit, or harm outweighs benefit (USPSTF-D) *Low quality evidence, desirable/undesirable effects closely balanced, uncertainty in preference or use of resources (e.g., GRADE-Weak)</p>	<p>→ Yes</p> <p>→ Yes</p> <p>→ Yes</p>	<p>Rate as High</p> <p>Rate as Moderate</p> <p>Rate as Low</p>
<p>↓ No - without QQC from SR, moderate is highest potential rating</p>				

		<p>6. Does the grade for the evidence or recommendation indicate: *High quality evidence (e.g., Table 1 - Quant:Mod/Hi; Qual:Hi; Consist:Hi; USPSTF - High certainty; GRADE-High quality) *Strong recommendation (e.g., GRADE -Strong; USPSTF-A)</p> <p>Answer NO if: *No grading of evidence and summary of QQC not provided *Not graded high quality or strong recommendation</p>	→ Yes		Rate as Moderate
↓ No		→No (moderate/weak quality or recommendation without QQC) →	→		Rate Low
6. Is empirical evidence submitted but without systematic review and grading of the evidence?	→ Yes	7. Does the empirical evidence that is summarized include all studies in the body of evidence? Answer NO if only selected studies included	→ Yes	8. Does the SC agree that the submitted evidence indicates high certainty that benefits clearly outweigh undesirable effects? (without SR, the evidence should be high-moderate quality and indicate substantial net benefit - See Table 1)	→ Yes Rate as Moderate
↓ No		No→→→→→→→→→→	→	No→→→→→→→→→→→→→→→→	→ Rate Low
9. Are there, OR could there be, performance measures of a related health outcome, OR evidence-based intermediate clinical outcome or process? <i>Example for YES: Propose to measure whether BP is assessed each visit instead of BP control or use of effective treatment</i>	→ No	10. Is there evidence of a systematic assessment of expert opinion (e.g., national/international consensus recommendation) that the benefits of what is being measured outweigh potential harms?	→ Yes	11. Does the SC agree that it is OK (or beneficial) to hold providers accountable for performance in the absence of empirical evidence of benefits to patients? (Consider potential detriments to endorsing the measure, e.g., focus attention away from more impactful practices, more costly without certainty of benefit; divert resources from developing more impactful measures.)	→ Yes Rate as Insufficient evidence with Exception
↓ Yes		↓ No		↓ No	Rate as Insufficient
No exception→→→→→	→	No exception→→→→→	→	No exception→→→→→→→→→→→→	→

174 Table 1: Evaluation of Quantity, Quality, and Consistency of Body of Evidence for
 175 Structure, Process, and Intermediate Outcome Measures

DEFINITION/ RATING	QUANTITY OF BODY OF EVIDENCE	QUALITY OF BODY OF EVIDENCE	CONSISTENCY OF RESULTS OF BODY OF EVIDENCE
Definition	Total number of studies (not articles or papers)	Certainty or confidence in the estimates of benefits and harms to patients across studies in the body of evidence related to study factors^a including: study design or flaws; directness/indirectness to the specific measure (regarding the population, intervention, comparators, outcomes); imprecision (wide confidence intervals due to few patients or events)	Stability in both the direction and magnitude of clinically/practically meaningful benefits and harms to patients (benefit over harms) across studies in the body of evidence
High	5+ studies ^b	Randomized controlled trials (RCTs) providing direct evidence for the specific measure focus, with adequate size to obtain precise estimates of effect, and without serious flaws that introduce bias	Estimates of clinically/practically meaningful benefits and harms to patients are consistent in direction and similar in magnitude across the preponderance of studies in the body of evidence
Moderate	2-4 studies ^b	<ul style="list-style-type: none"> • Non-RCTs with control for confounders that could account for other plausible explanations, with large, precise estimate of effect OR • RCTs without serious flaws that introduce bias, but with either indirect evidence or imprecise estimate of effect 	<p>Estimates of clinically/practically meaningful benefits and harms to patients are consistent in direction across the preponderance of studies in the body of evidence, but may differ in magnitude</p> <p>If only one study, then the estimate of benefits greatly outweighs the estimate of potential harms to patients (one study cannot achieve high consistency rating)</p>
Low	1 study ^b	<ul style="list-style-type: none"> • RCTs with flaws that introduce bias OR • Non-RCTs with small or imprecise estimate of effect, or without control for confounders that could account for other plausible explanations 	<ul style="list-style-type: none"> • Estimates of clinically/practically meaningful benefits and harms to patients differ in both direction and magnitude across the preponderance of studies in the body of evidence OR • wide confidence intervals prevent estimating net benefit <p>If only one study, then estimate of benefits do not greatly outweigh harms to patients</p>

DEFINITION/ RATING	QUANTITY OF BODY OF EVIDENCE	QUALITY OF BODY OF EVIDENCE	CONSISTENCY OF RESULTS OF BODY OF EVIDENCE
Insufficient to Evaluate (See Table 3 for exceptions.)	<ul style="list-style-type: none"> No empirical evidence OR Only selected studies from a larger body of evidence 	<ul style="list-style-type: none"> No empirical evidence OR Only selected studies from a larger body of evidence 	No assessment of magnitude and direction of benefits and harms to patients

176 ^aStudy designs that affect certainty of confidence in estimates of effect include: randomized controlled trials (RCTs), which
177 control for both observed and unobserved confounders, and non-RCTs (observational studies) with various levels of control for
178 confounders.
179 Study flaws that may bias estimates of effect include: lack of allocation concealment; lack of blinding; large losses to follow-up;
180 failure to adhere to intention to treat analysis; stopping early for benefit; and failure to report important outcomes.
181 Imprecision with wide confidence intervals around estimates of effects can occur in studies involving few patients and few
182 events.
183 Indirectness of evidence includes: indirect comparisons (e.g., two drugs compared to placebos rather than head-to head); and
184 differences between the population, intervention, comparator interventions, and outcome of interest and those included in the
185 relevant studies.¹⁵
186 ^bThe suggested number of studies for rating levels of quantity is considered a general guideline.

187 Measure Testing

188 The challenges related to implementing the 2011 guidance on evaluating measure testing for reliability
189 and validity (Appendix B) included:

- 190 • Lack of understanding of differences between testing using patient level data versus testing
191 using the computed performance measure score.
- 192 • Measure testing that is not consistent with the measure as specified (including data
193 specifications and level of analysis).
- 194 • No empirical statistical testing for reliability (e.g., descriptive statistics, only report that it is in
195 use with descriptive statistics on performance, report only a process for data management and
196 cleaning or computer programming; report only percent agreement for inter-rater reliability).
- 197 • The rating scale did not differentiate varying levels of confidence in the results, such as when
198 the scope of testing is narrow (e.g., 3-4 sites), or when the reliability statistic is only marginally
199 acceptable.
- 200 • Measures were submitted for endorsement with testing results that indicated the data or the
201 measure was not reliable or valid.
- 202 • Concerns about misclassification relate to reliability of the computed performance measure
203 score (given that validity is demonstrated), but current criteria allow for testing of the data
204 elements only (i.e., do not require testing at the measure score level).
- 205 • Confusion between clinical evidence for a process being measured versus validity of the
206 performance measure as specified.
- 207 • Complexity of concepts of reliability and validity, including measure testing methods, statistical
208 methods, and interpretation of results. Some may not be prepared to evaluate whether testing
209 used an appropriate method, with an adequate sample, and obtained sufficient results.

210 • The criteria allow face validity and many measures are submitted with only face validity.
211 Sometimes the same group of experts who helped develop the measure is used to establish face
212 validity, or the assessment did not address the primary validity issue of whether the
213 performance measure score from the measure as specified represents an accurate reflection of
214 quality of care. Therefore, face validity may be questioned, especially when threats to validity
215 such as exclusions are not adequately assessed.

216 The above issues also apply to eMeasures; but the most common challenges for eMeasures included:

- 217 • Measures were submitted without standard eMeasure specifications (HQMF and QDM).
- 218 • Testing that did not use electronic data (e.g., two manual abstractions).
- 219 • “Retooled” eMeasure specifications that could not be implemented.
- 220 • Difficulty recruiting test sites for testing and obtaining data from EHRs.

221 The CSAC and its subcommittee addressed two key questions.

- 222 • Should the rating scale reflect different levels of testing and different levels of confidence in the
223 results?
- 224 • Can the guidance be more explicit, with recommended methods and minimum thresholds for
225 samples and results?

226 In addition, the CSAC and HITAC addressed two key questions regarding eMeasures:

- 227 • Should specific thresholds for scope of testing or required type of testing be identified for
228 eMeasures?
- 229 • How can NQF facilitate progress with eMeasures while maintaining the same criteria for
230 endorsement as for other measures?

231 Testing Data Elements vs. Performance Measure Score

232 Data elements refer to the patient-level data used in constructing performance measures. For example,
233 if the performance measure is the percentage of patients 65 and older with a diagnosis of diabetes with
234 Hba1c>9 in the measurement year, then age, diagnosis (and possibly medications or lab values) are used
235 to identify the target population of patients with diabetes for the denominator as well as potential
236 exclusions (e.g., pregnant women) and the Hba1c lab value and date identify what is being measured for
237 the numerator. Reliability and validity of the data elements are different from that of the computed
238 performance score. Reliable and valid data are important building blocks for performance measures, but
239 ultimately the computed performance measure scores are what are used to make conclusions about the
240 quality of care provided. The question is whether the performance measure score can distinguish real
241 differences (signal) among providers from measurement error (noise) and whether that signal is a
242 reflection of the quality of care. These are relevant questions whether using the performance results to
243 identify areas for improvement activities, or for purposes of accountability. The CSAC and subcommittee
244 agreed that the rating scale should be modified slightly to reflect the difference between testing data
245 element and performance measure scores but in such a way that the “bar” for endorsement isn’t
246 changed. For example, face validity and testing at the level of data elements should continue to be
247 acceptable options.

248 **More explicit Guidance on Minimum Thresholds and Types of Testing**

249 Steering Committees often question what is considered an adequate sample for testing, and what is
250 considered an acceptable result. However, due to the various factors and context that should be
251 considered, the Measure Testing Task Force did not set minimum thresholds; nonetheless, they did
252 identify some basic principles (e.g., using a representative sample of a size that was sufficient for the
253 question and statistical method). This guidance provides much flexibility, but this flexibility can also
254 increase uncertainty in the evaluation process and can also increase the potential for inconsistency in
255 evaluation between measures and projects. While the CSAC and subcommittee would like to have
256 provided some guidance regarding minimum thresholds, they again noted the difficulties in determining
257 such thresholds and the need for steering committees to have flexibility to make judgments. For
258 example, 0.70 is most often cited a minimum threshold for most reliability statistics, however, a higher
259 threshold may be indicated for specific uses and 0.6 may be used for kappa.

260 Similarly, the Measure Testing Task Force report identified a variety of options for empirical testing and
261 did not prescribe a particular method. The CSAC and subcommittee suggested that proposed guidance
262 should reference the most common testing approaches but not limit measure developers from using
263 other approaches to address the same questions.

264 | *In response to ~~the CSAC is interested in receiving request for~~ comments on whether specific thresholds
265 for the reliability statistic or sample size used in measure testing should be specified in the rating scales
266 for reliability and validity, most commenters agreed that it is difficult or impossible to identify minimum
267 thresholds that are applicable to all testing situations. Three commenters suggested considering power
268 analysis to determine the appropriate sample size for the statistical test used. This will require further
269 exploration.*

270

271 **Proposed Guidance on Evaluating Reliability and Validity – Algorithms 2 and 3**

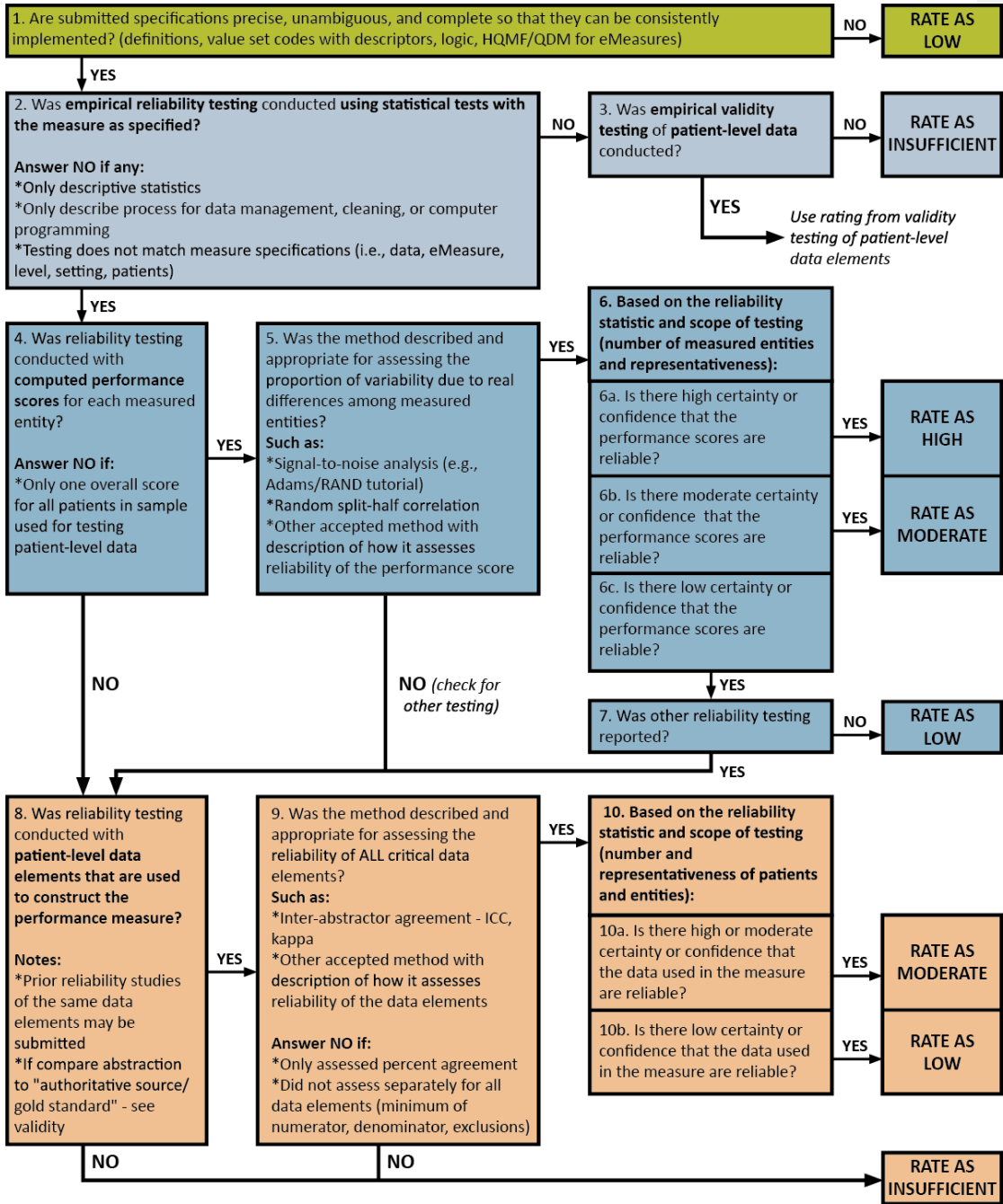
272 Algorithms 2 and 3 present modified approaches to guide steering committee evaluation of the
273 reliability (Algorithm 2) and validity (algorithm 3) for all measures (including eMeasures). They are
274 consistent with the prior guidance (Appendix B) but are intended to clarify and promote greater
275 consistency and transparency.

276 The key features of this proposed guidance include:

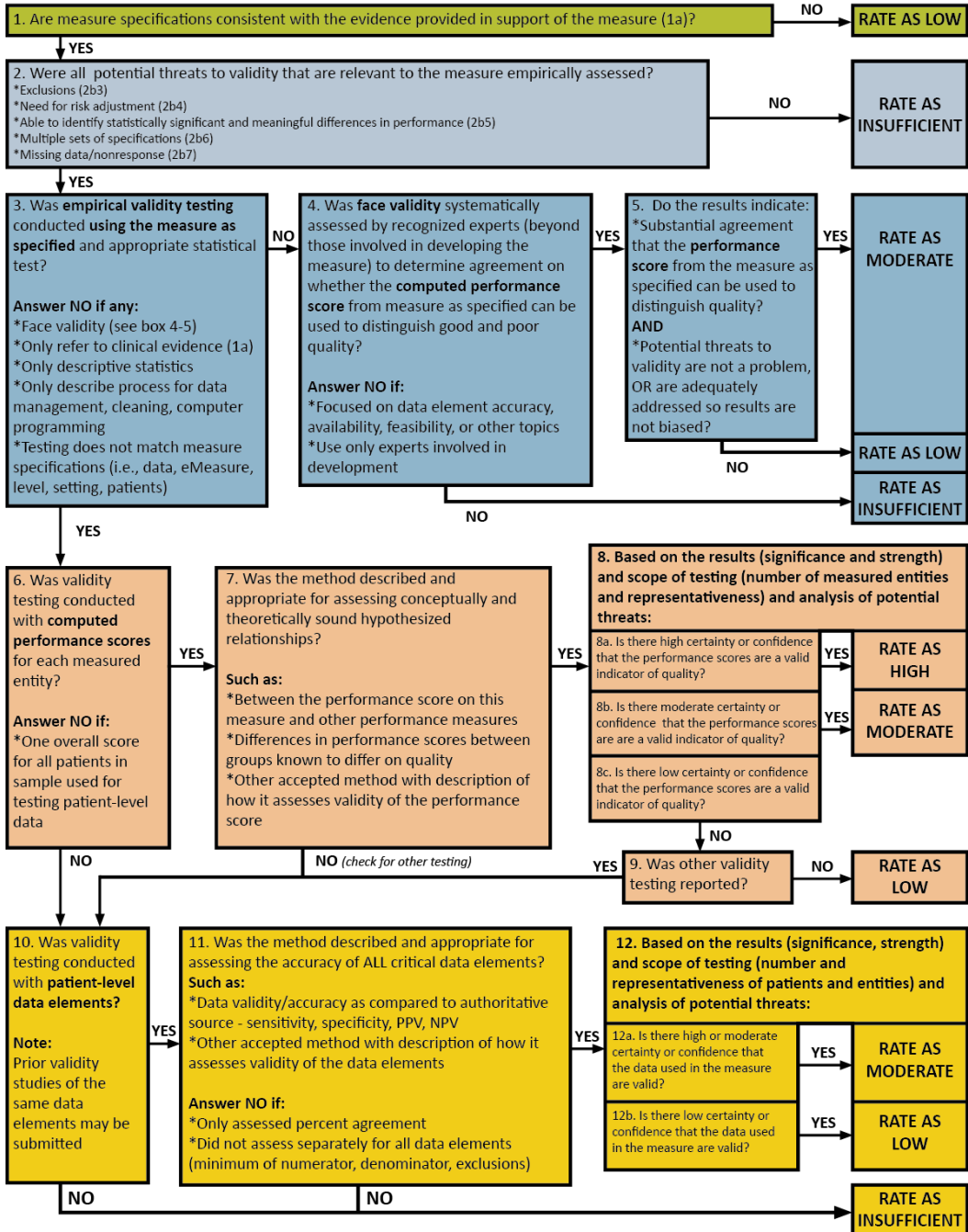
- 277 • Preserves most aspects of the 2011 rating scales:
- 278 ○ If tested at both levels, a measure would potentially receive a high rating depending on the
 - 279 assessment of results and scope of testing.
 - 280 ○ Testing only at the level of data elements would be rated as previously – the highest
 - 281 potential rating is moderate, depending on results and scope of testing.
 - 282 ○ Face validity of the performance measure score is eligible for a moderate rating if
 - 283 appropriate method, scope, and result.
- 284 • The main modification to the rating scales is that testing at the level of the performance measure
- 285 score alone could be eligible for a high rating, depending on result and scope of testing.

- 286 • Clarifies some common misunderstandings about testing (e.g., testing must be conducted with the
287 measure as specified; clinical evidence is not a substitute for validity testing of the measure; data
288 element level refers to patient-level data).
- 289 • Reinforces that testing of patient level data elements should include all critical data elements, but at
290 minimum must include a separate assessment and results for numerator, denominator, and
291 exclusions.
- 292 • Preserves the option to use data element validity testing for meeting both reliability and validity at
293 the data element level.
- 294 • Reinforces that if empirical testing was not conducted or an inappropriate method was used, there
295 is no information about reliability or validity, leading to a rating of insufficient. This preserves the
296 distinction between insufficient information versus demonstrating low reliability or validity.
297
298

299 Algorithm 2. Guidance for Evaluating Reliability



300 Algorithm 3. Guidance for Evaluating Validity



Discussion and Key Questions – Face Validity

Although empirical validity testing is preferable, NQF’s criteria allow for use of face validity in lieu of empirical testing. Because face validity is the weakest form of validity, the updated guidance proposed that systematic assessment of face validity involve experts who were not involved in measure development. Commenters requested more clarity about defining “involvement in measure development” and two commenters specifically disagreed with requiring experts beyond those involved in measure development and also noted this would “raise the bar” for meeting NQF criteria.

- Recognizing that the updated guidance could represent a “higher bar” and the variety of roles that people may play in measure development, do you agree to make no change to the current guidance? (i.e., developers choose who to use for systematic assessment of face validity)

However, the current criterion and its implementation by developers seems to represent an unnecessary exercise to “check the box”. Should we assume that a developer would not submit a measure for potential endorsement if its experts did not think it had face validity? Does a systematic assessment of face validity offer any additional information, especially at the time of initial endorsement?

- Should NQF consider revising guidance on face validity?

Two potential options:

- Assume the measure has face validity with the developer and its experts at the time of initial endorsement and don’t require a systematic assessment. Then require empirical validity testing at the time of endorsement maintenance.
This is consistent with the approved recommendation from Measure Testing Task Force that testing be expanded on endorsement maintenance if a measure did not already achieve a high rating, so empirical testing would be required at that time. However, currently such expanded testing is often not conducted and NQF has not implemented this requirement to date.)
- Require empirical validity testing at the time of initial endorsement.
This would represent a higher bar and would need to be phased in.

302

Revision to Algorithm 3

<p>3. Was empirical validity testing conducted using the measure as specified and appropriate statistical test?</p> <p>Answer NO if any: *Face validity (see box 4-5) *Only refer to clinical evidence (1a) *Only descriptive statistics *Only describe process for data management, cleaning, computer programming *Testing does not match measure specifications (i.e., data, eMeasure, level, setting, patients)</p>	<p>→ No</p>	<p>4. Was face validity systematically assessed by recognized experts (beyond those involved in developing the measure) to determine agreement on whether the computed performance measure score from measure as specified can be used to distinguish good and poor quality?</p> <p>Answer NO if: *Focused on data element accuracy, availability, feasibility, or other topics *Use only experts involved in development</p>	<p>→ Yes</p>	<p>5. Do the results indicate *Substantial agreement that the performance measure score from the measure as specified can be used to distinguish quality; AND *Potential threats to validity are not a problem, OR are adequately addressed so results are not biased?</p>	<p>→ Yes</p>	<p>Rate as Moderate</p>
		<p>↓No</p>			<p>→No → → → → → → → → → → →</p>	<p>Rate Low</p>
<p>↓ Yes</p>				<p>→ → → → → → → → → → → → → → → → → → → → → → →</p>		<p>Rate as Insufficient</p>

303

Applying NQF Criteria for Endorsement to eMeasures

305 EMeasures are subject to the same evaluation criteria as other performance measures. The unique
 306 aspect of eMeasures is the measure specifications, which require the health quality measure format
 307 (HQMF) and quality data model (QDM). However, these requirements pose two significant challenges.
 308 First, the HQMF and QDM may not accommodate all types or components of performance measures
 309 (e.g., PRO-PMs, risk adjustment, composites). Second, the HQMF does not prescribe where data must be
 310 located in EHRs, usually requiring additional programming to identify where the data can be found.
 311 Therefore, it may be difficult to test eMeasures to the extent necessary to meet NQF endorsement
 312 criteria—at least until they are implemented more widely. At the same time, there is interest in
 313 developing eMeasures for use in federal programs and obtaining NQF endorsement for these
 314 eMeasures. NQF endorsement may provide the impetus to implement measures; however if a
 315 submitted measure with very limited testing does not meet NQF endorsement criteria, it could be
 316 prematurely abandoned. Some other standard-setting organizations have instituted a process to
 317 approve standards for trial use; at present, such an alternative pathway may be desirable for
 318 eMeasures.

319 The following guidance first addresses the criteria for endorsement of eMeasures and then offers a
 320 proposed optional alternative pathway for those eMeasures that do not meet the requirements for
 321 reliability and validity testing. As described below, the proposed alternative pathway is NOT time-limited

322 endorsement and also is NOT the previously piloted two-stage CDP, but it does contain a few of the
323 elements of those efforts.

324 Clarification of Requirements for Endorsing eMeasures

325 The following is a consolidation and clarification of the requirements for testing eMeasures submitted to
326 NQF for endorsement (initial or endorsement maintenance). These requirements would apply to both
327 new (de novo) eMeasures and previously endorsed measures (retooled).

- 328 • EMeasures must be specified in the accepted standard of HQMF format, and must use the Quality
329 Data Model (QDM) and value sets vetted through the National Library of Medicine's Value Set
330 Authority Center (VSAC). Output from the Measure Authoring Tool (MAT) ensures that an eMeasure
331 is in the HQMF format and uses the QDM (however, the MAT is not required to produce HQMF).
332 Alternate forms of "e-specifications" other than HQMF are not considered eMeasures. *However, if*
333 *HQMF or QDM does not support all aspects of a particular measure construct, those may be*
334 *specified outside HQMF with an explanation and plans to request expansion of those standards. If a*
335 *value set not vetted by the VSAC explain why and plans to submit for approval. Please contact NQF*
336 *staff to discuss format for measure specifications that the standards do not support.*
- 337 • A new requirement for a feasibility assessment will be implemented with projects beginning after
338 July 1, 2013 (see the [eMeasure Feasibility Report](#)). The feasibility assessment addresses the data
339 elements as well as the measure logic. (See Appendix C for feasibility criteria and example
340 scorecard).
- 341 • All measures (including eMeasures) are subject to meeting the same evaluation criteria that are
342 current at the time of initial submission or endorsement maintenance (regardless of meeting prior
343 criteria and prior endorsement status). Algorithms 1, 2, and 3 apply to eMeasures.
 - 344 ○ Importance to Measure and Report (clinical evidence, performance gap, priority)
 - 345 ○ Scientific Acceptability of Measure Properties (reliability, validity)
 - 346 ○ Feasibility
 - 347 ○ Usability and Use (Accountability/transparency, improvement)
 - 348 ○ Related and competing measures
- 349 • To be considered for NQF endorsement, All measures (including eMeasures) must be tested for
350 reliability and validity using the data sources that is-are specified. Therefore, eMeasures, whether
351 new (de novo), previously respecified (retooled) but without eMeasure testing, or newly respecified,
352 must be submitted with testing using the eMeasure specifications with the specified data source
353 (e.g., EHRs, registry).
 - 354 ○ In the information provided on description of the data sample used for testing, indicate how the
355 eMeasure specifications were used to obtain the electronic data. Often eMeasures cannot be
356 directly applied to EHRs or databases from EHRs and additional programming is needed to
357 identify the location of the standardized data elements. However, in some instances, the
358 eMeasure specifications might be used directly with EHRs.
- 359 • If testing of eMeasures occurs in a small number of sites, it may be best accomplished by focusing
360 on patient-level data element validity (comparing data used in the measure to the authoritative
361 source). However, as with other measures, testing at the level of the performance measure score is
362 acceptable if data can be obtained from enough measured entities. The use of EHRs and the
363 potential access to robust clinical data provides opportunities for other approaches to testing.
 - 364 ○ If the testing is focused on validating the accuracy of the electronic data, analyze
365 agreement between the electronic data obtained using the eMeasure specifications and
366 those obtained through abstraction of the entire electronic record (not just the fields

- 367 | used to obtain the electronic data), using statistical analysis such as sensitivity and
 368 | specificity, positive predictive value, negative predictive value. The guidance on measure
 369 | testing allows this type of validity testing to also satisfy reliability of patient-level data
 370 | elements (see Algorithms 2 and 3).
- 371 | ○ Note that testing at the level of data elements requires that all critical data elements be
 372 | tested (not just agreement of one final overall computation for all patients). ~~At~~ a
 373 | minimum the numerator, denominator, and exclusions (or exceptions) must be assessed
 374 | and reported separately.
 - 375 | ○ Use of a simulated data set is no longer suggested for testing validity of data elements
 376 | and is best suited for checking that the measure specifications and logic are working as
 377 | intended.
 - 378 | ○ NQF's guidance has some flexibility; therefore, measure developers should consult with
 379 | NQF staff if they think they have another reasonable approach to testing reliability and
 380 | validity.
- 381 | ● For eMeasures, the sample for testing the patient-level data used in constructing the eMeasures
 382 | should include a **minimum of three sites, each with a different EHR system, each with three sites**
 383 | ~~(total of 9 sites). This requirement is consistent with ONC's [Office of the National Coordinator for~~
 384 | ~~Health Information Technology] requirement. Given the proposed optional path of approval for trial~~
 385 | ~~use, the HITAC subcommittee agreed that for NQF endorsement, this should be the minimum~~
 386 | ~~requirement.~~
 - 387 | ● The general guidance on samples for testing any measure also is relevant for eMeasures:
 - 388 | ○ Testing may be conducted on a sample of the accountable entities (e.g., hospital,
 389 | physician). The analytic unit specified for the particular measure (e.g., physician,
 390 | hospital, home health agency) determines the sampling strategy for scientific
 391 | acceptability testing.
 - 392 | ○ The sample should represent the variety of entities whose performance will be
 393 | measured. The Measure Testing Task Force recognized that the samples used for
 394 | reliability and validity testing often have limited generalizability because measured
 395 | entities volunteer to participate. Ideally, however, all types of entities whose
 396 | performance will be measured should be included in reliability and validity testing.
 - 397 | ○ The sample should include adequate numbers of units of measurement and adequate
 398 | numbers of patients to answer the specific reliability or validity question with the
 399 | chosen statistical method.
 - 400 | ○ When possible, units of measurement and patients within units should be randomly
 401 | selected.
 - 403 | ● The following subcriteria under Scientific Acceptability of Measure Properties also apply to
 404 | eMeasures.
 - 405 | ○ Exclusion analysis (2b3). If exclusions (or exceptions) are not based on the clinical evidence,
 406 | analyses should identify the overall frequency of occurrence of the exclusions as well as
 407 | variability across the measured entities to demonstrate the need to specify exclusions.
 - 408 | ○ Risk adjustment (2b4). Outcome and resource use measures require testing of the risk
 409 | adjustment approach.
 - 410 | ○ Differences in performance (2b5). This criterion is about using the measure as specified to
 411 | distinguish differences in performance across the entities that are being measured. The
 412 | performance measure scores should be computed for all accountable entities for which ~~you~~
 413 | have eMeasure data are available (not just those on which reliability/validity testing was
 414 | conducted) and then analyzed to identify differences in performance.

Formatted

- 415 ○ ~~eMeasures should be submitted as a separate measure even if the same or similar measure~~
416 ~~exists for another data source (e.g., claims). Therefore, Comparability of performance~~
417 ~~score/measure scores~~ if specified for multiple data sources (2b6) ~~would not apply. (e.g., EHRs,~~
418 ~~claims). (NQF will explore alternatives for linking measures that are the same except for data~~
419 ~~source). If a performance measure is specified for more than one data source, it should be tested~~
420 ~~with each. Unless empirical analyses demonstrate comparability of computed scores computed,~~
421 ~~assume non-comparability and submit as separate measures. Measures are endorsed only for~~
422 ~~the data specificatins and levels of analyses for which they are tested.~~ The measures specified
423 for different data sources will be evaluated as competing measures ~~(unless they apply to~~
424 ~~different care settings or levels of analysis)~~ to determine whether one is superior to the other or
425 whether there is justification for endorsing multiple measures.
426 ○ Analysis of missing data (2b7). Approved recommendations from the 2012 projects on
427 eMeasure feasibility assessment, composites, and patient-reported outcomes call for an
428 assessment of missing data or nonresponses.

429 ~~The HITAC and CSAC are interested in receiving requested comments on the minimum number of testing~~
430 ~~sites when conducting validity testing of the data elements for eMeasures and whether a similar~~
431 ~~requirement should apply to all measures. Most commenters did not agree with setting a minimum for~~
432 ~~testing of 3EHRs with 3 sites each. They cited burden with finding sites and costs and thought it~~
433 ~~represented a “higher bar” for endorsement. Kevin Larsen clarified that ONC requires a minimum of 3~~
434 ~~EHR systems, but no minimum number of sites. The CSAC and HITAC agreed to make the NQF~~
435 ~~requirement consistent with ONC and require 3 EHRs as reflected above.~~

436 Proposed Approval for of Trial Use for eMeasures for Implementation and Testing

437 This optional path of **approval for trial use** is intended for eMeasures that are ready for implementation
438 but cannot yet be adequately tested to meet NQF endorsement criteria. For such eMeasures, NQF
439 proposes to utilize the multi-stakeholder consensus process to evaluate and approve eMeasures for trial
440 use that address important areas for performance measurement and quality improvement, though they
441 may not have the requisite testing needed for NQF endorsement. These eMeasures must be assessed to
442 be technically acceptable for implementation. The goal of approving eMeasures for trial use is to
443 promote implementation and the ability to conduct more robust reliability and validity testing that can
444 take advantage of the clinical data in EHRs.

445 Approval for trial use is NOT time-limited endorsement as it carries no endorsement label. Also, this is
446 not a required two-stage review process: eMeasures that meet endorsement criteria do not need to
447 first go through an approval for trial use.

448 To be clear, eMeasures that are approved by NQF for trial use would differ from eMeasures that are
449 endorsed.

450 *NQF Endorsement* means that the eMeasure has been judged to meet all NQF evaluation criteria
451 and is suitable for use in accountability applications as well as performance improvement.

452 *NQF Approval ~~for of~~ Trial Measures Use for Implementation and Testing* means the eMeasure
453 has been judged to meet the criteria ~~indicating that indicate its readiness for implementation it~~
454 ~~is ready to be implemented~~ in real-world settings ~~in order~~ to generate the data required to
455 assess reliability and validity ~~in the future. It Such measures~~ also could be used for internal

456 | performance improvement. However, ~~it such measures has not yet~~would not have been judged
457 | to meet all the criteria indicating it is suitable for use in accountability applications.

458 | Criteria for Approval for Trial Use

- 459 | • Such measures will be considered Approved ~~for as~~ Trial Measures for Use Implementation and
460 | Testing, NOT Endorsed
461 | • When sufficient data have been accumulated for adequate reliability and validity testing, the
462 | eMeasure can be submitted to NQF for potential endorsement (not all may progress to
463 | endorsement).
464 |

465 | Criteria for Approval ~~for of~~ Trial Measures for Use Implementation and Testing

- 466 | ~~• The following are the proposed requirements for Approval for Trial Use:~~
467 | • Must be eMeasures, meaning the measures are specified in the accepted standard of HQMF
468 | format, and must use the Quality Data Model (QDM). Output from the Measure Authoring Tool
469 | (MAT) ensures that an eMeasure is in the HQMF format and uses the QDM (however, the MAT
470 | is not required to produce HQMF). Alternate forms of “e-specifications” other than HQMF are
471 | not considered eMeasures. *However, if HQMF or QDM does not support all aspects of a*
472 | *particular measure construct, those may be specified outside HQMF. Please contact NQF staff to*
473 | *discuss format for measure specifications.)*
474 | • Must use value sets vetted through the National Library of Medicine’s Value Set Authority
475 | Center. This will help ensure appropriate use of codes and code systems and will help minimize
476 | value set harmonization issues in submitted eMeasures. If particular value sets are not vetted by
477 | the VSAC, explain why they are used in the measure and describe plans to submit them to VSAC
478 | for approval.
479 | • Must meet all criteria under Importance to ~~m~~Measure and ~~r~~Report (clinical evidence,
480 | performance gap, priority).
481 | • The feasibility assessment must be completed.
482 | • Results from testing with a simulated (or test) data set demonstrate that the QDM and HQMF
483 | are used appropriately and that the measure logic performs as expected.
484 | • There is a plan for use and discussion of how the measure will be useful for accountability and
485 | improvement.
486 | • Related and competing measures are identified with plan for harmonization or justification for
487 | developing a competing measure.

488 | Process

- 489 | • Measures submitted for approval as Trial Measures for Implementation and Testing will be
490 | submitted and evaluated as other measures that are submitted for endorsement using the
491 | multistakeholder Consensus Development Process (CDP).
492 | • Measures submitted as Trial Measures will be evaluated on Importance to Measure and Report,
493 | eMeasure specifications, Feasibility, Usability and Use (planned), and Related and Competing
494 | Measures (identified, harmonized, or justification for developing competing measure).
495 | • Trial Measure designation automatically expires 3 years after approved if not submitted for
496 | endorsement prior to that time.
497 | o The time to submit for endorsement is driven by success with testing. There is no
498 | expectation that every trial measure will be submitted for endorsement – some may fail
499 | during testing.

Formatted: Heading 4, Space After: 0 pt

- A developer could request an extension of status as an approved Trial Measure if they document that they are actively engaged in the process of testing.
- When submitted for endorsement, the measure will be evaluated through the multistakeholder process. Ideally, standing committees and/or more flexible schedules for submitting measures will prevent delays for the endorsement process.
- If submitted for endorsement prior to the 3-year expiration, the developer can select from the following options for evaluation and endorsement:
 - Option 1: Submit and evaluate only Scientific Acceptability of Measure Properties, including the final eMeasure specifications and all testing. If endorsed, endorsement maintenance will be due **3 years from the date approved as a trial measure**, at which time it will be submitted for endorsement maintenance and subject to evaluation on all criteria.
 - Option 2: Submit and evaluate on all criteria. If endorsed, a new endorsement date will be identified and endorsement maintenance will be due 3 years from the new endorsement date, at which time it will be submitted for endorsement maintenance and subject to evaluation on all criteria.
- If submitted for endorsement three or more years after the date of approval as a trial measure, the measure must be submitted and evaluated on all criteria just as any measure being submitted for initial endorsement or maintenance.

Formatted: Line spacing: single, Bulleted + Level: 2 + Aligned at: 0.75" + Indent at: 1"

Formatted: Font:

Formatted: Line spacing: single, Bulleted + Level: 2 + Aligned at: 0.75" + Indent at: 1"

Formatted: Font: Bold, Italic

Formatted: Font: Bold, Italic

Formatted: Font: Bold, Italic

Formatted: Font: Bold, Italic

Formatted: List Paragraph, Space After: 11 pt, Bulleted + Level: 2 + Aligned at: 0.75" + Indent at: 1"

Formatted: List Paragraph, Bulleted + Level: 1 + Aligned at: 0.25" + Indent at: 0.5"

Formatted: Font:

519 *The CSAC and HITAC are especially interested in receiving comments about approval for trial use.*

Discussion – Trial Measures

- On the 9/10 CSAC/HITAC call, the question was raised of whether NQF approval as Trial Measure will facilitate implementation and testing. If not, then this path could delay use of eMeasures.
- Other questions were raised regarding use of Trial Measures in federal programs and whether HHS contracts for measure development could accommodate development and submission for approval as Trial Measures vs. endorsement. Trial Measures would not be tested and thus do not meet endorsement criteria as suitable for accountability applications until satisfactory testing was submitted.
- Commenters suggested approval as Trial Measures be available for any measure, not just eMeasures. NQF has been phasing out “time-limited endorsement” for untested measures, but still grants occasionally.
- A more flexible approach to accepting measure submissions so that they could be submitted whenever ready may minimize the need for such an interim pathway.

Key Questions

- Should NQF proceed with a pilot for approval of trial measures as outlined in the criteria and process above? (For example, now through June 2014)?
- Do you have any suggested changes to criteria or process?
- Should it be limited to eMeasures or open for all types of measures?
- If open to all types of measures is there a need for any additional criteria as exist for time-limited endorsement (time-sensitive legislative mandate, not competing with existing measure, not complex – i.e., outcome, resource use, composite). Because approval as Trial Measure does not imply endorsement along with fact that measures will be evaluated on Importance to Measure and Report, planned use, and identification/justification for related or competing measures these additional criteria probably are not necessary.

Formatted Table

1. Evidence, Performance Gap, and Priority (Impact)—Importance to Measure and Report

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 (high-impact) aspect of healthcare where there is variation in or overall less-than-optimal performance. *Measures must be judged to meet all subcriteria to pass this criterion and be evaluated against the remaining criteria.* Yes No

1a. Evidence to Support the Measure Focus Yes No

Quantity: H M L I Quality: H M L I Consistency: H M L I

The measure focus is evidence-based, demonstrated as follows:

- **Health outcome:**³ a rationale supports the relationship of the health outcome to processes or structures of care. Applies to patient-reported outcomes (PRO), including health-related quality of life/functional status, symptom/symptom burden, experience with care, health-related behavior.
- **Intermediate clinical outcome:** a systematic assessment and grading of the quantity, quality, and consistency of the body of evidence⁴ that the measured intermediate clinical outcome leads to a desired health outcome.
- **Process:**⁵ a systematic assessment and grading of the quantity, quality, and consistency of the body of evidence⁴ that the measured process leads to a desired health outcome.
- **Structure:** a systematic assessment and grading of the quantity, quality, and consistency of the body of evidence⁴ that the measured structure leads to a desired health outcome.
- **Efficiency:**⁶ evidence not required for the resource use component.

AND

1b. Performance Gap H M L I

Demonstration of quality problems and opportunity for improvement, i.e., data⁷ demonstrating

- considerable variation, or overall less-than-optimal performance, in the quality of care across providers; and/or
- disparities in care across population groups.

AND

1c. High Priority (previously referred to as High Impact) H M L I

The measure addresses:

- a specific national health goal/priority identified by DHHS or the National Priorities Partnership convened by NQF;

OR

- a demonstrated high-priority (high-impact) 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).
- For patient-reported outcomes, there is evidence that the target population values the PRO and finds it meaningful.

1d. For composite performance measures, the following must be explicitly articulated and logical:

1d1. The quality construct, including the overall area of quality; included component measures; and the

relationship of the component measures to the overall composite and to each other; and
1d2. The rationale for constructing a composite measure, including how the composite provides a distinctive or additive value over the component measures individually; and
1d3. How the aggregation and weighting of the component measures are consistent with the stated quality construct and rationale.

Notes

3. 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.
4. The preferred systems for grading the evidence are the U.S. Preventive Services Task Force (USPSTF) [grading definitions](#) and [methods](#), or Grading of Recommendations, Assessment, Development and Evaluation ([GRADE](#)) [guidelines](#).
5. 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. Note: A measure focused only on collecting PROM data is not a PRO-PM.
6. Measures of efficiency combine the concepts of resource use and quality (NQF's [Measurement Framework: Evaluating Efficiency Across Episodes of Care](#); [AQA Principles of Efficiency Measures](#)).
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.

523 **Table A-1: Evaluation of Quantity, Quality, and Consistency of Body of Evidence for Structure, Process,**
 524 **and Intermediate Outcome Measures**

DEFINITION/ RATING	QUANTITY OF BODY OF EVIDENCE	QUALITY OF BODY OF EVIDENCE	CONSISTENCY OF RESULTS OF BODY OF EVIDENCE
Definition	Total number of studies (not articles or papers)	Certainty or confidence in the estimates of benefits and harms to patients across studies in the body of evidence related to study factors^a including: study design or flaws; directness/indirectness to the specific measure (regarding the population, intervention, comparators, outcomes); imprecision (wide confidence intervals due to few patients or events)	Stability in both the direction and magnitude of clinically/practically meaningful benefits and harms to patients (benefit over harms) across studies in the body of evidence
High	5+ studies ^b	Randomized controlled trials (RCTs) providing direct evidence for the specific measure focus, with adequate size to obtain precise estimates of effect, and without serious flaws that introduce bias	Estimates of clinically/practically meaningful benefits and harms to patients are consistent in direction and similar in magnitude across the preponderance of studies in the body of evidence
Moderate	2-4 studies ^b	<ul style="list-style-type: none"> • Non-RCTs with control for confounders that could account for other plausible explanations, with large, precise estimate of effect OR • RCTs without serious flaws that introduce bias, but with either indirect evidence or imprecise estimate of effect 	<p>Estimates of clinically/practically meaningful benefits and harms to patients are consistent in direction across the preponderance of studies in the body of evidence, but may differ in magnitude</p> <p>If only one study, then the estimate of benefits greatly outweighs the estimate of potential harms to patients (one study cannot achieve high consistency rating)</p>
Low	1 study ^b	<ul style="list-style-type: none"> • RCTs with flaws that introduce bias OR • Non-RCTs with small or imprecise estimate of effect, or without control for confounders that could account for other plausible explanations 	<ul style="list-style-type: none"> • Estimates of clinically/practically meaningful benefits and harms to patients differ in both direction and magnitude across the preponderance of studies in the body of evidence OR • wide confidence intervals prevent estimating net benefit <p>If only one study, then estimate of benefits do not greatly outweigh harms to patients</p>

DEFINITION/ RATING	QUANTITY OF BODY OF EVIDENCE	QUALITY OF BODY OF EVIDENCE	CONSISTENCY OF RESULTS OF BODY OF EVIDENCE
Insufficient to Evaluate (See Table 3 for exceptions.)	<ul style="list-style-type: none"> No empirical evidence OR Only selected studies from a larger body of evidence 	<ul style="list-style-type: none"> No empirical evidence OR Only selected studies from a larger body of evidence 	No assessment of magnitude and direction of benefits and harms to patients

525 ^aStudy designs that affect certainty of confidence in estimates of effect include: randomized controlled trials (RCTs), which
526 control for both observed and unobserved confounders, and non-RCTs (observational studies) with various levels of control for
527 confounders.
528 Study flaws that may bias estimates of effect include: lack of allocation concealment; lack of blinding; large losses to follow-up;
529 failure to adhere to intention to treat analysis; stopping early for benefit; and failure to report important outcomes.
530 Imprecision with wide confidence intervals around estimates of effects can occur in studies involving few patients and few
531 events.
532 Indirectness of evidence includes: indirect comparisons (e.g., two drugs compared to placebos rather than head-to head); and
533 differences between the population, intervention, comparator interventions, and outcome of interest and those included in the
534 relevant studies.¹⁵
535 ^bThe suggested number of studies for rating levels of quantity is considered a general guideline.

536 **Table A-2: Evaluation of Subcriterion 1a Based on the Quantity, Quality, and Consistency of the Body**
537 **of Evidence**

QUANTITY OF BODY OF EVIDENCE	QUALITY OF BODY OF EVIDENCE	CONSISTENCY OF RESULTS OF BODY OF EVIDENCE	PASS SUBCRITERION 1A
Moderate-High	Moderate-High	Moderate-High	Yes
Low	Moderate-High	Moderate (if only one study, high consistency not possible)	Yes, but only if it is judged that additional research is unlikely to change conclusion that benefits to patients outweigh harms; otherwise, No
Moderate-High	Low	Moderate-High	Yes, but only if it is judged that potential benefits to patients clearly outweigh potential harms; otherwise, No
Low-Moderate-High	Low-Moderate-High	Low	No
Low	Low	Low	No
Exception to Empirical Body of Evidence for Health Outcome For a health outcome measure: A rationale supports the relationship of the health outcome to at least one healthcare structure, process, intervention, or service			Yes, if it is judged that the rationale supports the relationship of the health outcome to at least one healthcare structure, process, intervention, or service
Potential Exception to Empirical Body of Evidence for Other Types of Measures If there is no empirical evidence, expert opinion is systematically assessed with agreement that the benefits to patients greatly outweigh potential harms.			Yes, but only if it is judged that potential benefits to patients clearly outweigh potential harms; otherwise, No

538

539 **Appendix B: Current Criteria and Guidance related to Reliability and Validity**
540 **Testing**

2. Reliability and Validity—Scientific Acceptability of Measure Properties

Extent to which the measure, as specified, produces consistent (reliable) and credible (valid) results about the quality of care when implemented. **Measures must be judged to meet the subcriteria for both reliability and validity to pass this criterion and be evaluated against the remaining criteria.**

Yes No

2a. Reliability H M L I

2a1. The measure is well defined and precisely specified ⁸ so it can be implemented consistently within and across organizations and allows for comparability. eMeasures should be specified in the Health Quality Measures Format (HQMF) and the quality data model (QDM). ²

2a2. Reliability testing ¹⁰ 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. For **PRO-PMs and composite performance measures**, reliability should be demonstrated for the computed performance score.

2b. Validity H M L I

2b1. The measure specifications ⁸ are consistent with the evidence presented to support the focus of measurement under criterion 1a. The measure is specified to capture the most inclusive target population indicated by the evidence, and exclusions are supported by the evidence.

2b2. Validity testing ¹¹ 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. . For **PRO-PMs and composite performance measures**, reliability should be demonstrated for the computed performance score.

2b3. Exclusions are supported by the clinical evidence; otherwise, they are supported by evidence of sufficient frequency of occurrence so that results are distorted without the exclusion; ¹²

AND

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). ¹³

2b4. For outcome measures and other measures when indicated (e.g., resource use):

- 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; ^{14,15} and has demonstrated adequate discrimination and calibration

OR

- rationale/data support no risk adjustment/ stratification.

2b5. 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 ¹⁶ differences in performance;

OR

there is evidence of overall less-than-optimal performance.

2b6. If multiple data sources/methods are specified, there is demonstration they produce comparable results.

2b7. For **eMeasures, composites, and PRO-PMs** (or other measures susceptible to missing data), analyses identify the extent and distribution of missing data (or nonresponse) and demonstrate that performance results are not biased due to systematic missing data (or differences between responders and nonresponders) and how the specified handling of missing data minimizes bias.

2c. Disparities (*Disparities should be addressed under subcriterion 1b*)

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);

2d. For composite performance measures, empirical analyses support the composite construction approach and demonstrate the following:

2d1. the component measures fit the quality construct and add value to the overall composite while achieving the related objective of parsimony to the extent possible; and

2d2. the aggregation and weighting rules are consistent with the quality construct and rationale while achieving the related objective of simplicity to the extent possible.

(if not conducted or results not adequate, justification must be submitted and accepted)

Notes

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.

Specifications for **PRO-PMs** also include: specific PROM(s); standard methods, modes, and languages of administration; whether (and how) proxy responses are allowed; standard sampling procedures; handling of missing data; and calculation of response rates to be reported with the performance measure results.

Specifications for **composite performance measures** include: component measure specifications (unless individually endorsed); aggregation and weighting rules; handling of missing data; standardizing scales across component measures; required sample sizes.

9. Some eMeasures may have aspects that cannot be specified in HQMF or the QDM (e.g., risk adjustment, composite aggregation and weighting rules) and should be specified in HQMF and QDM to the extent possible, with the additional specifications and an explanation why cannot be represented in HQMF or QDM. eMeasure specifications include data type from the QDM, value sets and attributes, measure logic, original source of the data and recorder.

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 measures 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.

541

542 Table B-1: Evaluation Ratings for Reliability and Validity

RATING	RELIABILITY	VALIDITY
<p>High</p>	<p>All measure specifications (e.g., numerator, denominator, exclusions, risk factors, scoring, etc.) are unambiguous and likely to consistently identify who is included and excluded from the target population and the process, condition, event, or outcome being measured; how to compute the score, etc.;</p> <p>AND</p> <p>Empirical evidence of reliability of <u>BOTH data elements AND computed performance measure score within acceptable norms:</u></p> <ul style="list-style-type: none"> • Data element: appropriate method, scope, and reliability statistics for critical data elements within acceptable norms (new testing, or prior evidence for the same data type); • OR commonly used data elements for which reliability can be assumed (e.g., gender, age, date of admission); • OR <i>may forego data element reliability testing if data element validity was demonstrated;</i> <p>AND</p> <ul style="list-style-type: none"> • Performance measure score: appropriate method, scope, and reliability statistic within acceptable norms 	<p>The measure specifications (numerator, denominator, exclusions, risk factors) are consistent with the evidence cited in support of the measure focus (1a) under <i>Importance to Measure and Report</i>;</p> <p>AND</p> <p>Empirical evidence of validity of <u>BOTH data elements AND computed performance measure score within acceptable norms:</u></p> <ul style="list-style-type: none"> • Data element: appropriate method, scope, and statistical results within acceptable norms (new testing, or prior evidence for the same data type) for critical data elements; <p>AND</p> <ul style="list-style-type: none"> • Performance measure score: appropriate method, scope, and validity testing result within acceptable norms; <p>AND</p> <p>Identified threats to validity (lack of risk adjustment/stratification, multiple data types/methods, systematic missing or “incorrect” data) are empirically assessed and adequately addressed so that results are not biased</p>
<p>Moderate</p>	<p>All measure specifications are unambiguous as noted above</p> <p>AND</p> <p>Empirical evidence of reliability <u>within acceptable norms</u> for <u>either critical data elements OR performance measure score</u> as noted above</p>	<p>The measure specifications reflect the evidence cited under <i>Importance to Measure and Report</i> as noted above;</p> <p>AND</p> <p>Empirical evidence of validity <u>within acceptable norms</u> for <u>either critical data elements OR performance measure score</u> as noted above; OR</p> <p>Systematic assessment of face validity of performance <u>measure score</u> as a quality indicator explicitly addressed and found substantial agreement that <i>the scores obtained from the measure as specified will provide an accurate reflection of quality and can be used to distinguish good and poor quality</i></p> <p>AND</p> <p>Identified threats to validity noted above are empirically assessed and adequately addressed so that results are not biased</p>

RATING	RELIABILITY	VALIDITY
Low	<p>One or more measure specifications (e.g., numerator, denominator, exclusions, risk factors, scoring) are <u>ambiguous</u> with potential for confusion in identifying who is included and excluded from the target population, or the event, condition, or outcome being measured; or how to compute the score, etc.;</p> <p>OR</p> <p>Empirical evidence (using appropriate method and scope) of <u>unreliability</u> for <u>either data elements OR performance measure score</u>, i.e., statistical results outside of acceptable norms</p>	<p>The measure specifications <u>do not</u> reflect the evidence cited under <i>Importance to Measure and Report</i> as noted above;</p> <p>OR</p> <p>Empirical evidence (using appropriate method and scope) of <u>invalidity</u> for <u>either data elements OR performance measure score</u>, i.e., statistical results outside of acceptable norms</p> <p>OR</p> <p>Identified threats to validity noted above are empirically assessed and determined to bias results</p>
Insufficient Evidence	<p>Inappropriate method or scope of reliability testing</p>	<p>Inappropriate method or scope of validity testing (including inadequate assessment of face validity as noted above);</p> <p>OR</p> <p>Threats to validity as noted above are likely and are NOT empirically assessed</p>

543

544 **Table B-2: Evaluation of Reliability and Validity of eMeasures**

RATING	RELIABILITY DESCRIPTION AND EVIDENCE	VALIDITY DESCRIPTION AND EVIDENCE
High	<p>Specified in HQMF and QDM and all specifications are unambiguous⁺; AND</p> <p>Empirical evidence of reliability of <u>both data element AND computed performance measure score within acceptable norms</u>:</p> <ul style="list-style-type: none"> • <u>Data element</u>: reliability (repeatability) assured with computer programming—must test data element validity <p>AND</p> <ul style="list-style-type: none"> • <u>Performance measure score</u>: appropriate method, scope, and reliability statistic within acceptable norms 	<p>The measure specifications (numerator, denominator, exclusions, risk factors) are consistent with the evidence cited in support of the measure focus (1a) under <i>Importance to Measure and Report</i>;</p> <p>AND</p> <p>Empirical evidence of validity of <u>both data elements AND computed performance measure score within acceptable norms</u>:</p> <ul style="list-style-type: none"> • <u>Data element</u>: validity demonstrated by analysis of agreement between data elements obtained using the eMeasure as specifications and data elements manually abstracted from the <u>entire</u> EHR with statistical results within acceptable norms; OR complete agreement between data elements and computed measure scores obtained by applying the eMeasure specifications to a simulated test EHR data set with known values for the critical data elements; <p>AND</p> <ul style="list-style-type: none"> • <u>Performance measure score</u>: appropriate method, scope, and validity testing result within acceptable norms; <p>AND</p> <p>Identified threats to validity (lack of risk adjustment/stratification, multiple data types/methods, systematic missing or “incorrect” data) are empirically assessed and adequately addressed so that results are not biased</p>
Moderate	<p>Specified in HQMF and QDM and all specifications are unambiguous⁺ and include only data elements from the QDM;* OR new data elements are submitted for inclusion in the QDM; AND</p> <p>Empirical evidence of reliability <u>within acceptable norms for either data elements OR performance measure score</u> as noted above</p>	<p>The measure specifications reflect the evidence cited under <i>Importance to Measure and Report</i> as noted above;</p> <p>AND</p> <p>Empirical evidence of validity <u>within acceptable norms for either data elements OR performance measure score</u> as noted above; OR</p> <p>Systematic assessment of face validity of <u>performance measure score</u> as a quality indicator explicitly addressed and found substantial agreement that <i>the scores obtained from the measure as specified will provide an accurate reflection of quality and can be used to distinguish good and poor quality</i></p> <p>AND</p> <p>Identified threats to validity noted above are empirically assessed and adequately addressed so that results are not biased</p>

RATING	RELIABILITY DESCRIPTION AND EVIDENCE	VALIDITY DESCRIPTION AND EVIDENCE
Low	One or more eMeasure specifications are ambiguous ⁺ or <u>do not</u> use HQMF and QDM*; OR Empirical evidence of <u>unreliability</u> for <u>either data elements OR performance measure score</u> —i.e., statistical results outside of acceptable norms	The EHR measure specifications do not reflect the evidence cited under <i>Importance to Measure and Report</i> as noted above; OR Empirical evidence (using appropriate method and scope) of <u>invalidity</u> for <u>either data elements OR performance measure score</u> — i.e., statistical results outside of acceptable norms OR Identified threats to validity noted above are empirically assessed and determined to bias results
Insufficient evidence	Inappropriate method or scope of reliability testing	Inappropriate method or scope of validity testing (including inadequate assessment of face validity as noted above) OR Threats to validity as noted above are likely and are NOT empirically assessed

545 ⁺Specifications are considered unambiguous if they are likely to consistently identify who is included and excluded from the
546 target population and the process, condition, event, or outcome being measured; how to compute the score, etc.

547 *HQMF and QDM (formerly called the QDS) should be used when available. When quality data elements are needed but are
548 not yet available in the QDM, they will be considered for addition to the QDM.

549 **Table B-3: Evaluation of Scientific Acceptability of Measure Properties Based on Reliability and**
550 **Validity Ratings**

VALIDITY RATING	RELIABILITY RATING	PASS SCIENTIFIC ACCEPTABILITY OF MEASURE PROPERTIES FOR INITIAL ENDORSEMENT*	
		Yes	No
High	Moderate-High	Yes	Evidence of reliability and validity
	Low	No	Represents inconsistent evidence—reliability is usually considered necessary for validity
Moderate	Moderate-High	Yes	Evidence of reliability and validity
	Low	No	Represents inconsistent evidence—reliability is usually considered necessary for validity
Low	Any rating	No	Validity of conclusions about quality is the primary concern. If evidence of validity is rated low, the reliability rating will usually also be low. Low validity and moderate-high reliability represents inconsistent evidence.

551 *A measure that does not pass the criterion of *Scientific Acceptability of Measure Properties* would not be recommended for
552 endorsement.
553

554 **Table B-4: Scope of Testing Required at the Time of Review for Endorsement Maintenance**

	FIRST ENDORSEMENT MAINTENANCE REVIEW	SUBSEQUENT REVIEWS
Reliability	<p>Measure In Use</p> <ul style="list-style-type: none"> • Analysis of data from entities whose performance is measured • Reliability of measure scores (e.g., signal to noise analysis) <p>Measure Not in Use</p> <ul style="list-style-type: none"> • Expanded testing in terms of scope (number of entities/patients) and/or levels (data elements/measure score) 	Could submit prior testing data, if results demonstrated that reliability achieved a high rating
Validity	<p>Measure in Use</p> <ul style="list-style-type: none"> • Analysis of data from entities whose performance is measured • Validity of measure score for making accurate conclusions about quality • Analysis of threats to validity <p>Measure Not in Use</p> <ul style="list-style-type: none"> • Expanded testing in terms of scope (number of entities/patients) and/or levels (data elements/measure score) 	Could submit prior testing data, if results demonstrated that validity achieved a high rating

555

556

Appendix C: Current Criteria and Guidance related to Feasibility

557

3. Feasibility:

Extent to which the specifications, including measure logic, require data that are readily available or could be captured without undue burden and can be implemented for performance measurement.

H M L I

3a. For clinical measures, the required data elements are routinely generated and used during care delivery (e.g., blood pressure, lab test, diagnosis, medication order). H M L I

3b. 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. H M L I

3c. Demonstration that the data collection strategy (e.g., data source/availability, timing, frequency, sampling, patient-reported data, patient confidentiality, ¹⁷ costs associated with fees/licensing for proprietary measures or elements such as risk model, grouper, instrument) can be implemented (e.g., already in operational use, or testing demonstrates that it is ready to put into operational use). For eMeasures, a feasibility assessment addresses the data elements and measure logic ¹⁸ and demonstrates the eMeasure can be implemented or feasibility concerns can be adequately addressed.

H M L I

Note

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.

18. The feasibility assessment uses a standard score card or a fully transparent alternative that includes at a minimum: a description of the assessment, feasibility scores for all data elements, and explanatory notes for all data element components scoring a “1” (lowest rating); measure logic can be executed; with rationale and plan for addressing feasibility concerns.

558

Guidance on Evaluating Feasibility

559

560

Table C-1: Generic Scale for Rating Feasibility Subcriteria

RATING	DEFINITION
High	Based on the information submitted, there is high confidence (or certainty) that the criterion is met
Moderate	Based on the information submitted, there is moderate confidence (or certainty) that the criterion is met
Low	Based on the information submitted, there is low confidence (or certainty) that the criterion is met
Insufficient	There is insufficient information submitted to evaluate whether the criterion is met (e.g., blank, incomplete, or not relevant, responsive, or specific to the particular question)

561

562 Table C-2. Data Element Feasibility Scorecard

DATA ELEMENT:			
Measure Title:			
Data element definition:			
Who performed the assessment:			
Type of setting or practice, i.e., solo practice, large group, academic hospital, safety net hospital, integrated system:			
EHR system used:			
	Current (1-3)	Future* (1-3)	Comments
Data Availability – Is the data readily available in structured format? Scale: 3 – Data element exists in structured format in this EHR. [2] – Not defined as this time. Hold for possible future use. 1 – Data element is not available in structured format in this EHR.			
Data Accuracy – Is the information contained in the data element correct? Are the data source and recorder specified? Scale: 3 – The information is from the most authoritative source and/or is highly likely to be correct. (e.g., laboratory test results transmitted directed from the laboratory information system into the EHR). 2 – The information may not be from the most authoritative source and/or has a moderate likelihood of being correct. (e.g., self-report of a vaccination). 1 – The information may not be correct. (e.g., a check box that indicates medication reconciliation was performed).			
Data Standards – Is the data element coded using a nationally accepted terminology standard? Scale: 3 – The data element is coded in nationally accepted terminology standard. 2 – Terminology standards for this data element are currently available, but is not consistently coded to standard terminology in the EHR, or the EHR does not easily allow such coding. 1 – The EHR does not support coding to the existing standard.			

DATA ELEMENT:			
<p>Workflow – To what degree is the data element captured during the course of care? How does it impact the typical workflow for that user?</p> <p>Scale:</p> <p>3 – The data element is routinely collected as part of routine care and requires no additional data entry from clinician solely for the quality measure and no EHR user interface changes. Examples would be lab values, vital signs, referral orders, or problem list entry.</p> <p>2 – Data element is not routinely collected as a part of routine care and additional time and effort over and above routine care is required, but perceived to have some benefit.</p> <p>1 – Additional time and effort over and above routine care is required to collect this data element without immediate benefit to care</p>			
DATA ELEMENT FEASIBILITY SCORE			

563
564

*For data elements that score low on current feasibility, indicate the anticipated feasibility score in 3-5 years based on a projection of the maturation of the EHR, or maturation of its use.

565 **Appendix D: Project Steering Committee and NQF Staff**

566 **Consensus Standards Approval Committee Member Roster**

567 **Frank Opelka, MD, FACS (Chair) ***

568 Vice President for Health Affairs and Medical Education
569 Louisiana State University, New Orleans, LA

570 **Cristie Upshaw Travis (Vice-Chair) ***

571 Chief Executive Officer
572 Memphis Business Group on Health, Memphis, TN

573 **Andrew Baskin, MD ***

574 National Medical Director for Quality and Provider Performance Measurement
575 Aetna, Blue Bell, PA

576 **Pamela Cipriano, PhD, RN NEA-BC, FAAN**

577 Senior Director
578 Galloway Consulting, Marietta, GA

579 **William Conway, MD**

580 Senior Vice President and Chief Quality Officer
581 Henry Ford Health System, Detroit, MI

582 **Robert Ellis**

583 Director of Operations and Online Services
584 Consumers' Checkbook, Ashburn, VA

585 **Lee Fleisher, MD ***

586 Robert D. Dripps Professor and Chair of Anesthesiology and Critical Care
587 University of Pennsylvania, Philadelphia, PA

588 **David Knowlton, MA**

589 President and Chief Executive Officer
590 The New Jersey Health Care Quality Institute, Pennington, NJ

591 **Philip E. Mehler, MD**

592 Chief Medical Officer and Director of Quality
593 Denver Health, Denver, CO

594 **Ann Monroe**

595 President
596 Health Foundation for Western & Central New York, Buffalo, NY

597 **Arden Morris, MD, MPH, FACS**

598 Associate Professor of Surgery
599 University of Michigan Health System, Ann Arbor, MI

- 600 **Lyn Paget, MPH**
601 Managing Partner
602 Health Policy Partners, Boston, MA
- 603 **Carolyn Pare**
604 President and Chief Executive Officer
605 Buyers Health Care Action Group, Bloomington, MI
- 606 **Lee Partridge ***
607 Senior Health Policy Advisor
608 National Partnership for Women & Families, Washington, DC
- 609 **Kyu Rhee, MD, MPP**
610 Vice President of Integrated Health Services
611 IBM Corporation, Somers, NY
- 612 **David Rhew, MD**
613 Chief Medical Officer and VP of Global Healthcare
614 Samsung SDS America, Moonachie, NJ
- 615 **Dana Gelb Safran, ScD ***
616 Senior Vice President for Performance Measurement and Improvement
617 Blue Cross Blue Shield of Massachusetts, Boston, MA
- 618 **David Shahian ***
619 Chair of 2010 Evidence Task Force
620 Consultant Surgeon
621 Massachusetts General Hospital (MGH), Boston, MA
- 622 *** Participated on subcommittee**

623 **Health Information Technology Advisory Committee Roster**

624 **Paul C. Tang, MD, MS (Chair) ***

625 Vice President and Chief Medical Information Officer
626 Palo Alto Medical Foundation, Palo Alto, CA

627 **J. Marc Overhage, MD, PhD (Vice-Chair) ****

628 Chief Medical Informatics Officer
629 Siemens Healthcare, USA, Malvern, PA

630 **Kristine Martin Anderson, MBA ****

631 Senior Vice President
632 BoozAllenHamilton, Rockville, MD

633 **David W. Bates, MD, MSc**

634 Medical Director of Clinical and Quality Analysis
635 Partners Healthcare System, Inc., Boston, MA

636 **Zahid Butt, MD, FACC ***

637 President and CEO, Medisol, Inc.
638 Columbia, MD

639 **Ian Z. Chuang, MD, MS ***

640 Senior Vice President, Healthcare Informatics, and Chief Medical Officer
641 Netsmart, Overland Park, KS

642 **John Derr, RPh**

643 Health Information Strategy Consultant, Golden Living, LLC
644 Anacortes, WA

645 **Richard Dutton, MD, MBA ***

646 Executive Director, Anesthesia Quality Institute
647 Park Ridge, IL

648 **Jamie Ferguson ***

649 Vice President, Health Information Technology Strategy & Policy
650 Kaiser Permanente, Oakland, CA

651 **Paul Fu, MD, MPH**

652 Chief Medical Information Officer, Harbor - UCLA Medical Center
653 Torrance, CA

654 **Leslie Kelly Hall**

655 Senior Vice President
656 Healthwise, Inc., Boise, Idaho

657 **Allison Jackson, MS**
658 Project Manager/Epidemiologist
659 Intel Corporation, Chandler, AZ

660 **Caterina E.M. Lasome, PhD, MSN, MBA, MHA, RN, CPHIMS**
661 President and CEO
662 iON Informatics, LLC, Dunn Loring, VA

663 **Russell Leftwich, MD ***
664 Chief Medical Informatics Officer, Office of eHealth Initiatives
665 State of Tennessee, Nashville, TN

666 **Michael Lieberman, MD ****
667 Associate Chief Health Information Officer
668 Oregon Health Science University, Portland, OR

669 **Andrew Litt, MD**
670 Chief Medical Officer
671 Dell Healthcare and Life Sciences, Park City, UT

672 **Erik Pupo, CPHIMS**
673 Senior Manager, Deloitte Consulting LLP
674 Health Sciences and Government, Alexandria, VA

675 **Christopher Queram, MA**
676 President and CEO
677 Wisconsin Collaborative for Healthcare Quality, Madison, WI

678 **Carol Raphael, MPA, M.Ed.**
679 Advanced Leadership Fellow at Harvard; Former President and Chief Executive Officer
680 Visiting Nurse Service of New York (VNSNY), New York, NY

681 **Deborah A. Reid, JD, MHA**
682 Senior Attorney
683 National Health Law Program, Washington, DC

684 **Joyce Sensmeier, MS, RN-BC, CPHIMS, FHIMSS, FAAN**
685 Vice President, Informatics
686 Healthcare Information and Management Systems Society (HIMSS), San Diego, CA

687 **Shannon Sims, MD, PhD**
688 Director of Clinical Informatics
689 Rush University Medical Center, Chicago, IL

690 **Christopher Snyder, DO**
691 Chief Medical Information Officer/Chief Quality Officer
692 Peninsula Regional Medical Center, Salisbury, MD

693 **Christopher Tonozzi, MD ***
694 Chief Medical Information Officer
695 Colorado Associated Community Health Information Enterprise, Denver, CO

696 **Madhavi Vemireddy, MD**
697 Chief Medical Office and Head of Product Management
698 ActiveHealth Management, New York, NY

699 **Judith Warren, PhD, RN, BC, FAAN, FACMI ***
700 Retired, Professor
701 University of Kansas School of Nursing, Kansas City, KS

702 **Federal Liaisons**

703 **Joseph Francis, MD, MPH**
704 Director of Health Performance Measurement, Office of Informatics and Analytics
705 Veterans Health Administration, Washington, DC

706 **Erin Grace, MHA**
707 Senior Manager, Health IT
708 Agency for Healthcare Research and Quality, Rockville, MD

709 **Christopher Lamer, PharmD**
710 Medical Informaticist, Office of Information Technology
711 Indian Health Service, Rockville, MD

712 **Kevin Larsen, MD**
713 Medical Director Meaningful Use
714 Office of the National Coordinator for Health IT, Washington, DC

715 **Martin Rice, MS, RN-BC**
716 Deputy Director, Office of HIT and Quality
717 Health Resources and Services Administration, Rockville, MD

718 * Participated on subcommittee

719 ** Also participated on CSAC subcommittee

720 **NQF Staff**

721 **Helen Burstin, MD, MPH**

722 Senior Vice President

723 Performance Measurement

724 **Karen Beckman Pace, PhD, RN**

725 Senior Director

726 Performance Measurement

727 **Christopher Millet, MS**

728 Senior Director

729 Performance Measurement

730 **Karen Johnson, MS**

731 Senior Director

732 Performance Measurement

733 **Reva Winkler, MD, MPH**

734 Senior Director

735 Performance Measurement

736 **Taroon Amin, MA, MPH**

737 Senior Director

738 Performance Measurement

739 **Evan Williamson, MPH, MS**

740 Project Manager

741 Performance Measurement

742 **Jessica Weber, MPH**

743 Project Manager

744 Performance Measurement