

NATIONAL QUALITY FORUM

TO: NQF Board of Directors
FR: CSAC
SU: Overview of Evidence and Measure Testing Task Force Guidance Reports
DA: September 13, 2010

BOARD ACTION

The CSAC approved the following guidance documents and they are now presented to the Board for final approval.

- Guidance for Evaluating the Evidence Related to the Focus of Quality Measurement and Importance to Measure and Report
- Guidance for Measure Testing and Scientific Acceptability of Measure Properties

Once approved by the NQF Board, CSAC will work with staff to implement the new Reports' recommendations effective January 2011.

BACKGROUND

Last October the Board directed NQF to strengthen guidance to consistently apply the measure evaluation criteria. To that end, NQF convened two task forces to review the criteria and develop guidance to clarify and apply the measure evaluation criteria. One task force, chaired by Dr. David Shahian, focused on the evidence supporting the measure focus, as well as the criterion of Importance to Measure and Report. The other task force, chaired by Dr. Timothy Ferris, focused on measure testing for reliability and validity, as well as the criterion of Scientific Acceptability of Measure Properties.

Process

The task forces met in-person once, which was followed by several conference calls and email discussions to develop the draft recommendations. The draft recommendations were shared with the CSAC for comment prior to posting for public comment, as well as after the comment period. The task forces reviewed and responded to the comments received resulting in some clarifications and modifications to the guidance reports. Additional clarifications were made as a result of the CSAC final review.

Overview

The purpose of these reports is to provide guidance to NQF Steering Committees and others evaluating measures for potential NQF endorsement, as well as measure developers who submit measures to NQF. The recommendations provide greater clarity on how to apply the criteria to strengthen the measure evaluation process and resulted in only modest changes to the evaluation criteria. Although the recommendations provide more explicit guidance on how to evaluate measures, they do not (and were not intended to) create an automatic scoring and decision about recommending measures for endorsement. They do not supplant the need for expert judgment and multi-stakeholder involvement. Neither can they substitute for the expertise needed for measure development.

Implementation of these recommendations should be monitored to assess if they result in the intended effect and do not adversely affect submission of measures to NQF.

GUIDANCE FOR EVALUATING THE EVIDENCE RELATED TO THE FOCUS OF QUALITY MEASUREMENT AND IMPORTANCE TO MEASURE AND REPORT

Following are the key features of the guidance.

- The guidance document identifies the type of evidence that is needed for various types of measures – primarily the quantity, quality, and consistency of a body of evidence related to the relevant structure-process-outcome linkages (see Table 3).
- Ratings for evaluating the quantity, quality, and consistency of the body of evidence on a scale of high, moderate, and low were developed (Table 4), as well as how to use those ratings to determine if a measure has met the evidence criterion (see Table 5).
- Two potential exceptions to the requirement for empirical evidence are addressed: 1) when expert opinion might be used, and 2) for outcome measures (see Table 5).
- The preferred evidence grading systems were identified (USPSTF and GRADE); however, evidence graded using other systems may be submitted in support of a measure. Regardless of the evidence grading system, the goal is transparency so that a summary of the quantity, quality, and consistency of the body of evidence needs to be submitted for review.
- The guidance does not direct that measure developers conduct primary reviews and grade the evidence; rather, they should utilize existing evidence reviews to the extent possible, such as those in guidelines or other systematic reviews and summarize the body of evidence and conclusions about the strength of the evidence when submitting a measure.
- The recommendations also indicate that all three subcriteria under *Importance to Measure and Report* (high impact, opportunity for improvement, and evidence) must be met to pass this threshold criterion (see Table 5).
- At the time of review for endorsement maintenance, overall high performance with little variation should result in removal of endorsement unless there is a strong justification to continue endorsement.
- The evidence required for NQF-endorsed practices should parallel what is required for a process measure.

Comments Received

The key issues raised in the comments included the following.

- Burden for measure developers to conduct primary evidence reviews
- Expert opinion should be distinguished from evidence
- Concern about the identification of preferred evidence grading systems

- Requirement for evidence related to outcome measures may stifle submissions
These issues were discussed and resulted in clarifications in the final report.

GUIDANCE FOR MEASURE TESTING AND SCIENTIFIC ACCEPTABILITY OF MEASURE PROPERTIES

Following are the key features of the guidance.

- Reliability and validity need to be demonstrated through empirical evidence for all types of measures and data types.
- Ratings for reliability and validity on a scale of high, moderate, and low (Table 2) were developed, as well as how to use those ratings to determine if a measure meets the criterion for *Scientific Acceptability of Measure Properties* (Table 3). Failure to pass the criterion of *Scientific Acceptability of Measure Properties* should result in no recommendation for endorsement.
- The recommendations allow flexibility and ways to mitigate some of the burden of testing to achieve a moderate rating, which is necessary to pass the criterion.
- The same criteria and guidance is applicable to measures specified for EHRs, however, that was detailed in a separate table (Table 4).
- Examples of types of testing are provided in the Appendix.
- Untested measures that meet the conditions to be considered for endorsement in an NQF project must also meet requirements for specifications to be ready for testing (Table 5).
- Reliability and validity testing requirements for endorsement maintenance are indicated (Table 6).

Comments Received

The key issues raised in the comments included the following.

- Burden of testing
- Question of applicability to all measures/data types (e.g., claims, EHR)
- Scope of testing (sample size)
- Ratings should incorporate scope and appropriateness
- Disagreement with requirement for QDS specifications for EHR measures
- Questions regarding the requirements at the time of review for endorsement maintenance
- Provide Examples, references

These issues were discussed and resulted in either clarifications or explanations in the final report.

NATIONAL QUALITY FORUM

Guidance for Measure Testing and Evaluating Scientific Acceptability of Measure Properties

September 13, 2010

NATIONAL QUALITY FORUM

Guidance for Measure Testing and Evaluating Scientific Acceptability of Measure Properties

TABLE OF CONTENTS

INTRODUCTION AND CHARGE	2
Task Force Charge.....	2
BACKGROUND.....	3
Table 1. Measure Evaluation Criteria and Type of Evidence	4
Reliability and Validity	4
Reporting of Measure Scores and Scientific Acceptability.....	7
Measure Testing Issues Identified with Measures Submitted to NQF.....	7
Electronic Health Records and Electronic Measures	9
Summary of Background	9
RECOMMENDATIONS.....	10
I. Recommendations for Empirical Evidence of Reliability and Validity	10
II. Recommendations for the Type of Testing and Results Needed to Demonstrate Scientific Acceptability of Measure Properties	13
Table 2. Evaluation Ratings for Reliability and Validity	16
Table 3. Evaluation of Scientific Acceptability of Measure Properties Based on Reliability and Validity Ratings.....	17
III. Recommendations for Measures Specified for EHRs	18
Table 4. Evaluation of Reliability and Validity of Measures Specified for EHRs	21
IV. Recommendations Related to Untested Measures.....	22
Table 5. Minimum Requirements for Untested Measures under Scientific Acceptability of Measure Properties	23
V. Recommendations for Testing Required for Maintenance of Endorsement	23
Table 6. Scope of Testing Required at the Time of Review for Endorsement Maintenance.....	24
VI. Recommendations for Modifications to the NQF Evaluation Criteria.....	24
VII. Recommendations for the Measure Submission.....	28
REFERENCES.....	31
APPENDIX A – COMMON APPROACHES TO TESTING.....	33
Table A-1 Reliability Testing at the Level of the Computed Performance Measure Score.....	33
Table A-2. Reliability Testing at the Level of the Data Elements	34
Table A-3. Validity Testing at the Level of the Computed Performance Measure Score	35
Table A-4. Validity Testing at the Level of Data Elements	36
Table A-5 Testing Related to Threats to Validity	37
Table A-6. Interpretation of Statistical Results.....	38
APPENDIX B – TASK FORCE MEMBERS.....	39
APPENDIX C – GLOSSARY.....	40
APPENDIX D – MEASURE EVALUATION CRITERIA	43

1 INTRODUCTION AND CHARGE

2 The National Quality Forum (NQF) relies on [four criteria](#) for evaluating the suitability of quality
3 measures for endorsement as voluntary consensus standards: Importance to Measure and
4 Report, Scientific Acceptability of Measure Properties, Usability, and Feasibility. The second
5 criterion, *Scientific Acceptability of Measure Properties*, is an important aspect of the successful use
6 of publicly reported measures to improve performance. Scientific acceptability of measure
7 properties refers to the reliability and validity of measures. The use of measures that are
8 unreliable or invalid undermines confidence in measures among both the providers of
9 healthcare and the consumers of the information. The goal of this document is to provide
10 recommendations on what constitutes scientific acceptability of measures to assist those
11 participants in the measure evaluation process, including steering committees and technical
12 advisory panel members, as well as measure developers. Guidance on scientific acceptability
13 will facilitate a shared understanding of this complex and highly specialized subject.

14
15 In evaluating a measure, both empirical evidence and expert judgment play a role. However,
16 judgment can best be applied when those evaluating a measure have a thorough understanding
17 of the evidence of scientific acceptability that does or does not exist. Evidence that a clearly
18 specified measure produces credible results on performance comes from the basic measurement
19 principles of reliability and validity. Although reliability and validity have always been
20 included in NQF evaluation criteria, the criteria have not included specific guidance on 1) the
21 scope of testing, 2) what tests of reliability and validity could be performed, and 3) how to
22 weigh the results of this testing.

23
24 Task Force Charge

25 The NQF Task Force on Measure Testing was asked to address the following tasks.
26 • Identify the type of testing for scientific acceptability that should be conducted for various
27 types of measures and data sources, and determine whether there are any acceptable
28 alternatives to formal testing.

- 29 • Identify the type of testing that should be required prior to endorsement of measures
30 specified for electronic health records (EHRs) – both measures originally developed using
31 other data sources besides the EHR and new measures developed specifically for EHRs.
- 32 • Develop guidance for measure stewards/ developers and NQF technical advisors and
33 steering committees on adequate measure testing, interpretation of results, and information
34 about testing that should be provided in the measure submission.
- 35 • Make recommendations for potential enhancements to the evaluation criteria.

36
37

38 BACKGROUND

39 NQF endorses quality measures intended for quality improvement as well as public reporting.
40 Measure scores are used to make decisions about selecting and rewarding healthcare providers
41 (e.g., by consumers and purchasers) and to identify opportunities for quality improvement (e.g.,
42 by providers). The level of confidence one can have in conclusions about quality based on the
43 measure scores is a function of the reliability and validity of measurement.

44

45 The NQF measure evaluation criteria can be viewed as a hierarchy that guides the sequential
46 process for evaluating measures. As described in some of the foundational work for NQF
47 processes:

48 “If a measure is not important, its other characteristics are less meaningful. If a
49 measure is not scientifically acceptable, its results may be at risk for improper
50 interpretation. If a measure is not interpretable [usable] we probably do not care if it is
51 feasible. If a measure is not feasible, alternative approaches to acquiring important
52 information should be considered (p. I-40).”¹

53 Once a measure has been determined to meet the criterion of *Importance to Measure and Report*, it
54 is evaluated on the criterion, *Scientific Acceptability of Measure Properties*. This criterion addresses
55 the basic measurement principles of reliability and validity. The NQF evaluation criteria
56 parallel best practices for measure development, which include testing reliability and validity. ²

57 ³

58

59 NQF’s measure [evaluation criteria](#) include a variety of types of evidence as indicated in Table 1.
 60 The criterion, *Scientific Acceptability of Measure Properties*, addresses *how* the healthcare quality
 61 concept is measured. This criterion includes reliability (2b) and validity (2c), as well as precision
 62 of specifications (2a) and potential threats to valid conclusions about quality related to
 63 exclusions (2d), risk adjustment for outcome and resource use measures (2e), and comparability
 64 of results from different data sources (2g). The other subcriteria include identification of
 65 differences in performance (2f) and specifications to detect disparities (2h).

67 Table 1. Measure Evaluation Criteria and Type of Evidence

Evaluation Criteria	Type of Evidence
1. Importance to measure and report 1a. High impact 1b. Opportunity for improvement 1c. Evidence that supports the focus of measurement	Epidemiologic data Resource use data Health services research Clinical research
2. Scientific acceptability of measure properties (reliability, validity, etc.)	Psychometric testing - reliability and validity, adequacy of risk adjustment, etc.
3. Usability 3a. Demonstration of understanding and usefulness for public reporting and quality improvement	Data and/or qualitative information demonstrating usefulness for public reporting and quality improvement
4. Feasibility 4e. Demonstration the measure can be implemented	Data and/or qualitative information demonstrating the measure can be implemented

68

69 **Reliability and Validity**

70 A quality measure is a numeric quantification of the relatively abstract construct of quality of
 71 healthcare, which is measured imperfectly. Reliability refers to the *repeatability or precision of*
 72 *measurement*. Validity refers to the *correctness of measurement*. The concepts of reliability and
 73 validity can be applied to the individual data elements used in a measure (e.g., diagnosis,
 74 medication, admission date, birth date), as well as the computed performance measure score
 75 (e.g., rate, proportion, average).

76

77 Reliability of data elements refers to repeatability and reproducibility of the data elements for
 78 the same population in the same time period. Validity of data elements *refers* to the correctness
 79 of the data elements as compared to an authoritative source.

80

81 Reliability of the measure score refers to the proportion of variation in the performance scores
82 due to systematic differences across the measured entities in relation to random error or noise.

83 Validity of the measure score refers to the correctness of conclusions about the quality of
84 measured entities that can be made based on the measure scores (i.e., a higher score on a quality
85 measure reflects higher quality)

86

87 Over the past four to five decades numerous methods have been devised to test measures and
88 thus address the measure properties inherent to all measurement. These approaches provide
89 empirical evidence of the properties of reliability and validity. Examples of approaches to
90 reliability and validity testing can be found in Tables A-1 through A-5 ([Appendix A](#)).

91

92 A measure score is an approximation of a theoretical “true” score plus error: The more error, the
93 less reliable and valid is the measurement. Random or chance errors affect the reliability or
94 repeatability of measurement and systematic errors affect the validity or correctness of the
95 conclusions one can make based on the measure score. Threats to reliability include ambiguous
96 measure specifications (including definitions, codes, data collection, and scoring) and small case
97 volume or sample size. Threats to validity include other aspects of the measure specifications
98 such as inappropriate exclusions, lack of appropriate risk adjustment or risk factors for
99 outcomes and resource use, specifications for multiple data sources or methods that result in
100 different scores and conclusions about quality, and systematic missing or “incorrect” data. Most
101 importantly, a measure may be invalid because the measurement has not correctly captured the
102 concept of quality it was intended to measure.

103

104 Reliability and validity are not all-or-none properties; rather, measures of reliability and validity
105 produce graduated results. Therefore, results of measure testing always require interpretation.
106 Reliability and validity are not static; they are influenced by the conditions under which the
107 measures are implemented (e.g., local documentation and coding practices, structures of
108 records, etc.). Evidence of validity, in particular, is accumulated over time. A discussion of
109 measurement concepts can be accessed in an online [research methods knowledge base](#).⁴ [Rubin](#)

110 [et al.](#)³ and others⁵ describe reliability and validity testing in quality measure development.
111 Examples of validity testing of healthcare quality measures also are reported in the literature.^{6,7}

112
113 Reliability is often considered to be necessary, but not sufficient, for achieving validity. That is,
114 if a measure is not reliable, a valid conclusion about quality would not be possible; and a
115 measure could be reliable, but wrong leading to incorrect (invalid) conclusions. However, this
116 relationship between reliability and validity is not universally held^{8,9} and may depend on how
117 a measure is defined. For example, if a measure is mean systolic blood pressure (BP), the mean
118 could be accurate even if the individual BP readings are unreliable (i.e., with substantial
119 random error). On the other hand, if the numerator of a measure is defined as systolic BP over
120 140, then unreliability of the measure can lead to assigning to the wrong category and hence
121 loss of validity.

122
123 Evaluation of the scientific acceptability of a measure does not occur in a vacuum. The Task
124 Force was aware of factors within the current environment affecting their deliberations. The
125 recommendations of the Task Force would have implications for both measure developers and
126 healthcare providers. For example, some observers have suggested that existing measure
127 evaluation criteria are too stringent (allowing “the perfect to be the enemy of the good”) while
128 others have suggested that the criteria are not rigorous enough. Some contend that providers
129 use adherence to the measure evaluation criteria as a barrier to making performance
130 information available; others maintain that unless a measure has adequate measure properties it
131 cannot provide useful information. Nonetheless, the consequences of using unreliable or
132 invalid measures can at times be significant for those being measured as well as those who use
133 the information to select a healthcare provider. Resources may be wasted or misdirected; and
134 there is potential for invalid measures to result in misinformation and misdirection of patients
135 or potential unintended harmful consequences. As the stakes around quality measurement are
136 raised, the potential for conflicts among these perspectives increases. The Task Force therefore
137 made a deliberate attempt to make recommendations that balanced the requirement for
138 insuring that NQF endorsed measures would be both sufficiently reliable and valid to make
139 them meaningful and minimize unintended consequences, with requirements for testing that
140 were not so high as to stifle measure development and innovation.

141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171

Reporting of Measure Scores and Scientific Acceptability

NQF does not determine the specific use or reporting formats of the measures it endorses. Nonetheless, the confidence in a measure can be related to the context in which the measure is used and the choices made in reporting performance measure scores. For example, Kaplan and colleagues¹⁰ demonstrated that the number of categories chosen for performance reporting (e.g., high/medium/low) influences the likelihood of misclassification. Misclassification is, by definition, an invalid reporting of performance. Reporting performance from highest to lowest, without information on margin of error and meaningful differences, limits and may misrepresent the knowledge to be gained from measures. Further, those choosing to report measures may decide to combine the measures into a composite in order to simplify reporting, making the metrics more usable for consumers and providing another way for providers to view performance. These combined composite measures also have potential to be misleading.¹¹ On the other hand, confidence intervals or other technical explanations could render the information incomprehensible to some audiences. Finding the right balance is important. Because NQF endorsement does not dictate how the measures are used, the Task Force was not asked to make recommendations on reporting but these issues are highlighted for further discussion and assessment.

Measure Testing Issues Identified with Measures Submitted to NQF

The Task Force understood their charge as emerging from several years of NQF experience with measure evaluation. This experience, enumerated below in six points, informed the Task Force’s recommendations. First, the NQF portfolio of endorsed measures shows considerable variation in the level of rigor used in measure testing. Measure developers are currently expected to address these requirements in a way that is most appropriate and feasible for the measure and data source involved. Nonetheless, some developers submit limited information on reliability or validity testing perhaps due to a lack of expertise or resources. On the other hand, other measure developers have conducted formal reliability and validity testing and have demonstrated that a proposed measure generates reproducible results and credible conclusions about quality.

172 Second, when reliability and validity testing results have been submitted, there has been
173 variability in the scope of testing and the rigor of methods and statistical analysis. For example,
174 reliability of categorical data elements may be assessed only as the percentage of agreement
175 between raters versus using the kappa statistic, which adjusts for chance agreement. In some
176 cases, the testing was conducted with a particular data source, such as the paper medical
177 record, while the measure was specified using a different data source, such as electronic health
178 record.

179
180 Third, there also has been some confusion regarding what is considered testing of scientific
181 acceptability. Terms such as “measure testing,” “pilot testing,” and “field testing” are
182 commonly used in the discipline of measure development and include reliability and validity
183 testing, as well as other aspects of measure development. For example, measure submissions
184 may include descriptive statistics that demonstrate the data are available and can be analyzed to
185 produce scores, but do not specifically address reliability or validity.

186
187 Fourth, some submissions rely on an assumption of reliability and validity. This assumption
188 may be based on prior use of the measure or some aspects of the measure specifications (e.g.,
189 diagnosis codes are relatively well defined and used in accordance with coding rules). In some
190 cases an argument is made that a data source would become more reliable and valid if a quality
191 measure was implemented and publicly reported.

192
193 Fifth, measure developers rarely submit analyses justifying exclusions or demonstrating
194 comparability of different methods of data collection.

195
196 Sixth, steering committees may variably weigh the strengths and weaknesses of the evidence for
197 reliability and validity in their recommendation for endorsement. In summary, while NQF has
198 been raising the bar of expectations and introducing greater rigor and standardization to the
199 evaluation process, the NQF portfolio of endorsed measures still includes varying levels of
200 methodological rigor.

201

202

203 **Electronic Health Records and Electronic Measures**

204 Development and implementation of electronic health record (EHR) systems hold great promise
205 for the efficient collection of clinical data that can be used for quality measurement. National
206 initiatives call for the adoption of electronic health records that include the capability for quality
207 measurement and NQF has made endorsing quality measures specified for EHRs an important
208 goal. Data stored in EHRs facilitate reporting of quality measures because EHR data 1) are
209 clinically specific, 2) include a large variety of data types including physiologic data such as
210 laboratory values, and 3) decrease the burden of the data collection through automated
211 collection and aggregation.

212

213 While the concepts of reliability and validity apply equally to measures derived from EHRs, the
214 electronic health record also presents additional issues related to measure testing. Widespread
215 EHR data are not yet available for measure development and testing. In addition, the numerous
216 vendors and home grown EHR systems present the additional challenge of insuring that the
217 selected data fields of interest for any particular measure are comparable among different
218 EHRs. Recommendations regarding testing and evaluation of EHR measures are addressed in
219 Section III.

220

221 **Summary of Background**

- 222 • There are no perfect quality performance measures and there will be some error in all
223 measurement. Performance measurement science is an imperfect science.
- 224 • Measurement principles of reliability and validity apply to quality performance measures
225 regardless of data source.
- 226 • Reliability and validity are not all-or-none properties and involve a matter of degree.
- 227 • Reliability and validity are not static properties and can vary under the conditions of
228 implementation.
- 229 • Reliability and validity can apply to individual data elements used in a measure, as well as
230 the computed measure score.
- 231 • Reliability does not guarantee validity.

- 232 • Variability in measure scores that is attributable to either random error (noise) or systematic
233 error (biased measurement) is misleading and leads to unwarranted conclusions about
234 quality.
- 235 • NQF is ultimately concerned with endorsing measures that produce scores from which
236 valid (i.e., correct) conclusions about the quality of care can be made.
- 237 • A measure that is not a valid indicator of quality is not useful for making decisions about
238 selecting healthcare providers based on quality or investing time and resources into
239 improvement.

240

241

242 RECOMMENDATIONS

243 The recommendations in this report are intended to provide additional guidance and
244 clarification regarding the NQF criteria related to measure testing and scientific acceptability.
245 However, the guidance does not address the unique aspects of testing for composite measures
246 as indicated in the composite measure evaluation criteria. The guidance is not intended to
247 provide a detailed primer on methods for measure testing. The recommendations also are not
248 intended as a definitive scoring system for measure evaluation; evaluation still requires
249 judgment regarding the adequacy of the empirical testing evidence. The recommendations
250 should promote greater consistency in applying the NQF criteria, while maintaining
251 consideration of multi-stakeholder perspectives during the evaluation. This guidance then
252 replaces any previous guidance on measure testing (e.g., field testing requirements in time-
253 limited endorsement policy).

254

255 I. Recommendations for Empirical Evidence of Reliability and Validity

256 Before developing guidance on the specific testing criteria, the Task Force was asked to consider
257 a fundamental question of whether reliability and validity need to be demonstrated empirically
258 or could be assumed or agreed upon through various review or consensus processes. The Task
259 Force recommended that *empirical evidence of reliability and validity should be expected for*
260 *all measures endorsed by NQF.*

261

262

263 **Rationale for Empirical Evidence**

264 Although reliability and validity are not static properties and can vary under different
265 conditions of implementation (e.g., local documentation and coding practices, structures of
266 paper or electronic records, etc.), the purpose of reliability and validity testing for consideration
267 of NQF endorsement is to demonstrate that a measure could be reliable and valid when
268 implemented as specified.

269
270 Although precise specifications provide a foundation for consistent implementation and thus
271 increase the likelihood of reliability, reliability cannot be assumed. Although evidence for the
272 measure focus (NQF criterion [1c](#)) provides a foundation for the validity (NQF criterion [2c](#)) of
273 the measure as an indicator of quality, the way a measure is specified can affect the validity of
274 the conclusions about quality.

275
276 Implementation and reporting of measures is expected to lead to improvements in
277 documentation, data coding, and data capture and thus reliability and validity. This assumption
278 of improved reliability and validity over time applies to all measures regardless of data type;
279 however, it does not negate the need for empirically demonstrating reliability and validity
280 when a measure is being considered for endorsement.

281
282 Recommendations for measures specified for EHRs are addressed in a separate section (Section
283 III) because they are newer and there are several differences from other data types. For example,
284 the clinician is often the source of data in EHRs and the data are intended for use in care
285 management. However, these distinctions are not absolute and the same requirement for
286 demonstrating scientific acceptability applies equally to EHR measures as to measures based on
287 other data types. Administrative claims data and EHR data may be viewed as complementary
288 sources of information, each with their own strengths and limitations.

289
290 **Strategies to Mitigate the Burden of Testing**

291 Although the Task Force was clear about the recommendation for empirical evidence of
292 reliability and validity, it also recognized the practical implications of this assertion for measure

293 developers. The Task Force therefore, further recommended some strategies that could
294 minimize the burden of testing as follows.

- 295 • Evidence for reliability and validity may be accumulated over time and evaluators
296 should remain flexible with regard to the extent of testing evidence submitted. The
297 scope of testing may be on a relatively small scale for initial endorsement, followed by
298 further analyses to support continued endorsement at the time of review for
299 maintenance of endorsement.
- 300 • Reliability and validity testing may be conducted on a sample of the measured entities.
301 The analytic unit of the particular measure (e.g., physician, hospital, home health
302 agency) determines the sampling strategy for scientific acceptability testing.
 - 303 ○ The sample should represent the variety of entities whose performance will be
304 measured. The Task Force recognized that the samples used for reliability and
305 validity testing often have limited generalizability because measured entities
306 volunteer to participate. Ideally, however, all types of entities whose
307 performance will be measured should be included in reliability and validity
308 testing.
 - 309 ○ The sample should include adequate numbers of units of measurement and
310 adequate numbers of patients to answer the specific reliability or validity
311 question with the chosen statistical method.
 - 312 ○ When possible, units of measurement and patients within units should be
313 randomly selected.
- 314 • Reliability and validity testing may be conducted for either the data elements used to
315 calculate the measure score or the computed measure score, to achieve an acceptable
316 rating for endorsement. Although ideally testing is conducted for both the critical data
317 elements and the computed measure score, only one level of testing would be required
318 for endorsement. See Tables A-1 to A-5 in [Appendix A](#) for examples of reliability and
319 validity testing of data elements and measure scores.
- 320 • Separate reliability testing of the data elements is not required if empirical validity
321 testing of the data elements (see [Table A-4](#)) is conducted (e.g., if the validity of ICD-9
322 codes in administrative claims data as compared to clinical diagnoses in the medical

323 record is demonstrated, then inter-coder or inter-abstractor reliability would not be
324 required).

- 325 • Prior evidence of reliability or validity of data elements (see Tables [A-2](#) and [A-4](#) in
326 [Appendix A](#)) for the data type specified in the measure (e.g., hospital claims) can be
327 used as evidence for those data elements. Prior evidence could include published or
328 unpublished testing that:
 - 329 ○ included the same data elements; and
 - 330 ○ used the same data type (e.g., claims, chart abstraction, etc.) ; and
 - 331 ○ was conducted on a sample as described above (i.e., representative, adequate
332 numbers, and randomly selected, if possible).
- 333 • Because validity testing of measure scores can be quite burdensome, a formal and
334 systematic testing of [face validity](#) as described in Table A-3 could be acceptable for a
335 moderate rating of measure score validity.^{12, 13} Key components include systematic and
336 transparent process, the inclusion of identified experts, and explicitly addressing
337 whether performance scores resulting from the measure as specified can be used to
338 distinguish good from poor quality.

339
340 The Task Force further acknowledged that there are degrees of reliability and validity and the
341 following guidance distinguishes ideal testing and evidence from what is acceptable for
342 endorsement by NQF. Measures without empirical testing of reliability and validity should be
343 considered untested measures and subject to NQF's [conditions for considering untested](#)
344 [measures](#) for endorsement. Untested measures are addressed in Section IV.

345 346 **II. Recommendations for the Type of Testing and Results Needed to Demonstrate Scientific** 347 **Acceptability of Measure Properties**

348 How should participants in the evaluation process assess the evidence provided when
349 measures are submitted? The Task Force chose to provide guidance on measure testing
350 through the development of rating categories for the reliability and validity of measures being
351 considered for endorsement. This approach requires well-defined descriptions of the rating
352 scheme to reduce ambiguity and miscommunication. While the Task Force has tried to achieve
353 this precision, it recognizes that there will inevitably be some ambiguity and room for

354 interpretation. In addition, the rating descriptions provided in this report may require further
355 clarification and/or revision. Finally, the Task Force was not able to fully assess the impact of
356 the proposed rating system on the measure endorsement process. So, this proposed approach to
357 evaluating scientific acceptability of measure properties should be monitored to ensure it
358 achieves the intent of endorsing reliable and valid measures and does not unduly impede
359 endorsement of measures.

360
361 The Task Force chose to provide guidance on evaluating *Scientific Acceptability of Measure*
362 *Properties* using a two-step process. First, guidance is provided on how to rate the evidence for
363 reliability and validity. Second, guidance is provided on how to use the ratings to determine if
364 the criterion of *Scientific Acceptability of Measure Properties* is met.

365
366 Table 2 provides the guidance for rating the level of evidence for reliability and validity, which
367 is classified as high, moderate, or low. The ratings depend on the level of testing conducted,
368 appropriateness of the selected method, scope of testing, and the results of testing meeting
369 acceptable norms. This table applies to all types of measures and data types; however, in Table
370 4, the rating scale is applied specifically to EHR measures.

371
372 The rating scheme is structured around a distinction between testing the data elements used to
373 calculate a measure (e.g., diagnosis, procedure, age) and the computed measure scores (e.g.,
374 rate, proportion, average). Some measures rely on many data elements. Testing at the data
375 element level does not necessarily need to be conducted for every single data element, but
376 should include those elements that are most critical to the computed score. The [critical data](#)
377 [elements](#) are those that contribute most to the computed measure score.

378
379 Testing at either the level of data elements or the computed measure score with appropriate
380 methods and scope and acceptable results is rated moderate and would be acceptable for
381 endorsement. Testing at both levels of data elements and computed measure score with
382 appropriate methods and acceptable results is rated high. The low rating represents evidence
383 that a measure has low reliability or validity. If the testing was conducted with an inappropriate
384 method or inadequate scope (i.e., representativeness, sample size), there would be inadequate

385 evidence to evaluate reliability and/or validity and the measure would be considered untested.
386 As noted previously, untested measures would not be rated on reliability and validity and
387 special considerations for untested measures are addressed in a separate section (see Section
388 IV).

389
390 The rating scale presented in Table 2 is not intended to provide a definitive scoring system. The
391 determination of adequate testing and results still requires judgment that incorporates a variety
392 of considerations including:

- 393 • whether the test was appropriate for the specified measure;
- 394 • whether the scope of testing (i.e., representativeness, sample size) was adequate ; and
- 395 • whether the results indicate acceptable level of reliability or validity.

396
397

398 Table 2. Evaluation Ratings for Reliability and Validity

Rating	Reliability	Validity
High	<p>All measure specifications (e.g., numerator, denominator, exclusions, risk factors, scoring) are unambiguous and likely to consistently identify who is included and excluded from the target population and the event, condition, or outcome being measured; how to compute the score, etc.;</p> <p>AND</p> <p>Empirical evidence of reliability of <u>both data elements</u> (Table A-2) <u>and measure score</u> (Table A-1):</p> <ul style="list-style-type: none"> • <u>Data element</u>: 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); <i>OR may forego data element reliability testing if data element validity (Table A-4) was demonstrated;</i> <p>AND</p> <ul style="list-style-type: none"> • <u>Measure score</u>: appropriate method, scope, and reliability statistic within acceptable norms 	<p>The measure specifications (numerator, denominator, exclusions, risk factors) reflect the quality of care problem (1a,1b) and evidence cited in support of the measure focus (1c) under <i>Importance to Measure and Report</i>;</p> <p>AND</p> <p>Empirical evidence of validity of <u>both data elements</u> (Table A-4) <u>and measure score</u> (Table A-3):</p> <ul style="list-style-type: none"> • <u>Data element</u>: appropriate method, scope, and statistical results within acceptable norms (new testing, or prior evidence for the same data type) for critical data elements; AND • <u>Measure score</u>: appropriate method, scope, and validity testing result within acceptable norms ; AND <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>All measure specifications are unambiguous as noted above</p> <p>AND</p> <p>Empirical evidence of acceptable reliability for <u>either critical data elements OR 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 acceptable validity for <u>either critical data elements OR measure score</u> as noted above; OR</p> <p><u>Systematic assessment of face validity</u> of <u>measure score</u> as a quality indicator (as described in Table A-3) 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>
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>low reliability</u> for <u>either data elements OR 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>low validity</u> for <u>either data elements OR 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>
Inadequate 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>

400 Table 3 presents the Task Force’s recommendation on how the ratings for reliability and validity
 401 are used to determine whether a measure adequately meets the criterion of *Scientific*
 402 *Acceptability of Measure Properties*. Moderate ratings for both validity and reliability as described
 403 in Table 2 (and Table 4) would be required to pass this criterion and be acceptable for
 404 endorsement. A high rating is not required for endorsement, but represents current thinking
 405 about best practices in measure development. **A measure that does not pass the criterion of**
 406 ***Scientific Acceptability of Measure Properties* would not be recommended for endorsement.**

408 Table 3. Evaluation of Scientific Acceptability of Measure Properties Based on Reliability and Validity
 409 Ratings

Validity Rating	Reliability Rating	Pass <i>Scientific Acceptability of Measure Properties</i> for initial endorsement *	
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 low, the reliability rating will usually also be low. If validity is low and reliability is moderate-high, it represents inconsistent evidence.

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

412
 413 Some common approaches to testing reliability and validity for the data elements as well as the
 414 computed measure score that can be applied to quality performance measures are listed in
 415 [Appendix A](#) (Tables A-1 through A-5). Measure developers should select the testing that is
 416 appropriate and feasible for the measure under consideration and that will at least meet the
 417 moderate rating as described in Table 2. Table [A-5](#) also addresses potential testing and analysis
 418 related to the threats to validity represented by other subcriteria under *Scientific Acceptability of*
 419 *Measure Properties*. Measure developers should identify the potential threats to validity for the
 420 specific measure and conduct analyses to demonstrate that the results are not biased.
 421 Information on interpretation of the common statistical tests used to demonstrate reliability and
 422 validity also are provided in [Table A-6](#); however, those norms provide only general guidelines
 423 and testing results must be interpreted within the unique context of the specific measure.

424

425 The information on approaches to testing is not meant to provide an exhaustive list of methods.
426 Other approaches to testing may be appropriate and could be used if the method and rationale
427 are explained and judged to be appropriate. For example, if agreement on data elements
428 between two time periods is proposed as a test of reliability (test/re-test), the rationale for
429 expecting stability (rather than change) over the time period is important to discuss. Calculation
430 of measures scores and descriptive statistics, or the fact that a measure has been in use do not
431 constitute empirical evidence of reliability or validity. Such information may be relevant to the
432 criteria of opportunity for improvement (1b), identification of differences in performance (2f),
433 usability of the measure (3a), and feasibility of implementation (4e); but alone does not address
434 the reliability or validity of the measure.

435
436

437 III. Recommendations for Measures Specified for EHRs

438 The EHR holds significant promise for improving the measurement of healthcare quality. The
439 availability of a broad range of reliable and valid data elements for quality measurement
440 without the burden of data collection is widely anticipated. Because clinical data can be entered
441 directly into standardized computer readable fields, the EHR will be considered the
442 authoritative source of clinical information. Quality measures based on EHRs use clinical
443 information recorded by healthcare clinicians in discrete computer readable fields; therefore,
444 measurement errors due to manual abstraction, coding by persons other than the originator, or
445 transcription could be eliminated. Despite these potential advantages over current data
446 sources, several potential sources of error pose threats to the reliability and validity of data
447 elements and measure scores for EHR measures including: 1) incorrect measure specifications,
448 including code lists, logic, or computer readable programming language; 2) EHR system
449 structure or programming that does not comply with standards for data fields, coding, or
450 exporting data; 3) difference in use of data fields by different users or entry into the wrong EHR
451 field; 4) entry of incorrect information; and 5) incorrect parsing of data by natural language
452 processing software used to analyze information from text fields. All of these potential errors
453 are analogous to sources of error with measures based on other data sources.

454

455 Table 4 provides the guidance for rating the level of evidence for reliability and validity of EHR
456 measures and it is analogous to the ratings in Table 2. Just as for other measures, Table 3

457 indicates how the ratings are used to make a determination if the criterion, *Scientific Acceptability*
458 *of Measure Properties* has been met for EHR measures. Testing approaches for reliability and
459 validity of the EHR measure score are the same as for any measure as noted in Tables [A-1](#) and
460 [A-3](#).

461
462 There are two differences highlighted between Table 2 and Table 4. First, EHR measures must
463 be specified in accordance with the Quality Data Set (QDS).¹⁴ The reason for requiring
464 specifications using the QDS is twofold: 1) the QDS can be translated to machine readable
465 specifications that can be applied to EHRs; and 2) the structure of QDS will fulfill the criterion
466 for precise specifications. The QDS will be updated on a regular basis, so if a measure needs a
467 quality data element not currently available, there will be a process to consider additional
468 quality data elements so that the measure could achieve a moderate or high rating.

469
470 Second, data elements for quality measures, which are extracted from EHRs using computer
471 programming, are by virtue of automation repeatable (reliable); however, they could be wrong.
472 Because different uses of an EHR data field by a clinician or different data extraction protocols
473 in different EHRs can produce different performance scores, testing at the data element level
474 should focus on validity as discussed below. This approach is consistent with the rating system
475 presented in Table 2, that is, if empirical validity testing of the data elements is conducted,
476 separate reliability testing of the data elements is not required.

477
478 An approach to testing validity of data elements analyzes agreement between data elements
479 and scores obtained with data exported electronically using the EHR measure specifications to
480 those obtained by review and abstraction of the entire EHR, preferably using EHRs that comply
481 with standards. This approach has been reported in the literature¹⁵⁻¹⁷ and by HealthPartners in
482 a [Commonwealth report](#)¹⁸ on performance measures and EHRs. As with measures for other
483 data types, testing may be conducted on a [sample of the measured entities](#) (see Section I).

484
485 Because EHR databases may not be available for such testing, another approach is to apply the
486 EHR measure to a simulated data set that reflects standards for EHRs and includes sample
487 patient data with the data elements needed for the specified measure. Because the simulated

488 data set is constructed, the values for the data elements and scores are known. When the EHR
489 specifications are applied to the simulated data set, they should return the known values of the
490 data elements and scores.

491
492 With either approach, when the results obtained for the EHR measure do not match the known
493 values in the simulated data set or the abstracted data, an analysis is conducted to determine
494 the source of error. If the error is related to the measure specifications, including code lists,
495 logic, and computer readable programming language, they would be corrected before
496 submission for endorsement. If the source of error is due to clinical data entry practices and
497 EHR structures unique to specific organizations, the error would not be mitigated by changes to
498 the EHR measure specifications but it could indicate the need for further evaluation such as
499 feasibility and whether alternative data fields could be used.

500
501 The recommended approach for evaluating reliability and validity of data elements for EHR
502 measures takes into account the current environment in which standards for EHRs and EHR
503 measures are under development and widespread adoption is not yet reality. Therefore, testing
504 sites are limited and testing in a sample of EHR systems may not be representative of others.
505 However, this is no different than testing of data elements for measures based on other data
506 sources in a sample of the measured entities. As noted in the background, reliability and
507 validity are not static properties and no one test is definitive.

508
509 Measure testing requirements should not impede the adoption of EHRs and EHR measures, but
510 should be true to the principles of scientific acceptability. EHRs and EHR measures are new and
511 will most likely require some adjustment of local EHR structures and recording practices to
512 meet standards. Therefore, providers should be encouraged to conduct their own internal
513 reliability studies.

514
515 Previously endorsed measures specified for chart abstraction or administrative claims data may
516 be appropriate for specification for EHRs. Although these endorsed measures should have
517 already been tested for reliability and validity, the EHR measure specifications require some
518 assessment of similarity to the original specifications, which also is addressed in Table 4. In

519 some cases, the EHR specifications will represent a substantive change to the measure so that an
 520 assessment of reliability and validity of the EHR measure is needed.

521

522 Table 4. Evaluation of Reliability and Validity of Measures Specified for EHRs

Rating	New Measure Specified for EHR		Modifications for Endorsed Measures <u>Re-specified</u> for EHRs
	Reliability Description and Evidence	Validity Description and Evidence	
High	<p>All EHR measure specifications are unambiguous and include only data elements from the Quality Data Set (QDS) * including quality data elements, code lists, and measure logic; OR new elements are submitted for inclusion to the QDS;</p> <p>AND</p> <p>Empirical evidence of reliability of <u>both data element and measure score</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>Measure score</u>: appropriate method, scope, and reliability statistic within acceptable norms 	<p>The measure specifications (numerator, denominator, exclusions, risk factors) reflect the quality of care problem (1a,1b) and evidence cited in support of the measure focus (1c) under <i>Importance to Measure and Report</i>;</p> <p>AND</p> <p>Empirical evidence of validity of <u>both data elements and measure score</u>:</p> <ul style="list-style-type: none"> • <u>Data element</u>: validity demonstrated by analysis of agreement between data elements exported electronically and data elements 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 EHR measure 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>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>	<p>The EHR measure specifications use only data elements from the Quality Data Set (QDS) * and include quality data elements, code lists, and measure logic;</p> <p>AND</p> <p>Crosswalk of the EHR measure specifications (QDS quality data elements, code lists, and measure logic) to the endorsed measure specifications demonstrates that they represent the original measure, which was judged to be a valid indicator of quality;</p> <p>AND</p> <p>Analysis of comparability of scores produced by the retooled EHR measure specifications with scores produced by the original measure specifications demonstrated similarity within tolerable error limits</p>
Moderate	<p>All EHR measure specifications are unambiguous and include only data elements from the QDS; * OR new elements are submitted for inclusion to the QDS as noted above;</p> <p>AND</p> <p>Empirical evidence of reliability for <u>either data elements OR 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 for <u>either data elements OR measure score</u> as noted above; OR <u>Systematic assessment of face validity</u> of <u>measure score</u> as a quality indicator (as described in Table A-3) 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>	<p>The EHR measure specifications use only data elements from the QDS as noted above</p> <p>AND</p> <p>Crosswalk of the EHR measure specifications as noted above demonstrates that they represent the original measure</p> <p>AND</p> <p>For measures with time-limited status, testing of the original measure and evidence ratings of moderate for reliability and validity as described in Table 2.</p>
Low	<p>One or more EHR measure specifications are ambiguous or <u>do not</u> use data elements from the QDS * as noted above;</p>	<p>The EHR measure specifications do not 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</p>	<p>The EHR measure specifications do <u>not</u> use only data elements from the QDS;</p> <p>OR</p> <p>Crosswalk of the EHR measure</p>

Rating	New Measure Specified for EHR		Modifications for Endorsed Measures <u>Re-specified</u> for EHRs
	Reliability Description and Evidence	Validity Description and Evidence	
	OR Empirical evidence of <u>low reliability</u> for <u>either data elements</u> OR <u>measure score</u> – i.e., statistical results outside of acceptable norms	scope) of <u>low validity</u> for <u>either data elements</u> OR <u>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	specifications as noted above identifies that they do NOT represent the original measure OR For measures with time-limited status, empirical evidence of low reliability or validity for original time-limited measure
Inade-quate	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	Crosswalk of the EHR measure specifications as noted above was not completed OR For measures with time-limited status, inappropriate method or scope of reliability or validity testing for original time-limited measure

523 *QDS elements should be used when available. When needed quality data elements are not yet available in the QDS,
524 they will be considered for addition to the QDS.

525

526 IV. Recommendations Related to Untested Measures

527 Measures without empirical evidence of reliability and validity are considered untested.

528 Untested measures are only eligible for time-limited endorsement if the conditions for
529 considering time-limited endorsement are met.

- 530
- An endorsed measure does not address the specific topic of interest in the proposed
 - 531 measure;
 - 532 • A critical timeline must be met (e.g., legislative mandate); and
 - 533 • The measure is not complex (e.g., composite, requires risk adjustment).

534 In addition to passing the criterion, *Importance to Measure and Report*, untested measures must
535 demonstrate an adequate foundation for both reliability and validity as follows. That is,
536 measures should be precisely specified and have at least minimal face validity. Measures that
537 do not meet these minimum requirements are not ready for testing and should not be
538 recommended for time-limited endorsement.

539

540

541

542

543

544

545 Table 5. Minimum Requirements for Untested Measures under Scientific Acceptability of Measure Properties

Foundation for Reliability	Foundation for Validity
<p>All measure specifications (e.g., numerator, denominator, exclusions, risk factors, scoring) are unambiguous and likely to consistently 1) identify who is included and excluded from the target population; 2) identify the event, condition, or outcome being measured; 3) compute the measure score; etc.</p> <p>All EHR measure specifications are unambiguous and include only data elements from the quality data set (QDS)* including quality data elements, code lists, and measure logic OR new elements are submitted for inclusion to the QDS</p>	<p>The measure specifications (e.g., numerator, denominator, exclusions, risk factors, scoring) reflect the quality of care problem (1a,1b) and evidence cited in support of the measure focus (1c) under <i>Importance to Measure and Report</i></p>

546 *QDS elements should be used when available. When needed quality data elements are not yet available in the QDS,
 547 they will be considered for addition to the QDS.
 548

549 V. Recommendations for Testing Required for Maintenance of Endorsement

550 The above guidance on testing and evidence of reliability and validity for initial endorsement
 551 decisions applies to testing required for endorsement maintenance with a few modifications.
 552 With the [current NOF system of endorsement cycles](#), endorsed measures will be reviewed for
 553 maintenance of endorsement every three years along with new measures. Both new and
 554 endorsed measures will be required to meet the measure evaluation criteria, including
 555 reliability and validity.

556
 557 The Task Force agreed that reliability and validity should be evaluated when measures are
 558 reviewed for maintenance of endorsement. Several considerations were relevant to the task
 559 force deliberations on this subject, including: recognizing that reliability and validity are not
 560 static properties, no one test is definitive, evidence accumulates over time, and the proposed
 561 rating system permits endorsement of measures that have limited evidence of reliability and
 562 validity (moderate rating). However, developers cannot be expected to monitor both reliability
 563 and validity indefinitely, once these measure properties have been well established.

564
 565 Table 6 outlines the expectations for reliability and validity testing for review at the time of
 566 endorsement maintenance. At the time of review for endorsement maintenance, reliability and
 567 validity testing should: a) use data from implementation of the endorsed measure as specified,
 568 and b) focus on the measure score rather than data elements. Of particular relevance to a
 569 measure in use is information on the accuracy of any classification based on the measure results.
 570 If an endorsed measure has not been implemented, expanded testing in terms of scope and

571 levels is required. The rating system provided in Table 2 and Table 3 also applies to the
 572 maintenance review. As with initial endorsement, all the other criteria also will be used to
 573 determine whether a measure warrants continued endorsement.

574

575 Table 6. 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 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 validity achieved a high rating

576

577 **VI. Recommendations for Modifications to the NQF Evaluation Criteria**

578 The recommendations of the Task Force as described above resulted in some wording changes
 579 to the NQF measure evaluation criteria, but the intent remains unchanged. Criterion 2, *Scientific*
 580 *Acceptability of Measure Properties*, is primarily about reliability and validity and threats to
 581 reliability and validity. This criterion can be simplified by focusing on the concepts of reliability
 582 and validity and arranging the subcriteria to reflect their relationship to reliability or validity as
 583 follows.

584 **2a. Reliability**

585 2a1. Precise specifications (previously 2a) including exclusions (previously 2d)

586 2a2. Reliability testing (previously 2b) – data elements or measure score

587 **2b. Validity**

588 2b1. Specifications consistent with evidence

- 589 2b2.Validity testing (previously 2c) – data elements or measure score
- 590 2b3.Justification of exclusions (previously 2d) – relates to evidence
- 591 2b4.Risk adjustment (previously 2e)
- 592 **2b5.**Identification of differences in performance (previously 2f)
- 593 2b6.Comparability of data sources/methods (previously 2g)
- 594 **2c.Disparities** (previously 2h)

595
596 Table 7. Current and Modified Measure Evaluation Criteria

Current Measure Evaluation Criteria	Modified Measure Evaluation Criteria
<p>2. Scientific acceptability of the measure properties: Extent to which the measure, <u>as specified</u>, produces consistent (reliable) and credible (valid) results about the quality of care when implemented.</p> <p><i>[See footnotes below the criteria]</i></p> <p>Reliability</p> <p>2a. The measure is well defined and precisely specified ⁶ so that it can be implemented consistently within and across organizations and allow for comparability. The required data elements are of high quality as defined by NQF's Health Information Technology Expert Panel (HITEP)).⁷</p> <p>2b. Reliability testing ⁸ demonstrates the measure results are repeatable, producing the same results a high proportion of the time when assessed in the same population in the same time period.</p> <p>Validity</p> <p>2c. Validity testing ⁹ demonstrates that the measure reflects the quality of care provided, adequately distinguishing good and poor quality. If face validity is the only validity addressed, it is systematically assessed.</p> <p>2e. For outcome measures and other measures (e.g., resource use) when indicated:</p> <ul style="list-style-type: none"> • an evidence-based risk-adjustment strategy (e.g., risk models, risk stratification) is specified and is based on patient clinical factors that influence the measured outcome (but not disparities in care) and are present at start of care ^{11,13} <p>OR</p> <ul style="list-style-type: none"> • rationale/data support no risk adjustment. 	<p>2. Scientific acceptability of the measure properties: Extent to which the measure, <u>as specified</u>, produces consistent (reliable) and credible (valid) results about the quality of care when implemented.</p> <p><i>[See footnotes below the criteria]</i></p> <p>2a. Reliability</p> <p>2a1. The measure is well defined and precisely specified ⁶ so that it can be implemented consistently within and across organizations and allow for comparability. EHR measure specifications are based on the quality data set (QDS).⁷</p> <p>2a2. Reliability testing ⁸ 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 the measure score is precise.</p> <p>2b. Validity</p> <p>2b1. The measure specifications ⁶ are consistent with the evidence presented to support the focus of measurement under criterion 1c. The measure is specified to capture the most inclusive target population indicated by the evidence and exclusions are supported by the evidence.</p> <p>2b2. Validity testing ⁹ demonstrates that the measure data elements are correct and/or the measure score correctly reflects the quality of care provided, adequately identifying differences in quality.</p> <p>2b3. Exclusions are supported by the clinical evidence otherwise, they are supported by evidence ¹⁰ of sufficient frequency of occurrence so that results are distorted without the exclusion;</p>

Current Measure Evaluation Criteria	Modified Measure Evaluation Criteria
<p>2f. Data analysis demonstrates that methods for scoring and analysis of the specified measure allow for identification of statistically significant and practically/clinically meaningful ¹⁴ differences in performance</p> <p>2g. If multiple data sources/methods are allowed, there is demonstration they produce comparable results.</p> <p>2d. Clinically necessary measure exclusions are identified and must be: supported by evidence ¹⁰ of sufficient frequency of occurrence so that results are distorted without the exclusion; AND</p> <ul style="list-style-type: none"> • a clinically appropriate exception (e.g., contraindication) to eligibility for the measure focus ¹¹; <p>AND</p> <ul style="list-style-type: none"> • precisely defined and specified: <ul style="list-style-type: none"> – if there is substantial variability in exclusions across providers, the measure is specified so that exclusions are computable and the effect on the measure is transparent (i.e., impact clearly delineated, such as number of cases excluded, exclusion rates by type of exclusion); – if patient preference (e.g., informed decision-making) is a basis for exclusion, there must be evidence that it strongly impacts performance on the measure and the measure must be specified so that the information about patient preference and the effect on the measure is transparent ¹² (e.g., numerator category computed separately, denominator exclusion category computed separately). <p>2h. 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); OR rationale/data justifies why stratification is not necessary or not feasible.</p> <p>Footnotes 6 Measure specifications include the target population (e.g., denominator) to whom the measure applies, identification of those from the target population who achieved the specific measure focus (e.g., numerator), measurement time window,</p>	<p>AND</p> <ul style="list-style-type: none"> – Measure specifications for scoring include computing exclusions so that the effect on the measure is transparent (i.e., impact clearly delineated, such as number of cases excluded, exclusion rates by type of exclusion); <p>AND</p> <ul style="list-style-type: none"> – If patient preference (e.g., informed decision-making) 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 ¹² (e.g., numerator category computed separately, denominator exclusion category computed separately). <p>2b4. For outcome measures and other measures when indicated (e.g., resource use):</p> <ul style="list-style-type: none"> • an evidence-based risk-adjustment strategy (e.g., risk models, risk stratification) is specified; is based on 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; ^{11,13} and has demonstrated adequate discrimination and calibration <p>OR</p> <ul style="list-style-type: none"> • rationale/data support no risk adjustment/stratification. <p>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 meaningful¹⁴ differences in performance; OR there is evidence of overall less than optimal performance.</p> <p>2b6. If multiple data sources/methods are specified, there is demonstration they produce comparable results.</p> <p>2c. 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); OR rationale/data justifies why stratification is not necessary or not feasible.</p> <p>Footnotes 6 Measure specifications include the target population</p>

Current Measure Evaluation Criteria	Modified Measure Evaluation Criteria
<p>exclusions, risk adjustment, definitions, data elements, data source and instructions, sampling, scoring/computation.</p> <p>7 The HITEP criteria for high quality data include: a) data captured from an authoritative/accurate source; b) data are coded using recognized data standards; c) method of capturing data electronically fits the workflow of the authoritative source; d) data are available in EHRs; and e) data are auditable. NQF. <i>Health Information Technology Expert Panel Report: Recommended Common Data Types and Prioritized Performance Measures for Electronic Healthcare Information Systems</i>. Washington, DC: NQF; 2008.</p> <p>8 Reliability testing may address the data items or final measure score. Examples of reliability testing 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.</p> <p>9 Examples of validity testing include, but are not limited to: determining if measure scores adequately distinguish between providers known to have good or poor quality assessed by another valid method; correlation of measure scores with another valid indicator of quality for the specific topic; ability of measure scores to predict scores on some other related valid measure; content validity for multi-item scales/tests. Face validity is a subjective assessment by experts of whether the measure reflects the quality of care (e.g., whether the proportion of patients with BP < 140/90 is a marker of quality). If face validity is the only validity addressed, it is systematically assessed (e.g., ratings by relevant stakeholders) and the measure is judged to represent quality care for the specific topic and that the measure focus is the most important aspect of quality for the specific topic.</p> <p>10 Examples of evidence that an exclusion distorts measure results include, but are not limited to: frequency of occurrence, sensitivity analyses with and without the exclusion, and variability of exclusions across providers.</p> <p>11 Risk factors that influence outcomes should not be specified as exclusions.</p> <p>12 Patient preference is not a clinical exception to eligibility and can be influenced by provider interventions.</p> <p>13 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, gender (e.g., poorer treatment outcomes of African American men with prostate cancer, inequalities in treatment for CVD risk factors between men and women). It is preferable to stratify measures by race and socioeconomic status rather than adjusting out differences.</p> <p>14 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% v. 75%) is clinically</p>	<p>(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.</p> <p>7 EHR measure specifications include data type from the QDS, code lists, EHR field, measure logic, original source of the data, recorder, and setting.</p> <p>8 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).</p> <p>9 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.</p> <p>10 Examples of evidence that an exclusion distorts measure results include, but are not limited to: frequency of occurrence, sensitivity analyses with and without the exclusion, and variability of exclusions across providers.</p> <p>11 Risk factors that influence outcomes should not be specified as exclusions.</p> <p>12 Patient preference is not a clinical exception to eligibility and can be influenced by provider interventions.</p> <p>13 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, gender (e.g., poorer treatment outcomes of African American men with prostate cancer, inequalities in treatment for CVD risk factors between men and women). It is preferable to stratify measures by race and socioeconomic status rather than adjusting out differences.</p> <p>14 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% v. 75%) is clinically</p>

Current Measure Evaluation Criteria	Modified Measure Evaluation Criteria
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 poor performance may not demonstrate much variability across providers.	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.

597

598 **VII. Recommendations for the Measure Submission**

599 The prior recommendations resulted in modest changes to the information that is currently
600 requested on the [measure submission form](#). The numbering system will need to be adjusted as
601 appropriate for the reorganization of the subcriteria noted above and the online submission and
602 measures database.

603

604 **Measure Specifications (Measure evaluation criterion 2a)**

605 **2a.1. Numerator Statement** *(Brief narrative description of the numerator - what is being measured about*
606 *the target population, e.g., target condition, event, or outcome)*

607 **2a.2. Numerator Time Window** *(The time period in which cases are eligible for inclusion in the numerator)*

608 **2a.3. Numerator Details** *(All information required to collect the data required to calculate the numerator,*
609 *including definitions and all codes with descriptors)*

610 **2a.4. Denominator Statement** *(Brief narrative description of the denominator - target population being*
611 *measured)*

612 **2a.5. Target Population Gender**

613 Female Male

614 **2a.6. Target Population Age Range**

615 **2a.7. Denominator Time Window** *(The time period in which cases are eligible for inclusion in the*
616 *denominator)*

617 **2a.8. Denominator Details** *(All information required to collect the data required to calculate the*
618 *denominator, including definitions and all codes with descriptors)*

619 **2a.9. Denominator Exclusions** *(Brief narrative description of exclusions from the target population)*

620 **2a.10. Denominator Exclusion Details** *(All information required to collect the data required for exclusions*
621 *to the denominator, including all definitions and codes with descriptors)*

622 **2a.11. Stratification Details/Variables** *(All information required to stratify the measure including the*
623 *stratification variables, all definitions and codes with descriptors)*

624 **2a.12. Risk Adjustment/Stratification Type**

625 No risk adjustment/stratification necessary - measure is not an outcome or resource use measure

626 No risk adjustment/stratification necessary - rationale and analysis provided in Section 2e

627 Stratification/analysis by subgroup - see variables in 2a.11

628 Statistical risk model - specifications 2a.14

629 Other (specify)

630 **2a.14. Specifications for Statistical Risk Model and Variables Included** *(Name the statistical method*
631 *(e.g., logistic regression) and list the risk model variables all definitions and codes with descriptors.*
632 *Development and testing are reported in Section 2e)*

633 **2a.15. Detailed Risk Model** *(Please provide a web page URL or attachment. NQF strongly prefers URLs.*
634 *Attach documents only if they are not available on a web page and keep attached file to 5 MB or less.)*

635 **2a.18. Type of Score**

636 count

637 frequency distribution

638 non-weighted score/ composite/scale

639 rate/proportion
640 ratio
641 weighted score/ composite/scale
642 categorical
643 continuous variable
644 Other (please indicate)

645 **2a.20. Interpretation of Score** (*Classifies interpretation of score according to whether better quality is*
646 *associated with a higher score, a lower score, a score falling within a defined interval, or a passing score)*
647 better quality= higher score
648 better quality = lower score
649 better quality = score within a defined interval
650 passing score defines better quality

651 **2a.21. Measure Score Calculation Algorithm** (*Describe the calculation of the measure score as a series of*
652 *steps, including identification of denominator, exclusions, identification of numerator, stratification or*
653 *adjustment, and classification category)*

654 **2a.22. Measure Algorithm or Flow Diagram** (*Please provide a web page URL or attachment. NQF strongly*
655 *prefers URLs. Attach documents only if they are not available on a web page and keep attached file to 5 MB*
656 *or less.*

657 **2a.23. Sampling (Survey) Methodology**
658 If measure is based on a sample (or survey), provide instructions for obtaining the sample, conducting the
659 survey, and guidance on minimum sample size (response rate).

660 **2a.24. Data Type** (*Check the sources for which the measure is specified and tested*)
661 Simplify, and allow for indication if electronic

662 Documentation of original self-assessment	Paper medical record/flow-sheet
663 Electronic administrative data/claims	Pharmacy data
664 Electronic clinical data	Public health data/vital statistics
665 Electronic Health/Medical Record	Registry data
666 External audit	Special or unique data
667 Lab data	Survey: Patient
668 Management data	Survey: Provider
669 Organizational policies and procedures	

670 **2a.25. Data Source or Collection Instrument** (*Name the specific data source or data collection instrument,*
671 *E.g. name of database, clinical registry, collection instrument, etc.)*

672 **2a.26. Data Source or Collection Instrument Reference** (*Please provide a web page URL or attachment.*
673 *NQF strongly prefers URLs. Attach documents only if they are not available on a web page and keep*
674 *attached file to 5 MB or less.)*

675 **2a.29. Data Dictionary or Code Table** (*Please provide a web page URL or attachment. NQF strongly prefers*
676 *URLs. Attach documents only if they are not available on a web page and keep attached file to 5 MB or*
677 *less.)*

678 **2a.32. Level of Measurement/Analysis** (*Check the level for which the measure is specified and tested*)

679 Clinicians	689 Regional/network
680 Individual	690 States
681 Group	691 Counties or cities
682 Other	692 Prescription drug plan
683 Facility/agency	693 Program
684 Health plan	694 Disease management
685 Integrated delivery system	695 Quality improvement organization (QIO)
686 Multi-site/corporate chain	696 Other
687 Population	697 Can be measured at all levels
688 National	698 Other

699 **2a.36. Care Setting** (*Check the settings for which the measure is specified and tested; check all that*
700 *apply.)*

701 Ambulatory Care Home	705 Emergency Department Nursing home (NH) /skilled nursing
702 Ambulatory surgery center Hospice	706 facility (SNF)
703 Office Hospital	707 Hospital Outpatient Rehabilitation facility
704 Clinic Long term acute care hospital	708 Assisted living
	709 Behavioral health/psychiatric unit All settings

710	Dialysis facility Unspecified or "not applicable"	712	Group homes
711	Emergency medical services/ambulance	713	Other
714	2a.38. Clinical Services (<i>Healthcare services being measured; check all that apply.</i>)		
715	Behavioral Health	728	Nurses
716	Mental health	729	Optometrist
717	Substance use treatment	730	PA/NP/Advanced Practice Nurse
718	Other	731	Pharmacist
719	Clinicians (<i>Continued</i>)	732	Physicians (MD/DO)
720	Podiatrist	733	Respiratory Therapy
721	Psychologist/LCSW	734	Other
722	PT/OT/Speech	735	Dialysis
723	Clinicians	736	Home health
724	Audiologist	737	Hospice/palliative care
725	Chiropractor	738	Imaging
726	Dentist/Oral surgeon	739	Laboratory
727	Dietician/Nutritional professional	740	Other

- 741
- 742 **Reliability Testing (Measure evaluation criterion 2b)**
- 743 **2b.1. Data/Sample** (*Description of data/sample and size*)
- 744 **2b.2. Analytic Methods** (*Method of reliability testing and rationale*)
- 745 **2b.3. Testing Results** (*Reliability statistics, assessment of adequacy in the context of norms for the test conducted*)
- 746

- 747
- 748 **Validity Testing (Measure evaluation criterion 2c)**
- 749 **2c.1. Data/Sample** (*Description of data/sample and size*)
- 750 **2c.2. Analytic Method** (*Method of validity testing and rationale*)
- 751 **2c.3. Testing Results** (*Statistical results, assessment of adequacy in the context of norms for the test conducted*)
- 752

- 753
- 754 **Measure Exclusions (Measure evaluation criterion 2d)**
- 755 **2d.1. Summary of Evidence Supporting Exclusion(s)**
- 756 **2d.2. Citations for Evidence**
- 757 **2d.3. Data/Sample** (*Description of data/sample and size*)
- 758 **2d.4. Analytic Method** (*Type of analysis and rationale*)
- 759 **2d.5. Testing Results** (*e.g., frequency, variability, sensitivity analyses of impact on measure scores*)
- 760

- 761 **Risk Adjustment Strategy (Measure evaluation criterion 2e)**
- 762 **2e.1. Data/Sample** (*Description of data/sample and size used for development and validation*)
- 763 **2e.2. Analytic Method** (*Describe methods for development and testing of risk model including selection of risk factors*)
- 764
- 765 **2e.3. Testing Results** (*Quantitative assessment of relative contribution of model risk factors; Risk model performance metrics including cross-validation calibration and discrimination statistics, and assessment of adequacy in the context of norms for risk models. Provide calibration curve and risk decile plot in attachment.*)
- 766
- 767
- 768
- 769 **2e.4. If outcome or resource use measure is not risk adjusted, provide rationale**
- 770

- 771 **Identification of Meaningful Differences in Performance (Measure evaluation criterion 2f)**
- 772 **2f.1. Data/Sample from Testing or Current Use** (*Description of data/sample and size*)
- 773 **2f.2. Methods to Identify Statistically Significant and Practical or Meaningful Differences in Performance** (*Type of analysis and rationale*)
- 774
- 775 **2f.3. Measure Scores from Testing or Current Use** (*Description of scores, e.g., distribution by quartile, mean, median, SD, etc.; identification of statistically significant and meaningful differences in performance. If no variability, discuss rationale for performance measurement, e.g., benchmark for determining overall poor performance.*)
- 776
- 777
- 778
- 779
- 780

781 **Comparability of Multiple Data Sources/Methods (Measure evaluation criterion 2g)**
782 **2g.1. Data/Sample** (*Description of data/sample and size*)
783 **2g.2. Analytic Method** (*Type of analysis and rationale*)
784 **2g.3. Testing Results** (*Statistical results, assessment of adequacy in the context of norms for the test*
785 *conducted*)

786
787 **Disparities in Care (Measure evaluation criterion 2h)**
788 **2h.1. If measure is stratified to identify disparities, provide stratified results** (*Scores by stratified*
789 *categories/cohorts*)
790 **2h.2. If disparities have been reported/identified but measure is not specified to detect disparities,**
791 **provide follow-up plans**

792
793

794 REFERENCES

- 795 1. McGlynn EA. Selecting common measures of quality and system performance. *Med Care.*
796 2003;41(1 Suppl):I39-I47.
- 797 2. McGlynn EA, Asch SM. Developing a clinical performance measure. *Am J Prev Med.*
798 1998;14(3 Suppl):14-21.
- 799 3. Rubin HR, Pronovost P, Diette GB. From a process of care to a measure: the development
800 and testing of a quality indicator. *Int J Qual Health Care.* 2001;13(6):489-496.
- 801 4. Trochim WMK. Research methods knowledge base. *Web Center for Social Research Methods*
802 2006; Available at: <http://www.socialresearchmethods.net/kb/index.php>. Last accessed
803 May 2010.
- 804 5. Physician Consortium for Performance Improvement. *Measure Testing Protocol for Physician*
805 *Consortium for Performance Improvement Performance Measures.* Chicago, IL: American
806 Medical Association; 2007.
- 807 6. Bhattacharyya T, Freiberg AA, Mehta P et al. Measuring the report card: the validity of pay-
808 for-performance metrics in orthopedic surgery. *Health Aff (Millwood).* 2009;28(2):526-532.
- 809 7. Schneider EC, Nadel MR, Zaslavsky AM et al. Assessment of the scientific soundness of
810 clinical performance measures: a field test of the National Committee for Quality
811 Assurance's colorectal cancer screening measure. *Arch Intern Med.* 2008;168(8):876-882.
- 812 8. Moss PA. Can There Be Validity without Reliability? *Educational Researcher.* 1994;23(2):5-12.
- 813 9. Salvucci S, Walter E, Conley V, Fink S, Saba M. *Measurement Error Studies at the National*
814 *Center for Education Statistics.* Washington, DC: =U.S. Department of Education; 1997.
- 815 10. Kaplan SH, Griffith JL, Price LL et al. Improving the reliability of physician performance
816 assessment: identifying the "physician effect" on quality and creating composite measures.
817 *Med Care.* 2009;47(4):378-387.
- 818 11. Reeves D, Campbell SM, Adams J et al. Combining multiple indicators of clinical quality: an
819 evaluation of different analytic approaches. *Med Care.* 2007;45(6):489-496.
- 820 12. Fitch K, Bernstein SJ, Aguilar MS et al. *The RAND/UCLA Appropriateness Method User's*
821 *Manual.* Santa Monica, CA: RAND Health; 2000. Available at
822 http://www.rand.org/pubs/monograph_reports/MR1269/.
- 823 13. Spertus JA, Eagle KA, Krumholz HM et al. American College of Cardiology and American
824 Heart Association methodology for the selection and creation of performance measures
825 for quantifying the quality of cardiovascular care. *Circulation.* 2005;111(13):1703-1712.
- 826 14. National Quality Forum. *Health Information Technology Expert Panel II - Health IT Enablement*
827 *of Quality Measurement.* Washington, DC: NQF; 2009.

- 828 15. Baker DW, Persell SD, Thompson JA et al. Automated review of electronic health records to
829 assess quality of care for outpatients with heart failure. *Ann Intern Med.* 2007;146(4):270-
830 277.
- 831 16. Persell SD, Wright JM, Thompson JA et al. Assessing the validity of national quality
832 measures for coronary artery disease using an electronic health record. *Arch Intern Med.*
833 2006;166(20):2272-2277.
- 834 17. Weiner M, Stump TE, Callahan CM et al. Pursuing integration of performance measures
835 into electronic medical records: beta-adrenergic receptor antagonist medications. *Qual Saf*
836 *Health Care.* 2005;14(2):99-106.
- 837 18. Briggs JB, Kind EA, Awwad S et al. *Performance Measures Using Electronic Health Records: Five*
838 *Case Studies.* New York, NY: The Commonwealth Fund; 2008. Report No.: 1132. Available
839 at www.commonwealthfund.org.
- 840 19. Austin PC. The reliability and validity of Bayesian measures for hospital profiling: a Monte
841 Carlo assessment. *Journal of Statistical Planning and Inference.* 2005;128(1):109-122.
- 842 20. Kahn JM, Iwashyna TJ. Accuracy of the discharge destination field in administrative data for
843 identifying transfer to a long-term acute care hospital. *BMC Res Notes.* 2010;3:205.
- 844 21. Quan H, Parsons GA, Ghali WA. Validity of procedure codes in International Classification
845 of Diseases, 9th revision, clinical modification administrative data. *Med Care.*
846 2004;42(8):801-809.
- 847 22. McGinn T, Wyer PC, Newman TB et al. Tips for learners of evidence-based medicine: 3.
848 Measures of observer variability (kappa statistic). *CMAJ.* 2004;171(11):1369-1373.
- 849 23. Viera AJ, Garrett JM. Understanding interobserver agreement: the kappa statistic. *Fam Med.*
850 2005;37(5):360-363.
- 851 24. Landis J, Koch G. The measurement of observer agreement for categorical data. *Biometrics.*
852 1977;33:159-174.
- 853 25. Zaslavsky AM. Statistical issues in reporting quality data: small samples and casemix
854 variation. *Int J Qual Health Care.* 2001;13(6):481-488.
- 855 26. Nunnally J, Bernstein I. *Psychometric Theory.* Third Edition ed. New York: McGraw-Hill;
856 1994.
857
858

859 APPENDIX A – COMMON APPROACHES TO TESTING

860 Tables A-1 through A-5 provide information on the various types of reliability and validity
 861 testing that *could* be performed. The information in the following tables is not meant to provide
 862 an exhaustive list of methods. Other approaches to testing may be appropriate and could be
 863 used if the method and rationale are explained and judged to be appropriate. Measure
 864 developers should select the testing that is appropriate and feasible for the measure being
 865 developed and that will meet at least the moderate rating as described in Table 2. Likewise,
 866 measure developers should identify the potential threats to validity for the specific measure and
 867 conduct analyses to demonstrate adequate control.

- 868 [Table A-1 Reliability Testing at the Level of the Computed Performance Measure Score](#).. **Error!**
 869 **Bookmark not defined.**
 870 [Table A-2. Reliability Testing at the Level of the Data Elements](#) **Error! Bookmark not defined.**
 871 [Table A-3. Validity Testing at the Level of the Computed Performance Measure Score](#) ... **Error!**
 872 **Bookmark not defined.**
 873 [Table A-4. Validity Testing at the Level of Data Elements](#) **Error! Bookmark not defined.**
 874 [Table A-5 Testing Related to Threats to Validity](#) **Error! Bookmark not defined.**
 875 [Table A-6. Interpretation of Statistical Results](#)..... **Error! Bookmark not defined.**
 876

877 Table A-1 Reliability Testing at the Level of the Computed Performance Measure Score

Reliability Testing – Measure Score	
Data	Aspect of Reliability/Test
Reliability testing of the computed <u>measure score</u> does not vary by type of data or type of measure Requires data for the computed measure scores and the individual patient-level data for the measured entities	<p>Statistical reliability (precision) of sample average as an estimate of the underlying population average</p> <p>Analysis of the relative value of variation in measure scores due to signal (i.e., variation between measured entities) versus noise (i.e., variation within measured entities) using statistical tests such as Analysis of Variance (ANOVA), Intraclass Correlation Coefficient (ICC), or variance components from a multi-level mixed model {references}</p> <p>Monte Carlo simulation to test Bayesian measures ¹⁹</p> <p>Generalizability analysis based on generalizability theory on the sources of variation {reference}</p> <p>Other: Other methods may be appropriate and rationale for method chosen should be provided</p>

878

879 Table A-2. Reliability Testing at the Level of the Data Elements

Reliability Testing – Data elements	
Separate reliability testing of the data elements is not required if validity testing conducted on the data elements.	Empirical validity testing of the data elements (see Table A-4) is conducted and demonstrates the data elements are valid
Prior evidence of reliability of data elements can be used for evidence of reliability of data elements.	Prior evidence could include published or unpublished testing that: <ul style="list-style-type: none"> • included the same data elements; and • used the same data type; and • was conducted on a sample as described above (i.e., representative, adequate numbers, and randomly selected, if possible).
Data Type	Aspect of Reliability/Test
<u>Retrospective chart abstraction (including registry data abstracted retrospectively from medical records)</u>	Inter-rater reliability between abstractors Analysis of agreement using appropriate statistical tests (e.g., kappa, ICC) with 2 nd abstractor on each critical data element and computed measure score
Administrative claims data where codes that are used to represent the primary clinical data (ICD, CPT, CPT-II/G)	Inter-rater reliability between coders Analysis of agreement using appropriate statistical tests (e.g., kappa, ICC) with a 2 nd coder on each critical data element and computed measure score
<u>Standardized clinical patient information (MDS, OASIS, registry, potentially some aspects of EHRs) collected by an authoritative source concurrently with care delivery (not abstracted, coded, or transcribed by another person)</u>	Inter-rater reliability between assessors Analysis of agreement using appropriate statistical tests (e.g., kappa, ICC) with 2 nd assessor on each critical data element and computed measure score
EHR clinical record information	Data elements obtained with EHR specifications and data exported electronically from EHRs according to standards are repeatable (reliable) when applied to the same population in the same time period
Survey – single items	Test-retest reliability Analysis of agreement between two administrations of the same items (time frame long enough so as not to remember and short enough so as not to have changed)
Instrument/scale	If patient scores from an instrument/scale are used in constructing a performance measure, generally the reliability of the scale has already been tested and documented and can be used as evidence of data element reliability. Internal consistency reliability (Cronbach’s alpha) Analysis of the extent to which item responses obtained at the same time correlate highly with each other Generalizability analysis based on generalizability theory on the sources of variation {reference}
Other data type	Rationale should be provided for method chosen to demonstrate reliability

880

881

882 Table A-3. Validity Testing at the Level of the Computed Performance Measure Score

Validity Testing—Measure Score	
Data	Aspect of Validity/Test
<p>Validity testing of the computed <u>measure score</u> does not vary by type of data or type or type of measure</p> <p>Requires data for the computed measure scores for the measured entities and other data as necessary for the chosen validity study</p>	<p>Evidence that supports the intended interpretation of measure scores for the intended purpose – making conclusions about the quality of care</p> <p>Systematic testing of face validity of the <u>measure score</u> as a quality indicator by identified experts, explicitly addressed the question of whether <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> (using a systematic and transparent process, e.g., modified Delphi, formal consensus process, RAND Appropriateness Method¹², ACC/AHA method)¹³ with methods and results reported for review.</p> <p>Criterion Validity: Studies to assess the correlation of the computed measure score against some criterion determined to be valid.</p> <p>Concurrent – Correlation with another measure of the same construct measured at the same time</p> <p>Predictive – Correlation with another measure of the same construct or an outcome measured at some time in the future</p> <p>Construct Validity: Studies to assess how the measure performs based on the theory of the construct.</p> <p>Contrasted Groups – Study to assess the ability of the measure score to distinguish between groups that it should theoretically be able to distinguish between</p> <p>Convergent – Study to examine the degree to which the measure score is similar to (converges on) other measures of the same construct or that it theoretically should be similar to</p> <p>Discriminative – Study to examine the degree to which the measure score is not similar to (diverges from) other measures that it theoretically should not be similar to</p> <p>Other: Other methods may be appropriate and rationale for method chosen should be provided</p>

883

884

885 Table A-4. Validity Testing at the Level of Data Elements

Validity Testing – Data elements	
Prior evidence of validity of data elements can be used for evidence of validity of data elements.	<p>Prior evidence could include published or unpublished testing that:</p> <ul style="list-style-type: none"> • included the same data elements; and • used the same data type ; and • was conducted on a sample as described above (i.e., representative, adequate numbers, and randomly selected, if possible).
Data Type	Aspect of Validity/Test
Retrospective chart abstraction (including registry data abstracted retrospectively from medical records)	<p>Validity of data elements abstracted from medical record as compared to some criterion authoritative source of the same data</p> <p>Analysis of agreement using appropriate statistical tests (e.g., sensitivity, specificity, positive predictive value, negative predicted value ^{20, 21} with some other source of the same information considered to be valid (e.g., original data collection such as survey or observation, vital statistics)</p>
Administrative claims data where codes that are used to represent the primary clinical data (ICD, CPT, CPT-II/G)	<p>Validity of coded data from claims as compared to some criterion authoritative source of the same data</p> <p>Analysis of agreement using appropriate statistical tests (e.g., sensitivity, specificity, positive predictive value, negative predicted value) with manual abstraction from the <u>full</u> medical record as the authoritative source</p>
Standardized clinical patient information (MDS, OASIS, registry, potentially some aspects of EHRs) <u>collected by an authoritative source concurrently with care delivery</u> (not abstracted, coded, or transcribed by another person)	<p>Validity of data elements from standardized assessment instruments as compared to some criterion authoritative source of the same data</p> <p>Analysis of agreement using appropriate statistical tests (e.g., sensitivity, specificity, positive predictive value, negative predicted value) with “expert” assessor (conducted at approximately the same time)</p> <p>Predictive validity as described in Table A-3 (e.g., patient-level assessment item or score predicts a subsequent outcome of undisputed importance, such as death or permanent disability)</p>
EHR clinical record information	<p>Validity of data elements extracted from specified fields in EHRs as compared to some criterion authoritative source of the same data</p> <p>Analysis of agreement using appropriate statistical tests (e.g., sensitivity, specificity, positive predictive value, negative predicted value) with data elements abstracted from the <u>entire</u> EHR (not just the fields where the data are expected)</p> <p>Demonstration of agreement between data elements and scores obtained by applying the EHR measure specifications to a simulated test EHR data set that reflects standards for EHRs and includes sample patient data with known values for the data elements needed for the specified measure and computed measure score.</p>
Survey – single items	<p>Validity of data elements from survey as compared to some criterion authoritative source of the same data</p> <p>Analysis of agreement using appropriate statistical tests (e.g., sensitivity, specificity, positive predictive value, negative predicted value) with some other source of the same information considered to be valid (e.g., medical record, vital statistics)</p>
Instrument/scale	<p>If patient scores from an instrument/scale are used in constructing a performance measure, generally the validity of the scale has already been tested and documented and can be used as evidence of data element validity.</p>

Validity Testing – Data elements	
	<p>Validity of the content of the items in an instrument or scale Systematic assessment by subject matter experts that the content of the instrument/scale is representative of the domain being measure</p> <p>Validity of whether the instrument is consistent with the theoretical construct Confirmatory factor analysis</p> <p>Criterion or construct validity of the patient-level score as described in Table A-3 (e.g., patient-level score predicts a subsequent outcome of undisputed importance, such as death or permanent disability)</p>
Other data type	Rationale should be provided for method chosen to demonstrate validity

886

887 Table A-5 Testing Related to Threats to Validity

Threat to Validity	Testing/Analysis
Threat that differences in measure scores are due to differences in severity of conditions of patients served rather than differences in quality (confounding bias)	For outcome and resource use measures, empirical evidence for the adequacy of adjustment for patient factors (analysis of risk factors, discrimination and calibration of risk models); OR evidence that risk adjustment/ stratification is not necessary for fair comparisons (patient outcomes do not vary by patient characteristics)
Threat of bias from differences in data type and/or differences in data collection practices; (information bias)	If multiple data sources (e.g., medical record and claims) or methods (e.g., mail survey and interview) are specified, empirical evidence that resulting measure scores are comparable (analysis of agreement between scores based on different data sources)
Threat of bias from missing or “incorrect” data; or exclusions (selection/attrition bias)	Sensitivity analysis of the impact of missing or “incorrect” data on resulting measure scores (analysis of patterns of missing data; simulate missing data or “incorrect” data and analyze impact on measure scores) Analyses of frequency of exclusions, sensitivity analyses with and without the exclusion, and variability of exclusions across providers

888

889

890

891 Table A-6. Interpretation of Statistical Results

Test	Interpretation
Kappa ^{22, 23} Measure of agreement between two raters that adjusts for chance agreements for categorical data (nominal, ordinal)	Kappa values range between 0 and 1. 0 and are interpreted as degree of agreement beyond chance 0 No better than chance 0.01-0.20 Slight 0.21-0.40 Fair 0.41-0.60 Moderate 0.61-0.80 Substantial 0.81-1.0 Almost perfect ²⁴
ICC Alternative measure of agreement when more than two raters or quantitative data (interval, ratio)	ICC values range between 0 and 1.0 Interpretations are similar for kappa noted above ICC approaches 1.0 only if there is no variance due to raters
ANOVA or ICC Used for signal-to-noise analysis for estimated mean (or proportion) – analysis of variance <u>between</u> the measured entities (signal) to variance <u>within</u> the measured entities (noise)	F test of equality of means for measured entities; F-1 is an estimate of the ratio of signal to noise, and $[1-(1/F)]$ estimates the fraction of total variance that is due to signal (real variation among measured entities), referred to as interunit reliability (IUR). When F is large, IUR is close to 1 indicating almost all signal and no noise. Zaslavsky ²⁵ demonstrated that value of F should be 10 or greater.
Cronbach’s alpha Measure of the average correlation of the items comprising a scale or subscale	A widely-accepted cut-off is .70 or higher ²⁶ for a set of items to be considered a scale. Some use .75 or .80 while others are as lenient as .60. That .70 is as low as one may wish to go is reflected in the fact that when alpha is .70, the standard error of measurement will be over half (0.55) a standard deviation. {reference}
Pearson Correlation Measure of the degree of association (not agreement) between two quantitative variables	Values range from -1 to +1 The squared correlation represents the proportion of variance shared by the two variables (e.g., correlation of 0.5 represents 25% shared variance). Interpretation depends on statistical significance, size, and context. For example, two measures of the same thing using different methods would have very high correlations (>0.9).
Spearman (rank order) correlation Measure of the degree of association (not agreement) for rank-order variables	Values range from -1 to +1 A high positive value indicates a strong tendency for the paired ranks to be similar; a low negative indicates the paired ranks to be opposite.

892

APPENDIX B – TASK FORCE MEMBERS

893	Timothy G. Ferris, MD, Mphil, MPH	925	Rebecca S. Lipner, PhD
894	(Chair)	926	Vice President of Psychometrics and
895	Associate Professor of Medicine and	927	Research Analysis
896	Pediatrics	928	American Board of Internal Medicine
897	Massachusetts General Hospital/Institute	929	
898	for Health Policy	930	
899	CSAC	931	Jerod Loeb, PhD
900		932	Executive Vice President for Research
901	Andy Amster, MSPH	933	The Joint Commission
902	Director, Integrated Analytics	934	
903	Kaiser Permanente	935	Sean O'Brien, PhD
904		936	Assistant Prof., Dept. of Biostatistics and
905	Nancy Dunton, PhD	937	Bioinformatics
906	Research Professor	938	Duke University Medical Center
907	University of Kansas School of Nursing	939	
908		940	Patrick S. Romano, MD, MPH
909	Steven Findlay, MPH	941	Professor of Medicine and Pediatrics
910	Senior Health Policy Analyst	942	UC Davis School of Medicine
911	Consumers Union	943	
912		944	Amy K. Rosen, PhD
913	David S. P. Hopkins, MS, PhD	945	VA Research Career Scientist
914	Director of Quality Measurement	946	VA Boston Healthcare System
915	Pacific Business Group on Health	947	
916	CSAC	948	Jed Weissberg, MD
917		949	Senior Vice President, Quality and Care
918	Karen Kmetik, PhD	950	Delivery Excellence
919	Vice President for Performance	951	Kaiser Permanente
920	Improvement		
921	American Medical Association-Physician		
922	Consortium for Performance Improvement		
923			
924			
953			
954			

955 APPENDIX C – GLOSSARY

956 **Data element, critical:** Quality performance measures are based on many individual items of
957 information. Testing at the data element level should include those elements that contribute
958 most to the computed measure score (e.g., account for identifying the greatest proportion of the
959 target condition, event, or outcome being measured (numerator); the target population
960 (denominator); population excluded (exclusions); and when applicable, risk factors with largest
961 contribution to variability in outcome.

962

963 **Data element, quality:** A quality data element is a single piece of information that is used in
964 quality measures to describe part of the clinical care process, including both a clinical entity and
965 its context of use (e.g., diagnosis, active) ¹⁴

966

967 **Electronic Health Record (EHR)** (also electronic patient record, electronic medical record, or
968 computerized patient record): As [defined by Healthcare Information Management and Systems
969 Society \(HIMSS\)](#), the Electronic Health Record (EHR) is a secure, real-time, point-of-care,
970 patient-centric information resource for clinicians.

971

972 **EHR measure:** An EHR measure is specified for use with electronic health records; it is
973 composed of data elements from the quality data set (see below), including code lists and
974 measure logic, and can be translated to machine readable specifications.

975

976 **eMeasure:** As [defined by Health Level Seven \(HL7\)](#), an eMeasure is a health quality measure
977 encoded in the Health Quality Measures Format (HQMF) format is referred to as an
978 "eMeasure." The HQMF is a standard for representing a health quality measure as an electronic
979 document. Through standardization of a measure's structure, metadata, definitions, and logic,
980 the HQMF provides for quality measure consistency and unambiguous interpretation.

981

982 **Empirical evidence:** Analyses of data for the measure as specified, unpublished or published

983

984 **Measure Testing:** Empirical analysis to demonstrate the reliability (2b) and validity (2c) of the
985 measure as specified including analysis of issues that pose threats to the validity of conclusions

986 about quality of care such as exclusions (2d), risk adjustment/stratification for outcome and
987 resource use measures (2e), methods to identify differences in performance (2f), and
988 comparability of data sources/methods (2g).

989
990 **Quality Data Set (QDS):** Clinical data necessary to measure quality performance. The QDS
991 framework contains three levels of information: standard elements, quality data elements, and
992 data flow attributes. Standard elements (e.g., diagnosis) represent the atomic unit of data
993 identified by a data element name, a code set, and a code list composed of one or more
994 enumerated values. The quality data element includes the standard element plus quality data
995 type or context (e.g., diagnosis active). Data flow attributes include source (originator), recorder,
996 setting, and health record field. ¹⁴

997
998 **Reliability:** Reliability refers to the repeatability or precision of measurement. Reliability of
999 data elements refers to repeatability and reproducibility of the data elements for the same
1000 population in the same time period. Reliability of the measure score refers to the proportion of
1001 variation in the performance scores due to systematic differences across the measured entities
1002 (signal) in relation to random error or noise.

1003
1004 **Reliability testing:** Empirical analysis of the measure as specified that demonstrate
1005 repeatability and reproducibility of the data elements in the same population in the same time
1006 period and the precision of the computed measure scores. Reliability testing focuses on random
1007 error in measurement and generally involves testing the agreement between repeated
1008 measurements of data elements (often referred to as inter-rater or inter-observer, which also
1009 applies to abstractors and coders); and the amount of error associated with the computed
1010 measure scores.

1011
1012 **Reliability, threats:** Some aspects of the measure specifications or the specific topic of
1013 measurement can affect reliability. Ambiguous measure specifications can result in unreliable
1014 measures. Small case volume or sample size, or rare events can affect the precision (reliability)
1015 of the measure score.

1016

1017 **Untested Measure:** Measure without empirical evidence of both reliability and validity.
1018 Untested measures are only eligible for time-limited endorsement if the conditions for
1019 considering time-limited endorsement are met.

1020

1021 **Validation:** Activity (testing) to determine if a measure has the property of validity. The term
1022 validation is most often used in reference to the data elements.

1023

1024 **Validity:** Validity refers to the correctness of measurement. Validity of data elements *refers to*
1025 *the correctness of the data elements as compared to an authoritative source. Validity of the*
1026 *measure score refers to the correctness of conclusions about quality that can be made based on*
1027 *the measure scores (i.e., a higher score on a quality measure reflects higher quality).*

1028

1029 **Validity testing:** Empirical analysis of the measure as specified that demonstrate that data are
1030 correct and conclusions about quality of care based on the computed measure score are correct.
1031 Validity testing focuses on systematic errors and bias. It involves testing agreement between the
1032 data elements obtained when implementing the measure as specified and data from another
1033 source of known accuracy. Validity of computed measure scores involves testing hypotheses of
1034 relationships between the computed measure scores as specified and other known measures of
1035 quality or conceptually related aspects of quality. A variety of approaches can provide some
1036 evidence for validity. The specific terms and definitions used for validity may vary by
1037 discipline, including face, content, construct, criterion, concurrent, predictive, convergent, or
1038 discriminant validity. Therefore, the proposed conceptual relationship and test should be
1039 described. The hypotheses and statistical tests often are based on various correlations between
1040 measures or differences between groups known to vary in quality.

1041

1042 **Validity, threats to conclusions about quality:** In addition to unreliability, some aspects of
1043 measure specifications and data can affect the validity of conclusions about quality. Potential
1044 threats include patients excluded from measurement; differences in patient mix for outcome
1045 and resource use measures; measure scores generated with multiple data sources/methods; and
1046 systematic missing or “incorrect” data (unintentional or intentional).

1047

1048

NATIONAL QUALITY FORUM

Measure Evaluation Criteria

December 2009

Conditions for Consideration

Four conditions must be met before proposed measures may be considered and evaluated for suitability as voluntary consensus standards:

- A. The measure is in the public domain or an intellectual property agreement is signed.
- B. The measure owner/steward verifies there is an identified responsible entity and process to maintain and update the measure on a schedule that is commensurate with the rate of clinical innovation, but at least every 3 years.
- C. The intended use of the measure includes both public reporting and quality improvement.
- D. The requested measure submission information is complete. Generally, measures should be fully developed and tested so that all the evaluation criteria have been addressed and information needed to evaluate the measure is provided. Measures that have not been tested are only potentially eligible for a time-limited endorsement and in that case, measure owners must verify that testing will be completed within 12 months of endorsement.

Criteria for Evaluation

If all four conditions for consideration are met, candidate measures are evaluated for their suitability based on four sets of standardized criteria: importance to measure and report, scientific acceptability of measure properties, usability, and feasibility. Not all acceptable measures will be strong – or equally strong – among each set of criteria. The assessment of each criterion is a matter of degree; however, all measures must be judged to have met the first criterion, importance to measure and report, in order to be evaluated against the remaining criteria.

1. Importance to measure and report: Extent to which the specific measure focus is important to making significant gains in health care quality (safety, timeliness, effectiveness, efficiency, equity, patient-centeredness) and improving health outcomes for a specific high impact aspect of healthcare where there is variation in or overall poor performance. *Candidate measures must be judged to be important to measure and report in order to be evaluated against the remaining criteria.*

1a. The measure focus addresses:

- a specific national health goal/priority identified by NQF’s National Priorities Partners;
OR
- a demonstrated high impact aspect of healthcare (e.g., affects large numbers, leading cause of morbidity/mortality, high resource use (current and/or future), severity of illness, and patient/societal consequences of poor quality).

1b. Demonstration of quality problems and opportunity for improvement, i.e., data¹ demonstrating considerable variation, or overall poor performance, in the quality of care across providers and/or population groups (disparities in care).

1c. The measure focus is:

¹ Examples of data on opportunity for improvement include, but are not limited to: prior studies, epidemiologic data, measure data from pilot testing or implementation. If data are not available, the measure focus is systematically assessed (e.g., expert panel rating) and judged to be a quality problem.

- an outcome (e.g., morbidity, mortality, function, health-related quality of life) that is relevant to, or associated with, a national health goal/priority, the condition, population, and/or care being addressed²;
OR
- if an intermediate outcome, process, structure, etc., there is **evidence**³ that supports the specific measure focus as follows:
 - o Intermediate outcome – evidence that the measured intermediate outcome (e.g., blood pressure, Hba1c) leads to improved health/avoidance of harm or cost/benefit.
 - o Process – evidence that the measured clinical or administrative process leads to improved health/avoidance of harm and
if the measure focus is on one step in a multi-step care process⁴, it measures the step that has the greatest effect on improving the specified desired outcome(s).
 - o Structure – evidence that the measured structure supports the consistent delivery of effective processes or access that lead to improved health/avoidance of harm or cost/benefit.
 - o Patient experience – evidence that an association exists between the measure of patient experience of health care and the outcomes, values and preferences of individuals/ the public.
 - o Access – evidence that an association exists between access to a health service and the outcomes of, or experience with, care.
 - o Efficiency⁵ – demonstration of an association between the measured resource use and level of performance with respect to one or more of the other five IOM aims of quality.

If not important to measure and report, STOP.

2. Scientific acceptability of the measure properties: Extent to which the measure, as specified, produces consistent (reliable) and credible (valid) results about the quality of care when implemented.

2a. The measure is well defined and precisely specified⁶ so that it can be implemented consistently within and across organizations and allow for comparability. The required data elements are of high quality as defined by NQF's Health Information Technology Expert Panel (HITEP)⁷.

² Generally, rare event outcomes do not provide adequate information for improvement or discrimination; however, “never events” that are compared to zero are appropriate outcomes for public reporting and quality improvement.

³ The strength of the body of evidence for the specific measure focus should be systematically assessed and rated (e.g., USPSTF grading system – [grade definitions](#) and [methods](#)). If the USPSTF grading system was not used, the grading system is explained including how it relates to the USPSTF grades or why it does not. However, evidence is not limited to quantitative studies and the best type of evidence depends upon the question being studied (e.g., randomized controlled trials appropriate for studying drug efficacy are not well suited for complex system changes). When qualitative studies are used, appropriate qualitative research criteria are used to judge the strength of the evidence.

⁴ 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 multi-step process, the step with the greatest effect on the desired outcome should be selected as the focus of measurement. For example, although assessment of immunization status and recommending immunization are necessary steps, they are not sufficient to achieve the desired impact on health status – patients must be vaccinated to achieve immunity. This does not preclude consideration of measures of preventive screening interventions where there is a strong link with desired outcomes (e.g., mammography) or measures for multiple care processes that affect a single outcome.

⁵ Efficiency of care is a measurement construct of cost of care or resource utilization associated with a specified level of quality of care. It is a measure of the relationship of the cost of care associated with a specific level of performance measured with respect to the other five IOM aims of quality. Efficiency might be thought of as a ratio, with quality as the numerator and cost as the denominator. As such, efficiency is directly proportional to quality, and inversely proportional to cost. (NQF's [Measurement Framework: Evaluating Efficiency Across Episodes of Care](#); based on [AQA Principles of Efficiency Measures](#)).

2b. Reliability testing⁸ demonstrates the measure results are repeatable, producing the same results a high proportion of the time when assessed in the same population in the same time period.

2c. Validity testing⁹ demonstrates that the measure reflects the quality of care provided, adequately distinguishing good and poor quality. If face validity is the only validity addressed, it is systematically assessed.

2d. Clinically necessary measure exclusions are identified and must be:

- supported by evidence¹⁰ of sufficient frequency of occurrence so that results are distorted without the exclusion;

AND

- a clinically appropriate exception (e.g., contraindication) to eligibility for the measure focus¹¹;

AND

- precisely defined and specified:

- if there is substantial variability in exclusions across providers, the measure is specified so that exclusions are computable and the effect on the measure is transparent (i.e., impact clearly delineated, such as number of cases excluded, exclusion rates by type of exclusion);
- if patient preference (e.g., informed decision-making) is a basis for exclusion, there must be evidence that it strongly impacts performance on the measure and the measure must be specified so that the information about patient preference and the effect on the measure is transparent¹² (e.g., numerator category computed separately, denominator exclusion category computed separately).

⁶ Measure specifications include the target population (e.g., denominator) to whom the measure applies, identification of those from the target population who achieved the specific measure focus (e.g., numerator), measurement time window, exclusions, risk adjustment, definitions, data elements, data source and instructions, sampling, scoring/computation.

⁷ The HITEP criteria for high quality data include: a) data captured from an authoritative/accurate source; b) data are coded using recognized data standards; c) method of capturing data electronically fits the workflow of the authoritative source; d) data are available in EHRs; and e) data are auditable. NQF. *Health Information Technology Expert Panel Report: Recommended Common Data Types and Prioritized Performance Measures for Electronic Healthcare Information Systems*. Washington, DC: NQF; 2008.

⁸ Examples of reliability testing 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 may address the data items or final measure score.

⁹ Examples of validity testing include, but are not limited to: determining if measure scores adequately distinguish between providers known to have good or poor quality assessed by another valid method; correlation of measure scores with another valid indicator of quality for the specific topic; ability of measure scores to predict scores on some other related valid measure; content validity for multi-item scales/tests. Face validity is a subjective assessment by experts of whether the measure reflects the quality of care (e.g., whether the proportion of patients with BP < 140/90 is a marker of quality). If face validity is the only validity addressed, it is systematically assessed (e.g., ratings by relevant stakeholders) and the measure is judged to represent quality care for the specific topic and that the measure focus is the most important aspect of quality for the specific topic.

¹⁰ Examples of evidence that an exclusion distorts measure results include, but are not limited to: frequency of occurrence, sensitivity analyses with and without the exclusion, and variability of exclusions across providers.

¹¹ Risk factors that influence outcomes should not be specified as exclusions.

¹² Patient preference is not a clinical exception to eligibility and can be influenced by provider interventions.

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

- an evidence-based risk-adjustment strategy (e.g., risk models, risk stratification) is specified and is based on patient clinical factors that influence the measured outcome (but not disparities in care) and are present at start of care^{11,13}

OR

- rationale/data support no risk adjustment.

2f. Data analysis demonstrates that methods for scoring and analysis of the specified measure allow for identification of statistically significant and practically/clinically meaningful¹⁴ differences in performance.

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

2h. 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);

OR

rationale/data justifies why stratification is not necessary or not feasible.

3. Usability: Extent to which intended audiences (e.g., consumers, purchasers, providers, policy makers) can understand the results of the measure and are likely to find them useful for decision making.

3a. Demonstration that information produced by the measure is meaningful, understandable, and useful to the intended audience(s) for both public reporting (e.g., focus group, cognitive testing) and informing quality improvement (e.g., quality improvement initiatives)¹⁵. An important outcome that may not have an identified improvement strategy still can be useful for informing quality improvement by identifying the need for and stimulating new approaches to improvement.

3b. The measure specifications are harmonized¹⁶ with other measures, and are applicable to multiple levels and settings.

3c. Review of existing endorsed measures and measure sets demonstrates that the measure provides a

¹³ 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, gender (e.g., poorer treatment outcomes of African American men with prostate cancer, inequalities in treatment for CVD risk factors between men and women). It is preferable to stratify measures by race and socioeconomic status rather than adjusting out differences.

¹⁴ 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% v. 75%) 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 poor performance may not demonstrate much variability across providers.

¹⁵ Public reporting and quality improvement are not limited to provider-level measures – community and population measures also are relevant for reporting and improvement.

¹⁶ Measure harmonization refers to the standardization of specifications for similar measures on the same topic (e.g., *influenza immunization* of patients in hospitals or nursing homes), or related measures for the same target population (e.g., eye exam and HbA1c for *patients with diabetes*), or definitions applicable to many measures (e.g., age designation for children) so that they are uniform or compatible, unless differences are dictated by the evidence. The dimensions of harmonization can include numerator, denominator, exclusions, and data source and collection instructions. The extent of harmonization depends on the relationship of the measures, the evidence for the specific measure focus, and differences in data sources.

distinctive or additive value to existing NQF-endorsed measures (e.g., provides a more complete picture of quality for a particular condition or aspect of healthcare).

4. Feasibility: Extent to which the required data are readily available, retrievable without undue burden, and can be implemented for performance measurement.

4a. For clinical measures, required data elements are routinely generated concurrent with and as a byproduct of care processes during care delivery.

4b. The required data elements are available in electronic sources. If the required data are not in existing electronic sources, a credible, near-term path to electronic collection by most providers is specified and clinical data elements are specified for transition to the electronic health record.

4c. Exclusions should not require additional data sources beyond what is required for scoring the measure (e.g., numerator and denominator) unless justified as supporting measure validity.

4d. Susceptibility to inaccuracies, errors, or unintended consequences and the ability to audit the data items to detect such problems are identified.

4e. Demonstration that the data collection strategy (e.g., source, timing, frequency, sampling, patient confidentiality¹⁷, etc.) can be implemented (e.g., already in operational use, or testing demonstrates that it is ready to put into operational use).

If a measure meets the above criteria and there are competing measures (either endorsed measures, or other new submissions that also meet the criteria), compare measures on: Scientific acceptability of measure properties, Usability, and Feasibility to determine best-in-class.

5. Demonstration that the measure is superior to competing measures – new submissions and/or endorsed measures (e.g., is a more valid or efficient way to measure).

1055
1056
1057

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.