

NATIONAL QUALITY FORUM

TO: NQF Members

FR: Karen Pace, PhD, RN

RE: Review of *Guidance for Measure Testing and Evaluating Scientific Acceptability of Measure Properties*

DA: June 14, 2010

Background

The National Quality Forum (NQF) relies on four criteria for evaluating the suitability of quality measures for endorsement as voluntary consensus standards. The second criterion, *Scientific Acceptability of Measure Properties*, is an important aspect of the successful use of publicly reported measures to improve performance. Scientific acceptability of measure properties refers to the reliability and validity of measures. In evaluating a measure, both empirical evidence and expert judgment play a role. However, judgment can best be applied when those evaluating a measure have a thorough understanding of the evidence of scientific acceptability that does or does not exist. Evidence that a clearly specified measure produces credible results on performance comes from the basic measurement principles of reliability and validity.

The NQF Task Force on Measure Testing was asked to address the following tasks.

- Identify the type of testing that should be conducted for various types of measures and data sources, and determine whether there are any acceptable alternatives to formal testing.
- Identify the type of testing that should be required prior to endorsement of measures specified for electronic health records (EHRs), both measures originally developed using other data sources besides the EHR and measures developed de novo for EHRs.
- Develop guidance for measure stewards/developers and NQF technical advisors and steering committees on adequate measure testing, interpretation of results, and information about testing that should be provided in the measure submission.
- Make recommendations for potential enhancements to the evaluation criteria.

The Task Force's recommendations are included in the draft document, *Guidance for Measure Testing and Evaluating Scientific Acceptability of Measure Properties*. The draft report is posted on the NQF web site for review and comment only – not voting.

You may post your comments and view the comments of others on the NQF website.

NQF Member comments must be submitted no later than 6:00 PM ET, July 13, 2010; public comments are due 6:00 PM ET, July 6, 2010.

NQF is now using a program that facilitates electronic submission of comments on this draft report. **All comments must be submitted using the online submission process.**

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due by July 13, 2010, 6:00 PM ET; PUBLIC comments due by July 6, 2010 by 6:00 PM ET

Supporting documents related to your comments may be submitted by e-mail to performancemeasures@qualityforum.org with “Measure Testing Report” in the subject line and your contact information in the body of the e-mail.

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

NATIONAL QUALITY FORUM

Guidance for Measure Testing and Evaluating Scientific Acceptability of Measure Properties

Draft Report for Review and Comment

June 14, 2010

NATIONAL QUALITY FORUM

Guidance for Measure Testing and Evaluating Scientific Acceptability of Measure Properties Draft Report June 14, 2010

TABLE OF CONTENTS

INTRODUCTION AND CHARGE	1
Task Force Charge.....	1
BACKGROUND.....	2
Table 1. Measure Evaluation Criteria and Type of Evidence	3
Reliability and Validity	3
Reporting of Measure Scores and Scientific Acceptability.....	5
Measure Testing Issues Identified with Measures Submitted to NQF.....	5
Electronic Health Records and Electronic Measures	7
Summary of Background	7
RECOMMENDATIONS.....	8
I. Recommendations for Empirical Evidence of Reliability and Validity	8
II. Recommendations for the Type of Testing and Results Needed to Demonstrate Scientific Acceptability of Measure Properties	10
Table 2. Evaluation Ratings for Reliability and Validity	13
Table 3. Evaluation of Scientific Acceptability of Measure Properties Based on Reliability and Validity Ratings.....	14
III. Recommendations for Measures Specified for EHRs	15
Table 4. Evaluation of Reliability and Validity of Measures Specified for EHRs	18
IV. Recommendations Related to Untested Measures.....	19
Table 5. Minimum Requirements for Untested Measures under Scientific Acceptability of Measure Properties	19
V. Recommendations for Additional Testing Required for Maintenance of Endorsement	19
VI. Recommendations for Modifications to the NQF Evaluation Criteria.....	20
VII. Recommendations for the Measure Submission.....	23
REFERENCES	26
APPENDIX A – COMMON APPROACHES TO TESTING	28
Table A-1. Common Approaches to Testing Reliability Applied to Quality Measures.....	29
Table A-2. Common Approaches to Testing Validity Applied to Quality Measures	30
Table A-3. Interpretation of Statistical Results.....	32
APPENDIX B – TASK FORCE MEMBERS.....	33
APPENDIX C – GLOSSARY.....	34
APPENDIX D – MEASURE EVALUATION CRITERIA.....	37

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. The second criterion, *Scientific*
4 *Acceptability of Measure Properties*, is an important aspect of the successful use of publicly
5 reported measures to improve performance. Scientific acceptability of measure properties refers
6 to the reliability and validity of measures. The use of measures that are unreliable or invalid
7 undermines confidence in measures of both the providers of healthcare and the consumers of
8 the information. The goal of this document is to provide recommendations on what constitutes
9 scientific acceptability of measures to assist those participants in the measure evaluation
10 process, including steering committees and technical advisory panel members, as well as
11 measure developers. Guidance on scientific acceptability will facilitate a shared understanding
12 of this complex and highly specialized subject.

13
14 In evaluating a measure, both empirical evidence and expert judgment play a role. However,
15 judgment can best be applied when those evaluating a measure have a thorough understanding
16 of the evidence of scientific acceptability that does or does not exist. Evidence that a clearly
17 specified measure produces credible results on performance comes from the basic measurement
18 principles of reliability and validity. Although reliability and validity have always been
19 included in NQF evaluation criteria, the criteria have not included specific guidance on 1) the
20 scope of testing, 2) what tests of reliability and validity could be performed, and 3) how to
21 weigh the results of this testing.

22
23 Task Force Charge

24 The NQF Task Force on Measure Testing was asked to address the following tasks.

- 25 • Identify the type of testing for scientific acceptability that should be conducted for various
26 types of measures and data sources, and determine whether there are any acceptable
27 alternatives to formal testing.
- 28 • Identify the type of testing that should be required prior to endorsement of measures
29 specified for electronic health records (EHRs), both measures originally developed using
30 other data sources besides the EHR and measures developed de novo for EHRs.

- 31 • Develop guidance for measure stewards/ developers and NQF technical advisors and
32 steering committees on adequate measure testing, interpretation of results, and information
33 about testing that should be provided in the measure submission.
- 34 • Make recommendations for potential enhancements to the evaluation criteria.

37 BACKGROUND

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

44
45 The NQF measure evaluation criteria can be viewed as a hierarchy for evaluating measures.

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

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

55 ³

56
57 NQF’s measure [evaluation criteria](#) include a variety of types of evidence as indicated in Table 1.
58 The criterion, *Scientific Acceptability of Measure Properties*, addresses *how* the healthcare quality
59 concept is measured. This criterion includes reliability ([2b](#)) and validity ([2c](#)), as well as precision
60 of specifications ([2a](#)) and potential threats to valid conclusions about quality related to
61 exclusions ([2d](#)), risk adjustment for outcome and resource use measures ([2e](#)), methods to

62 discriminate performance (2f), and comparability of results from different data sources (2g). The
 63 other subcriteria include specifications to detect disparities (2h).

64

65 Table 1. Measure Evaluation Criteria and Type of Evidence

Evaluation Criteria	Sources of Evidence
1.Importance to Measure and Report 1a . High impact	Epidemiologic data Health services research
1.Importance to Measure and Report 1b . Opportunity for improvement	Epidemiologic data Health services research
1.Importance to Measure and Report 1c . Evidence that supports the focus of measurement	Clinical research Health services research
2.Scientific Acceptability of Measure Properties 2a-h . Reliability and validity	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 understanding and 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

66

67

68 **Reliability and Validity**

69 A quality measure is a numeric quantification of the relatively abstract construct of quality of
 70 healthcare, which is measured imperfectly. Reliability refers to the *repeatability and*
 71 *reproducibility of measure scores for the same population in the same time period*. Validity
 72 refers to the *correctness of conclusions about quality* that can be made based on the measure
 73 scores. Over the past four to five decades numerous methods have been devised to test
 74 measures and thus address the measure performance issues inherent to all measurement. These
 75 approaches provide empirical evidence of the properties of reliability and validity. Examples of
 76 approaches to reliability and validity testing can be found in Table A-1 (Appendix A).

77

78 A measure score reflects a theoretical ‘true score’ plus error (or noise): The more error, the less
 79 reliable and valid is the measurement. Random or chance errors affect the reliability or
 80 repeatability of measurement and systematic errors affect the validity or correctness of the

81 conclusions one can make based on the measure score. Threats to reliability include ambiguous
82 measure specifications (including definitions, codes, data collection, and scoring) and statistical
83 issues such as sampling and small case volume. Threats to validity include unreliability and
84 other aspects of the measure specifications and data such as inappropriate exclusions, lack of
85 appropriate risk adjustment or risk factors for outcomes and resource use, specifications for
86 multiple data sources or methods that result in different scores and conclusions about quality,
87 and systematic missing or “incorrect” data. Most importantly, a measure may be invalid
88 because the measurement has not correctly captured the concept of quality it was intended to
89 measure.

90
91 Reliability and validity are not all-or-none properties; rather, measures of reliability and validity
92 produce graduated results. Therefore, results of measure testing always require interpretation.
93 Reliability and validity are not static; they are influenced by the conditions under which the
94 measures are used. A discussion of measurement concepts can be accessed in an online [research](#)
95 [methods knowledge base](#).⁴ [Rubin et al.](#)³ describe reliability and validity testing in quality
96 measure development. Examples of validity testing of healthcare quality measures also are
97 reported in the literature.^{5,6}

98
99 The concept of reliability can be applied to both the computed performance measure score (e.g.,
100 rate, proportion, average) and the individual data elements used in a measure (e.g., diagnosis,
101 medication, admission date, birth date). Reliability is considered necessary, but not sufficient,
102 for achieving a valid measure. A measure could be repeatable, but lead to incorrect conclusions
103 and would therefore be invalid. In addition, if a measure is not repeatable, a valid conclusion
104 about quality would not be possible.

105
106 Evaluation of the scientific acceptability of a measure does not occur in a vacuum. The Task
107 Force was aware of factors within the current environment impacting their deliberations. The
108 recommendations of the Task Force would have implications for both measure developers and
109 healthcare providers. For example, some observers have suggested that existing measure
110 evaluation criteria are too stringent (the perfect being the enemy of the good) while others have
111 suggested that the criteria are not rigorous enough. Some contend that providers use adherence

112 to the measure evaluation criteria as a barrier to making performance information available;
113 others maintain that unless a measure has adequate measure properties it cannot provide useful
114 information. Nonetheless, the consequences of using unreliable or invalid measures can at
115 times be significant for those being measured as well as those who use the information to select
116 a healthcare provider. Resources may be wasted or misdirected; and there is potential for
117 invalid measures to result in misinformation and misdirection of patients or potential
118 unintended harmful consequences. As the stakes around quality measurement are raised, the
119 potential for conflicts among these perspectives increases. The Task Force therefore made a
120 deliberate attempt to make recommendations that balanced the requirement for insuring that
121 NQF endorsed measures would be both sufficiently reliable and valid to make them meaningful
122 and minimize unintended consequences, with requirements for testing that were not so high as
123 to stifle measure development and innovation.

124

125 **Reporting of Measure Scores and Scientific Acceptability**

126 NQF does not determine the specific use or reporting formats of the measures it endorses.
127 Nonetheless, the scientific acceptability of a measure can be related to the context in which the
128 measure is used and the choices made in reporting performance measure scores. For example,
129 Kaplan and colleagues ⁷ demonstrated that the number of categories chosen for performance
130 reporting (e.g., high/medium/low) influences the likelihood of misclassification.
131 Misclassification is, by definition, an invalid reporting of performance. Reporting performance
132 from highest to lowest without information on margin of error and meaningful differences
133 limits and may misrepresent the knowledge to be gained from measures. Further, those
134 choosing to report measures may decide to combine the measures in order to simplify
135 reporting, making the metrics more accessible and usable for consumers. These combined
136 measures also have potential to be misleading. ⁸ Because NQF endorsement does not dictate
137 how the measures are used, the Task Force was not asked to make recommendations on
138 reporting but these issues are highlighted for further discussion and assessment.

139

140 **Measure Testing Issues Identified with Measures Submitted to NQF**

141 The Task Force understood their charge as emerging from several years of NQF experience with
142 measure evaluation. This experience, enumerated below in six points, informed the Task Force's

143 recommendations. First, the NQF portfolio of endorsed measures shows considerable variation
144 in the level of rigor used in measure testing. Measure developers are currently expected to
145 address these requirements in a way that is most appropriate and feasible for the measure and
146 data source involved. Nonetheless, some developers submit limited information on reliability or
147 validity testing perhaps due to a lack of expertise or resources. On the other hand, other
148 measure developers have conducted formal reliability and validity testing and have
149 demonstrated that a proposed measure generates reproducible results and credible conclusions
150 about quality.

151
152 Second, when reliability and validity testing results have been submitted, there has been
153 variability in the scope of testing and the rigor of methods and statistical analysis. For example,
154 reliability of categorical data elements may be assessed only as the percentage of agreement
155 between raters versus using the kappa statistic, which adjusts for chance agreement. In some
156 cases, the testing was conducted with a particular data source, such as the paper medical
157 record, while the measure was specified using a different data source, such as electronic health
158 record.

159
160 Third, there also has been some confusion regarding what is considered testing of scientific
161 acceptability. Terms such as “measure testing,” “pilot testing,” and “field testing” are
162 commonly used in the discipline of measure development and include reliability and validity
163 testing, as well as other aspects of measure development. For example, measure submissions
164 may include descriptive statistics that demonstrate the data are available and can be analyzed to
165 produce scores, but do not specifically address reliability or validity.

166
167 Fourth, some submissions rely on an assumption of reliability and validity. This assumption
168 may be based on prior use of the measure or some aspects of the measure specifications (e.g.,
169 diagnosis codes are relatively well defined and used in accordance with coding rules). In some
170 cases an argument is made that a data source would become more reliable and valid if a quality
171 measure was implemented and publicly reported.

172

173 Fifth, measure developers rarely submit analyses justifying exclusions or demonstrating
174 comparability of different methods of data collection.

175
176 Sixth, steering committees may variably weigh the strengths and weaknesses of the evidence for
177 reliability and validity in their recommendation for endorsement. In summary, while NQF has
178 been raising the bar of expectations and introducing greater rigor and standardization to the
179 evaluation process, the NQF portfolio of endorsed measures still includes varying levels of
180 methodological rigor.

181 182 **Electronic Health Records and Electronic Measures**

183 Development and implementation of electronic health record (EHR) systems hold great promise
184 for the efficient collection of clinical data that can be used for quality measurement. National
185 initiatives call for the adoption of electronic health records that include the capability for quality
186 measurement and NQF has made endorsing quality measures specified for EHRs an important
187 goal. Data stored in EHRs facilitate reporting of quality measures because EHR data 1) are
188 clinically specific, 2) include a large variety of data types including physiologic data, and 3)
189 decrease the burden of the data collection through automated collection and aggregation.

190
191 While the concepts of reliability and validity apply equally to measures derived from EHRs, the
192 electronic health record also presents additional issues related to measure testing. Widespread
193 EHR data are not available for measure development and testing. In addition, the numerous
194 vendors and home grown EHR systems presents the additional challenge of insuring that the
195 selected data fields of interest for any particular measure are comparable among different
196 EHRs. Recommendations regarding testing and evaluation of EHR measures is addressed in
197 Section III.

198 199 **Summary of Background**

- 200 • There are no perfect quality performance measures and there will be some error in all
201 measurement. Performance measurement science is an imperfect science.
- 202 • Measurement principles of reliability and validity apply to quality performance measures
203 regardless of data source.

- 204 • Reliability and validity are not all-or-none properties and involve a matter of degree.
- 205 • Reliability and validity can apply to individual data elements used in a measure, as well as
206 the computed measure score.
- 207 • Reliability is necessary but not sufficient for validity (i.e., reliability is about repeatability of
208 data or scores (even if incorrect); validity is about the “true” or correct values for data
209 elements and conclusions about quality). Measure scores that are mainly determined by
210 random error (noise) are misleading and lead to unwarranted conclusions about quality.
- 211 • NQF is ultimately concerned with endorsing measures that produce scores from which
212 valid (i.e., correct) conclusions about the quality of care can be made.
- 213 • A measure that is not a valid indicator of quality is not useful for making decisions about
214 selecting healthcare providers based on quality or investing time and resources into
215 improvement.
- 216 • There is no one definitive test for reliability or validity because reliability and validity
217 testing results can vary under the conditions of implementation.

218
219

220 RECOMMENDATIONS

221 The recommendations in this report are intended to provide additional guidance and
222 clarification regarding the NQF criteria related to measure testing. They are not intended to
223 provide a detailed primer on methods for measure testing. The recommendations also are not
224 intended as a definitive scoring system for measure evaluation. Evaluation still requires
225 judgment regarding the adequacy of the empirical testing evidence. The recommendations
226 should promote greater consistency in applying the NQF criteria, while maintaining
227 consideration of multi-stakeholder perspectives during the evaluation.

228

229 I. Recommendations for Empirical Evidence of Reliability and Validity

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

235

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

241

242 Implementation and reporting of measures is expected to lead to improvements in
243 documentation, data coding, and data capture and thus reliability and validity. This assumption
244 of improved reliability and validity over time applies to all measures regardless of data source;
245 however, it does not negate the need for reliability and validity to be demonstrated empirically
246 when a measure is being considered for endorsement.

247

248 Recommendations for measures specified for EHRs are addressed in a separate section (Section
249 III) because they are newer and there are several differences from other data sources. For
250 example, the clinician is often the source of data in EHRs and the data are intended for use in
251 care management. However, these distinctions are not absolute and the same issues of
252 demonstrating scientific acceptability apply to EHR measures, as well as measures based on
253 other data sources. Administrative claims data and EHR data are often combined to assess
254 performance and may be viewed as complementary sources of information, each with their own
255 strengths and limitations.

256

257 Although the Task Force was clear about the recommendation for empirical evidence of
258 reliability and validity, it also recognized the practical implications of this assertion for measure
259 developers. The Task Force therefore, further recommended some strategies that could
260 minimize the burden of testing as follows.

- 261 • Evidence for reliability and validity may be accumulated over time and evaluators
262 should remain flexible with regard to the extent of evidence submitted. The scope of
263 testing may be on a relatively small scale for initial endorsement, followed by further
264 analyses to support continued endorsement at the time of review for maintenance of
265 endorsement.

- 266 • Reliability and validity testing may be conducted on a sample of the measured entities.
267 The analytic unit of the particular measure (e.g., physician, hospital, home health
268 agency) determines the sampling strategy for scientific acceptability testing. The sample
269 should represent the variety of entities whose performance will be measured. The
270 sample should include adequate numbers of units of measurement and patients to
271 answer the specific reliability or validity question with the chosen statistical method.
272 Testing requirements for maintenance of endorsement are addressed in Section V.
- 273 • Reliability and validity testing may be conducted for either the data elements used to
274 calculate the measure score or the computed measure score to achieve an acceptable
275 rating for endorsement. Although ideally testing is conducted for both the critical data
276 elements and the computed measure score, only one level of testing would be required
277 for endorsement.
- 278 • Reliability testing of data elements could be bypassed if empirical validity testing of the
279 data elements demonstrates acceptable validity.
- 280 • Prior evidence of reliability or validity of data elements for the data source specified in
281 the measure (e.g., hospital claims) can be used as evidence for those data elements.
- 282 • Because validity testing of measure scores can be quite burdensome, a formal and
283 [systematic testing of face validity](#) of the measure score as an indicator of quality (see
284 Table A-2) could be acceptable for validity of measure score.

285

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

291

292 II. Recommendations for the Type of Testing and Results Needed to Demonstrate Scientific 293 Acceptability of Measure Properties

294 As noted in the first set of recommendations, the Task Force recommended that empirical
295 evidence of reliability and validity are required for measures considered for NQF endorsement.
296 How should participants in the evaluation process assess the evidence provided when

297 measures are submitted? The Task Force chose to provide guidance on measure testing
298 through the development of criteria to rate the reliability and validity of measures being
299 considered for endorsement. This approach requires well-defined descriptions of the rating
300 scheme to reduce ambiguity and miscommunication. While the Task Force has tried to achieve
301 this precision, it recognizes that there will inevitably be some ambiguity and room for
302 interpretation. In addition, the descriptions may require further clarification and/or revision.
303 Finally, the Task Force was not able to fully assess the impact of the proposed rating system on
304 the measure endorsement process. So, this proposed approach to evaluating scientific
305 acceptability of measure properties should be monitored to ensure it achieves the intent of
306 endorsing reliable and valid measures and does not unduly impede endorsement of measures.

307

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

312

313 Table 2 provides the guidance for rating the level of evidence for reliability and validity, which
314 is classified as high, moderate, or low. The ratings depend on both the type of testing conducted
315 and the results of testing meeting acceptable norms.

316

317 The rating scheme is structured around a distinction between testing the data elements used to
318 calculate a measure (e.g., diagnosis, procedure, age) and the computed measure scores (e.g.,
319 rate, proportion, average). Some measures rely on many data elements. Testing at the data
320 element level does not necessarily need to be conducted for every single data element, but
321 should include those elements that are most critical to the computed score.

322

323 Testing at either the level of data elements or the computed measure score with acceptable
324 results is rated moderate and would be acceptable for endorsement. Testing at both levels of
325 data elements and computed measure score with acceptable results is rated high. The low rating
326 represents inadequate testing or inadequate testing results. As noted previously, untested

327 measures would not be rated on reliability and validity and special considerations for untested
328 measures are addressed in a separate section (see Section IV).

329

330 Table 3 presents the Task Force’s recommendation on how the ratings for reliability and validity
331 are used to determine whether a measure adequately meets the criterion of *Scientific*
332 *Acceptability of Measure Properties*. Moderate ratings for both validity and reliability as described
333 in Table 2 would be required to pass this criterion and be acceptable for endorsement.

334

335 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 and measure score</u>:</p> <ul style="list-style-type: none"> • <u>Data element</u> reliability statistics for critical data elements and measure score are within acceptable norms (<u>tested, or reported in the literature</u> for the same data source); OR commonly used data elements with little question of reliability (e.g., gender, age, date of admission); OR <i>may forego data element reliability testing if data element validity demonstrated</i>; AND • <u>Measure score</u> reliability (precision) 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 statistical testing results are within acceptable norms; AND • <u>Measure score</u> validity testing demonstrates a statistically significant result for the hypothesized performance of the measure score; AND <p>Identified threats to validity (lack of risk adjustment/stratification, multiple data sources/methods, systematic missing or “incorrect” data, statistical methods) are empirically assessed and adequately addressed in measure specifications</p>
Moderate	<p>All measure specifications are unambiguous 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</p> <p>Systematic assessment of face validity of <u>measure score</u> as a quality indicator explicitly addressed and found substantial agreement that <i>the scores obtained from the measure as specified will provide an accurate reflection of quality and can be used to distinguish good and poor quality</i></p> <p>AND</p> <p>Identified threats to validity noted above are empirically assessed and adequately addressed in measure specifications</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 of <u>low reliability</u> for <u>either data elements OR measure score</u> – reliability statistics outside of acceptable norms</p>	<p>The measure specifications (numerator, denominator, exclusions, risk factors) <u>do not</u> 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>OR</p> <p>Empirical evidence of <u>low validity</u> for <u>either data elements OR measure score</u></p> <p>OR Inadequate assessment of face validity:</p> <ul style="list-style-type: none"> • Face validity assessed only by the developer or one group; OR • Face validity did not explicitly address the question of whether 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 <p>OR</p> <p>Threats to validity as noted above are considered likely and are NOT empirically assessed</p>

336

337 Table 3. Evaluation of Scientific Acceptability of Measure Properties Based on Reliability and Validity
 338 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 necessary for validity
Moderate	Moderate-High	Yes	Evidence of reliability and validity
	Low	No	Represents inconsistent evidence – reliability is 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.

339
 340 Common approaches to testing reliability and validity for the data elements as well as the
 341 computed measure score that can be applied to quality performance measures are provided in
 342 Appendix A (Tables A-1 and A-2). The type of testing mentioned is those that *could* be
 343 performed – not what is required. Measure developers should select the testing that is
 344 appropriate and feasible for the measure under consideration and that will meet the moderate
 345 rating as described in Table 2. Table A-2 also addresses potential testing and analysis related to
 346 the threats to validity represented by other subcriteria under *Scientific Acceptability of Measure*
 347 *Properties*. Measure developers should identify the potential threats to validity for the specific
 348 measure and conduct analyses to demonstrate adequate control. Information on interpretation
 349 of the common statistical tests used to demonstrate reliability and validity also are provided in
 350 Appendix A (Table A-3); however, those norms provide only general guidelines and testing
 351 results must be interpreted within the unique context of the specific measure.

352
 353 The information on approaches to testing is not meant to provide an exhaustive list of methods.
 354 Other approaches to testing may be appropriate and could be used if the method and rationale
 355 are explained and judged to be appropriate. For example, if agreement on data elements
 356 between two time periods is proposed as a test of reliability, the rationale for expecting stability
 357 (rather than change) over the time period is important to discuss. Stability over time might be
 358 more plausible for the characteristics of patients served than for interventions provided.
 359 However, calculation of measures scores and descriptive statistics, or the fact that a measure has
 360 been in use do not constitute empirical evidence of reliability or validity. Such information may

361 be relevant to the criteria of opportunity for improvement (1b) and usability of the measure (3a);
362 but alone does not provide information about the reliability or validity of the measure scores.

363
364
365

III. Recommendations for Measures Specified for EHRs

366 The EHR holds significant promise for improving the measurement of healthcare quality. The
367 availability of a broad range of reliable and valid data elements for quality measurement
368 without the burden of data collection is widely anticipated. Because clinical data are entered
369 directly into standardized computer readable fields, the EHR will be considered the
370 authoritative source of clinical information and legal record of care. Quality measures based on
371 EHRs require exporting clinical information recorded by healthcare clinicians from discrete
372 computer readable fields; therefore, measurement errors due to manual abstraction, coding by
373 persons other than the originator, or transcription are eliminated. The potential sources of error
374 that pose threats to the reliability and validity of data elements and measure scores for EHR
375 measures include: 1) incorrect measure specifications, including code lists, logic, or computer
376 readable programming language; 2) EHR system structure or programming that does not
377 comply with standards for data fields, coding, or exporting data; 3) difference in use of data
378 fields by different users or entry into the wrong EHR field; and 4) incorrect parsing of data by
379 natural language processing software used to analyze information from text fields. All of these
380 potential errors except those from parsing technology are analogous to sources of error with
381 measures based on other data sources.

382

383 Table 4 provides the guidance for rating the level of evidence for reliability and validity of EHR
384 measures and it is analogous to the ratings in Table 2. Just as for other measures, Table 3
385 indicates how the ratings are used to make a determination if the criterion, *Scientific Acceptability*
386 *of Measure Properties* has been met for EHR measures. Testing approaches for reliability and
387 validity of the EHR measure score are the same as for any measure as noted in Tables A-1 and
388 A-2.

389

390 There are two differences highlighted in Table 4. First, EHR measures must be specified in
391 accordance with the Quality Data Set (QDS).⁹ The QDS will be updated on a regular basis so if

392 a measure needs a quality data element not currently available, there will be a process to
393 consider additional quality data elements.

394
395 Second, data elements for quality measures, which are extracted from EHRs using computer
396 programming, are by virtue of automation repeatable (reliable); therefore, testing at the data
397 element level should focus on validity as discussed below. This would be consistent with the
398 rating system presented for Table 2, where reliability of data items may be bypassed if validity
399 of data items is demonstrated.

400
401 An approach to testing validity of data elements analyzes agreement between data elements
402 and scores obtained with data exported electronically using the EHR measure specifications to
403 those obtained by review and abstraction of the entire EHR, preferably using EHRs that comply
404 with standards. This approach has been reported in the literature ¹⁰⁻¹² and by HealthPartners in
405 a [Commonwealth report](#) ¹³ on performance measures and EHRs. Because EHR databases may
406 not be available for such testing, another approach is to apply the EHR measure to a simulated
407 data set that reflects standards for EHRs and that should return the known values of the data
408 elements and scores in the simulated data set.

409
410 With either approach, when the results obtained for the EHR measure do not match the known
411 values in the simulated data set or the abstracted data, an analysis is conducted to determine
412 the source of error. If the error is related to the measure specifications, including code lists,
413 logic, and computer readable programming language, they would be corrected before
414 submission for endorsement. If the source of error is due to clinical data entry practices and
415 EHR structures unique to specific organizations, the error would not amenable to changes to the
416 EHR measure specifications but may necessitate the need for further evaluation and refinement
417 of the measure.

418
419 The recommended approach for evaluating reliability and validity of data elements for EHR
420 measures takes into account the current environment in which standards for EHRs and EHR
421 measures are under development and widespread adoption is not yet reality. Therefore, testing
422 sites are limited and testing in a sample of EHR systems may not be representative of others.

423 However, this is no different than testing of data elements for measures based on other data
424 sources in a sample of the measured entities. As noted in the background, reliability and
425 validity are not static properties and no one test is definitive.

426
427 Measure testing requirements should not impede the adoption of EHRs and EHR measures, but
428 should be true to the principles of scientific acceptability. EHRs and EHR measures are new and
429 will most likely require some adjustment of local EHR structures and recording practices to
430 meet standards. Therefore, two recommendations regarding implementation are offered.

- 431 • There should be a period of time when measure scores are reported only to providers
432 so that they correct any discrepancies in their systems and practices before public
433 reporting begins.
- 434 • Providers should be encouraged to conduct their own internal reliability studies.

435
436 Previously endorsed measures specified for chart abstraction or administrative claims data, may
437 be appropriate for specification for EHRs. Although these endorsed measures should have
438 already been tested for reliability and validity, the EHR measure specifications require some
439 assessment of similarity to the original specifications, which also is addressed in Table 4. In
440 some cases, the EHR specifications will represent a substantive change to the measure so that an
441 assessment of reliability and validity of the EHR measure is needed.

442

443 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>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>Empirical evidence of reliability of <u>both data elements and measure score</u>:</p> <ul style="list-style-type: none"> • Data element reliability (repeatability) assumed with computer programming – should test data element validity <p>AND</p> <ul style="list-style-type: none"> • <u>Measure score</u> 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 (and computed scores) exported electronically and data elements (and computed scores) 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 and computed measure score; <p>AND</p> <ul style="list-style-type: none"> • <u>Measure score</u> validity testing demonstrates a statistically significant result for the hypothesized performance of the measure score; <p>AND</p> <p>Identified threats to validity (lack of risk adjustment/stratification, multiple data sources/methods, systematic missing or “incorrect” data, statistical methods) empirically assessed and adequately addressed in measure specifications</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 confirms 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</p>
Moderate	<p>The EHR measure specifications use only data elements from 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 (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 for <u>either data elements OR measure score</u> as noted above;</p> <p>AND</p> <p>Identified threats to validity noted above empirically assessed and adequately addressed in measure specifications</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 (QDS data types, code lists, and measure logic) to the endorsed measure specifications confirms that they represent the original measure, which was judged to be a valid indicator of quality</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>The EHR measure specifications <u>do not</u> use only data elements from the QDS;</p> <p>OR</p> <p>Empirical evidence of <u>low reliability</u> for <u>either data elements OR measure score</u> – reliability statistics outside of acceptable norms</p>	<p>The EHR measure specifications (numerator, denominator, exclusions, risk factors) do not 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>OR</p> <p>Empirical evidence of <u>low validity</u> for <u>either data elements OR measure score</u></p> <p>OR</p> <p>Threats to validity as noted above are considered likely and are NOT empirically assessed</p>	<p>The EHR measure specifications <u>do not</u> use only data elements from the QDS;</p> <p>OR</p> <p>Crosswalk of the EHR measure specifications to the endorsed measure specifications identifies they <u>do NOT</u> represent the original measure</p> <p>OR</p> <p>For measures with time-limited status, empirical evidence of low reliability or validity for original time-limited measure</p>

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

446
447
448
449
450
451
452
453
454
455
456
457
458
459
460

IV. Recommendations Related to Untested Measures

Measures without empirical evidence of reliability and validity are considered untested. Untested measures are only eligible for time-limited endorsement if the conditions for considering time-limited endorsement are met.

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

In addition to passing the criterion, *Importance to Measure and Report*, untested measures must demonstrate an adequate foundation for both reliability and validity as follows. Measures that do not meet these minimum requirements are not ready for testing and should not be recommended for time-limited endorsement.

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

Foundation for Reliability	Foundation for Validity
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.	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>

461
462
463
464
465
466
467
468
469
470
471

V. Recommendations for Additional Testing Required for Maintenance of Endorsement

The above guidance on testing and evidence of reliability and validity is for initial endorsement decisions. With the [new system of endorsement cycles](#), currently endorsed measures will be reviewed for maintenance of endorsement every three years along with new measures. Both new and endorsed measures will be required to meet the measure evaluation criteria, including reliability and validity. Recognizing that reliability and validity are not static properties, no one test is definitive, evidence is accumulated over time, and the proposed rating system allows endorsement for measures with limited evidence of reliability and validity, the Task Force agreed that reliability and validity should be evaluated when measures are reviewed for

472 continued endorsement. For measures that are already endorsed, the evidence needed for
473 reliability and validity may vary based on the following considerations:

- 474 • whether data are available from implementation;
- 475 • and the extent of testing and evidence provided at the time of initial endorsement.

476

The Task Force requests comments on what evidence of reliability and validity should be required for review for maintenance of endorsement.

477

478 VI. Recommendations for Modifications to the NQF Evaluation Criteria

479 The suggested changes to the evaluation criteria reflect the recommended guidance. Criterion 2,
480 *Scientific Acceptability of Measure Properties*, is primarily about reliability and validity and threats
481 to reliability and validity. This criterion is simplified by focusing on the concepts of reliability
482 and validity and arranging the subcriteria to reflect their relationship to reliability or validity as
483 follows. The changes to the current criteria are redlined.

484 **Reliability**

485 Precise specifications (2a) including exclusions (2d)

486 Reliability testing (2b) – data elements or measure score

487 **Validity**

488 Validity testing (2c) – data elements or measure score

489 Justification of exclusions (2d) – relates to evidence

490 Risk adjustment (2e)

491 Comparability of data sources/methods (2g)

492 Identification of differences in performance (2f)

493 Disparities (2h)

494

495

496 **2. Scientific acceptability of the measure properties:** Extent to which the measure, as specified,
497 produces consistent (reliable) and credible (valid) results about the quality of care when
498 implemented. **[See footnotes below the criteria]**

499

500 **2a. Reliability**

501 **2a1.** The measure (with exclusions) is well defined and precisely specified ⁶ so that it can be implemented
502 consistently within and across organizations and allow for comparability. ~~The required data elements are~~
503 ~~of high quality as defined by NQF's Health Information Technology Expert Panel (HITEP) EHR measure~~
504 ~~specifications are based on the quality data set (QDS).~~⁷

505

506 | ~~2b2a2.~~ Reliability testing ⁸ demonstrates the measure results are repeatable, producing the same results a
507 | high proportion of the time when assessed in the same population in the same time period.
508 |

509 | **2b. Validity**

510 | ~~2b1.~~ The measure specifications ⁶ are consistent with the evidence presented to support the focus of
511 | measurement under criterion 1c. The measure is specified to capture the most inclusive target population
512 | indicated by the evidence.

514 | ~~2b2.~~ Validity testing ⁹ demonstrates that the measure correctly reflects the quality of care provided,
515 | adequately ~~distinguishing-identifying~~ good and poor quality. ~~If face validity is the only validity~~
516 | ~~addressed, it is systematically assessed.~~

518 | ~~2b3.~~ For outcome measures and other measures (e.g., resource use) when indicated:

- 519 | • an evidence-based risk-adjustment strategy (e.g., risk models, risk stratification) is specified; and is
520 | based on patient clinical factors that influence the measured outcome (but not disparities in care) and
521 | are present at start of care; ^{11,13} and has demonstrated adequate discrimination and calibration

522 | **OR**

- 523 | • rationale/data support no risk adjustment.

525 | ~~2b4.~~ Data analysis of computed measure scores demonstrates ~~that methods for scoring and analysis of~~
526 | ~~the specified measure allow for identification of~~ statistically significant and practically/clinically
527 | meaningful ¹⁴ differences in performance OR overall poor performance. ~~[does this fit with reliability~~
528 | ~~(small case volume, rare event issues), or validity, neither?]~~

530 | ~~2b5.~~ If multiple data sources/methods are ~~allowed~~ specified, there is demonstration they produce
531 | comparable results.

533 | ~~2b6.~~ If disparities in care have been identified, measure specifications, scoring, and analysis allow for
534 | identification of disparities through stratification of results (e.g., by race, ethnicity, socioeconomic status,
535 | gender);

536 | **OR**

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

539 | ~~2d.~~ Clinically necessary-If measure exclusions are ~~identified and must be:~~ not indicated by the clinical
540 | evidence (1c, 2b.1), they are supported by evidence ¹⁰ of sufficient frequency of occurrence so that results
541 | are distorted without the exclusion;

542 | **AND**

- 543 | • ~~a clinically appropriate exception (e.g., contraindication) to eligibility for the measure focus ¹¹;~~

544 | **AND**

- 545 | • ~~precisely defined and specified:~~

- 546 | – ~~if there is substantial variability in exclusions across providers, the measure is specified so that~~
547 | Specifications for scoring include computing exclusions ~~are computable and so that~~ the effect on the
548 | measure is transparent (i.e., impact clearly delineated, such as number of cases excluded, exclusion
549 | rates by type of exclusion);

550 | **AND**

- 551 | – ~~if~~ patient preference (e.g., informed decision-making) is a basis for exclusion, there must be evidence
552 | that it strongly impacts performance on the measure and the measure must be specified so that the
553 | information about patient preference and the effect on the measure is transparent ¹² (e.g., numerator
554 | category computed separately, denominator exclusion category computed separately).

555 |

556 ~~2h. If disparities in care have been identified, measure specifications, scoring, and analysis allow for~~
557 ~~identification of disparities through stratification of results (e.g., by race, ethnicity, socioeconomic status,~~
558 ~~gender);~~
559 ~~OR~~
560 ~~rationale/data justifies why stratification is not necessary or not feasible.~~
561

562

563 Evaluation Criteria Footnotes

564 6 Measure specifications include the target population (e.g., denominator) to whom the measure applies,
565 identification of those from the target population who achieved the specific measure focus (e.g., numerator, target
566 condition, event, outcome), measurement time window, exclusions, risk adjustment/stratification, definitions, ~~data~~
567 ~~elements~~, data source, code lists with descriptors, and instructions, sampling, scoring/computation.

568 7 EHR measure specifications include data type from the QDS, code lists, EHR field, measure logic data flow
569 diagram, original source of the data, recorder, and setting. The HITEP criteria for high quality data include: a) data
570 captured from an authoritative/accurate source; b) data are coded using recognized data standards; c) method of
571 capturing data electronically fits the workflow of the authoritative source; d) data are available in EHRs; and e) data
572 are auditable. NQF. *Health Information Technology Expert Panel Report: Recommended Common Data Types and Prioritized*
573 *Performance Measures for Electronic Healthcare Information Systems.* Washington, DC: NQF; 2008.

574 8 Reliability testing ~~may address~~ applies to both the ~~data item~~ data elements ~~or and final computed~~ measure score.
575 Examples of reliability testing for data elements include, but are not limited to: inter-rater/abstractor or intra-
576 rater/abstractor studies; internal consistency for multi-item scales; test-retest for survey items. Reliability testing of
577 the measure score addresses precision of measurement (e.g., signal-to-noise).

578 9 Validity testing applies to both the data elements and computed measure score. Validity testing of data elements
579 typically analyzes agreement with another authoritative source of the information. Examples of validity testing of the
580 measure score include, but are not limited to: testing hypotheses that the measures scores indicate quality of care,
581 e.g., determining if measure scores adequately distinguish between are different for providers/groups known to have
582 good or poor differences in quality assessed by another valid quality measure or method; correlation of measure
583 scores with another valid indicator of quality for the specific topic; or relationship to conceptually related measures
584 (e.g., scores on process measures to scores on outcome measures). ability of measure scores to predict scores on some
585 other related valid measure; content validity for multi-item scales/tests. Face validity is a subjective assessment by
586 experts of whether the measure reflects the quality of care (e.g., whether the proportion of patients with BP < 140/90
587 is a marker of quality). If face validity is the only validity addressed, it is systematically tested using a modified
588 Delphi or other process that is open to peer review, e.g., RAND Appropriateness Method, ACC/AHA method)
589 assessed (e.g., ratings by relevant stakeholders) and the measure is judged to represent quality care for the specific
590 topic and that the measure focus is the most important aspect of quality for the specific topic.

591 10 Examples of evidence that an exclusion distorts measure results include, but are not limited to: frequency of
592 occurrence, sensitivity analyses with and without the exclusion, and variability of exclusions across providers.

593 11 Risk factors that influence outcomes should not be specified as exclusions.

594 12 Patient preference is not a clinical exception to eligibility and can be influenced by provider interventions.

595 13 Risk models should not obscure disparities in care for populations by including factors that are associated with
596 differences/inequalities in care such as race, socioeconomic status, gender (e.g., poorer treatment outcomes of
597 African American men with prostate cancer, inequalities in treatment for CVD risk factors between men and
598 women). It is preferable to stratify measures by race and socioeconomic status rather than adjusting out differences.

599 14 With large enough sample sizes, small differences that are statistically significant may or may not be practically or
600 clinically meaningful. The substantive question may be, for example, whether a statistically significant difference of
601 one percentage point in the percentage of patients who received smoking cessation counseling (e.g., 74% v. 75%) is
602 clinically meaningful; or whether a statistically significant difference of \$25 in cost for an episode of care (e.g., \$5,000
603 v. \$5,025) is practically meaningful. Measures with overall poor performance may not demonstrate much variability
604 across providers.

605

606

607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656

VII. Recommendations for the Measure Submission

Following are suggested modifications to the information that is currently requested on the [measure submission form](#) based on the recommended guidance and questions and issues that have arisen with the current submission form. The changes to the current submission items are redlined.

Measure Specifications ([Measure evaluation criterion 2a](#))

- 2a.1. Numerator Statement (Brief ~~text narrative~~ description of the numerator - what is being measured about the target population, e.g., target condition, event, or outcome)
- 2a.2. Numerator Time Window (The time period in which cases are eligible for inclusion in the numerator)
- 2a.3. Numerator Details (All information required to collect ~~the data~~ required to calculate the numerator, including definitions and all codes with descriptors, logic, and definitions)
- 2a.4. Denominator Statement (Brief ~~narrative text~~ description of the denominator - target population being measured)
- 2a.5. Target Population Gender
 - Female
 - Male
- 2a.6. Target Population Age Range
- 2a.7. Denominator Time Window (The time period in which cases are eligible for inclusion in the denominator)
- 2a.8. Denominator Details (All information required to collect ~~the data~~ required to calculate the denominator, ~~the target population being measured~~ including definitions and all codes with descriptors, logic, and definitions)
- 2a.9. Denominator Exclusions (Brief ~~text narrative~~ description of exclusions from the target population)
- 2a.10. Denominator Exclusion Details (All information required to collect ~~the data~~ required for exclusions to the denominator, including all definitions and codes with descriptors, logic, and definitions)
- 2a.11. Stratification Details/Variables (All information required to stratify the measure including the stratification variables, all definitions and codes with descriptors, logic, and definitions)
- 2a.12. Risk Adjustment Type
 - No risk adjustment necessary - measure is not an outcome or resource use measure
 - No risk adjustment necessary - rationale and analysis provided in Section 2e
 - stratification/analysis by subgroup - see variables in 2a.11
 - statistical risk model - specifications 2a.14
 - case mix adjustment
 - paired data at patient level
 - risk adjustment devised specifically for this measure/condition
 - risk adjustment method widely or commercially available
 - Other (specify)
- 2a.14. **Specifications for Statistical Risk Model and Adjustment Methodology/Variables** (Name the statistical method (e.g., logistic regression) and list the risk adjustment model variables all definitions and codes with descriptors, and describe conceptual models, statistical models, or other aspects of model or method Development and testing are reported in Section 2e)
- 2a.15. Detailed Risk Model (Please provide a web page URL or attachment. NQF strongly prefers URLs. Attach documents only if they are not available on a web page and keep attached file to 5 MB or less.)
- 2a.18. Type of Score
 - count
 - frequency distribution
 - non-weighted score/ composite/scale
 - rate/proportion
 - ratio

657 weighted score/ composite/scale
658 categorical
659 continuous variable
660 Other (please indicate)
661 **2a.20. Interpretation of Score** (Classifies interpretation of score according to whether better quality is
662 associated with a higher score, a lower score, a score falling within a defined interval, or a passing score)
663 better quality= higher score
664 better quality = lower score
665 better quality = score within a defined interval
666 passing score defines better quality
667 **2a.21. Measure Score Calculation Algorithm** (Describe the calculation of the measure score as a flowchart
668 or series of steps, including identification of denominator, exclusions, identification of numerator,
669 stratification or adjustment, and classification category)
670 **2a.22. Measure Algorithm or Flow Diagram** (Please provide a web page URL or attachment. NQF strongly
671 prefers URLs. Attach documents only if they are not available on a web page and keep attached file to 5 MB
672 or less. ~~Describe the method for discriminating performance_ (E.g., significance testing)~~)
673 **2a.23. Sampling (Survey) Methodology**
674 If measure is based on a sample (or survey), provide instructions for obtaining the sample, conducting the
675 survey, and guidance on minimum sample size (response rate).
676 **2a.24. Data Source** (Check the sources for which the measure is specified and tested)
677 Documentation of original self-assessment Paper medical record/flow-sheet
678 Electronic administrative data/claims Pharmacy data
679 Electronic clinical data Public health data/vital statistics
680 Electronic Health/Medical Record Registry data
681 External audit Special or unique data
682 Lab data Survey: Patient
683 Management data Survey: Provider
684 Organizational policies and procedures
685 **2a.25. Data Source or Collection Instrument** (~~Identify Name~~ the specific data source or data collection
686 instrument, E.g. name of database, clinical registry, collection instrument, etc.)
687 **2a.26. Data Source or Collection Instrument Reference** (Please provide a web page URL or attachment.
688 NQF strongly prefers URLs. Attach documents only if they are not available on a web page and keep
689 attached file to 5 MB or less.)
690 **2a.29. Data Dictionary or Code Table** (Please provide a web page URL or attachment. NQF strongly prefers
691 URLs. Attach documents only if they are not available on a web page and keep attached file to 5 MB or
692 less.)
693 **2a.32. Level of Measurement/Analysis** (Check the level for which the measure is specified and tested)
694 **Clinicians** 704 Regional/network
695 Individual 705 States
696 Group 706 Counties or cities
697 Other 707 Prescription drug plan
698 Facility/agency 708 **Program**
699 Health plan 709 Disease management
700 Integrated delivery system 710 Quality improvement organization (QIO)
701 Multi-site/corporate chain 711 Other
702 **Population** 712 Can be measured at all levels
703 National 713 Other
714 **2a.36. Care Setting** (Check the settings for which the measure is specified and tested; check all that
715 apply.)
716 **Ambulatory Care** Home 723 Assisted living
717 Ambulatory surgery center Hospice 724 Behavioral health/psychiatric unit All settings
718 Office Hospital 725 Dialysis facility Unspecified or "not applicable"
719 Clinic Long term acute care hospital 726 Emergency medical services/ambulance
720 Emergency Department Nursing home (NH) /skilled nursing 727 Group homes
721 facility (SNF) 728 Other
722 Hospital Outpatient Rehabilitation facility

729 **2a.38. Clinical Services** (Healthcare services being measured; check all that apply.)

730 **Behavioral Health**

- 731 Mental health
- 732 Substance use treatment
- 733 Other

734 **Clinicians** (Continued)

- 735 Podiatrist
- 736 Psychologist/LCSW
- 737 PT/OT/Speech
- 738 **Clinicians**
- 739 Audiologist
- 740 Chiropractor
- 741 Dentist/Oral surgeon
- 742 Dietician/Nutritional professional

- 743 Nurses
- 744 Optometrist
- 745 PA/NP/Advanced Practice Nurse
- 746 Pharmacist
- 747 Physicians (MD/DO)
- 748 Respiratory Therapy
- 749 Other
- 750 Dialysis
- 751 Home health
- 752 Hospice/palliative care
- 753 Imaging
- 754 Laboratory
- 755 Other

756

757 **Reliability Testing** ([Measure evaluation criterion 2b](#))

758 **2b.1. Data Sample** (Description of data sample and size)

759 **2b.2. Analytic Methods** (~~Type-Method~~ of reliability testing and rationale, ~~method for testing~~)

760 **2b.3. Testing Results** (Reliability statistics, assessment of adequacy in the context of norms for the test conducted)

761

762

763 **Validity Testing** ([Measure evaluation criterion 2c](#))

764 **2c.1. Data Sample** (Description of data sample and size)

765 **2c.2. Analytic Method** (~~Type-Method~~ of validity testing and rationale, ~~method for testing~~)

766 **2c.3. Testing Results** (Statistical results, assessment of adequacy in the context of norms for the test conducted)

767

768

769 **Measure Exclusions** ([Measure evaluation criterion 2d](#))

770 **2d.1. Summary of Evidence Supporting Exclusion(s)**

771 **2d.2. Citations for Evidence**

772 **2d.3. Data Sample** (Description of data sample and size)

773 **2d.4. Analytic Method** (Type of analysis and rationale)

774 **2d.5. Testing Results** (~~E.g.~~, frequency, variability, sensitivity analyses of impact on measure scores)

775

776

777 **Risk Adjustment Strategy** ([Measure evaluation criterion 2e](#))

778 **2e.1. Data Sample** ~~from Testing or Current Use~~ (Description of data sample and size used for development and validation)

779 **2e.2. Analytic Method** (~~Type of risk adjustment, analysis and rationale~~ Describe methods for development and testing of risk model including selection of risk factors)

780 **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.)

781 **2e.4. If outcome or resource use measure is not risk adjusted, provide rationale**

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

Identification of Meaningful Differences in Performance ([Measure evaluation criterion 2f](#))

2f.1. Data Sample from Testing or Current Use (Description of data sample and size)

2f.2. Methods to Identify Statistically Significant and Practical or Meaningful Differences in Performance (Type of analysis and rationale)

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 meaningfully differences in performance. If no variability, discuss rationale for performance measurement, e.g., benchmark for determining overall poor performance.)

2f.4. Testing Results (Statistical results, assessment of adequacy in the context of norms for the test conducted)

Comparability of Multiple Data Sources/Methods ([Measure evaluation criterion 2g](#))

2g.1. Data Sample (Description of data sample and size)

798 **2g.2. Analytic Method** (Type of analysis and rationale)
799 **2g.3. Testing Results** (Statistical results, assessment of adequacy in the context of norms for the test
800 conducted E.g., correlation statistics, comparison of rankings)

801
802 **Disparities in Care** (Measure evaluation criterion 2h)

803 **2h.1. If measure is stratified to identify disparities, provide stratified results** (Scores by stratified
804 categories/cohorts)

805 **2h.2. If disparities have been reported/identified but measure is not specified to detect disparities,**
806 **provide follow-up plans**

807

808

809 **REFERENCES**

- 810 1. McGlynn EA. Selecting common measures of quality and system performance. *Med Care*.
811 2003;41(1 Suppl):I39-I47.
- 812 2. McGlynn EA, Asch SM. Developing a clinical performance measure. *Am J Prev Med*.
813 1998;14(3 Suppl):14-21.
- 814 3. Rubin HR, Pronovost P, Diette GB. From a process of care to a measure: the development
815 and testing of a quality indicator. *Int J Qual Health Care*. 2001;13(6):489-496.
- 816 4. Trochim WMK. Research methods knowledge base. *Web Center for Social Research Methods*
817 2006; Available at: <http://www.socialresearchmethods.net/kb/index.php>. Last accessed
818 May 2010.
- 819 5. Bhattacharyya T, Freiberg AA, Mehta P et al. Measuring the report card: the validity of pay-
820 for-performance metrics in orthopedic surgery. *Health Aff (Millwood)*. 2009;28(2):526-532.
- 821 6. Schneider EC, Nadel MR, Zaslavsky AM et al. Assessment of the scientific soundness of
822 clinical performance measures: a field test of the National Committee for Quality
823 Assurance's colorectal cancer screening measure. *Arch Intern Med*. 2008;168(8):876-882.
- 824 7. Kaplan SH, Griffith JL, Price LL et al. Improving the reliability of physician performance
825 assessment: identifying the "physician effect" on quality and creating composite measures.
826 *Med Care*. 2009;47(4):378-387.
- 827 8. Reeves D, Campbell SM, Adams J et al. Combining multiple indicators of clinical quality: an
828 evaluation of different analytic approaches. *Med Care*. 2007;45(6):489-496.
- 829 9. National Quality Forum. *Health Information Technology Expert Panel II - Health IT Enablement*
830 *of Quality Measurement*. Washington, DC: NQF; 2009.
- 831 10. Baker DW, Persell SD, Thompson JA et al. Automated review of electronic health records to
832 assess quality of care for outpatients with heart failure. *Ann Intern Med*. 2007;146(4):270-
833 277.
- 834 11. Persell SD, Wright JM, Thompson JA et al. Assessing the validity of national quality
835 measures for coronary artery disease using an electronic health record. *Arch Intern Med*.
836 2006;166(20):2272-2277.
- 837 12. Weiner M, Stump TE, Callahan CM et al. Pursuing integration of performance measures
838 into electronic medical records: beta-adrenergic receptor antagonist medications. *Qual Saf*
839 *Health Care*. 2005;14(2):99-106.
- 840 13. Briggs JB, Kind EA, Awwad S et al. *Performance Measures Using Electronic Health Records: Five*
841 *Case Studies*. New York, NY: The Commonwealth Fund; 2008. Report No.: 1132. Available
842 at www.commonwealthfund.org.
- 843 14. Fitch K, Bernstein SJ, Aguilar MS et al. *The RAND/UCLA Appropriateness Method User's*
844 *Manual*. Santa Monica, CA: RAND Health; 2000. Available at
845 http://www.rand.org/pubs/monograph_reports/MR1269/.

- 846 15. Spertus JA, Eagle KA, Krumholz HM et al. American College of Cardiology and American
847 Heart Association methodology for the selection and creation of performance measures
848 for quantifying the quality of cardiovascular care. *Circulation*. 2005;111(13):1703-1712.
- 849 16. McGinn T, Wyer PC, Newman TB et al. Tips for learners of evidence-based medicine: 3.
850 Measures of observer variability (kappa statistic). *CMAJ*. 2004;171(11):1369-1373.
- 851 17. Viera AJ, Garrett JM. Understanding interobserver agreement: the kappa statistic. *Fam Med*.
852 2005;37(5):360-363.
- 853 18. 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
856
857

858 APPENDIX A – COMMON APPROACHES TO TESTING

859

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

870

The Task Force requests suggestions for references and resources on testing approaches and interpretation of results.

871

872

873 Table A-1. Common Approaches to Testing Reliability Applied to Quality Measures

Reliability Testing – Data elements	
Data Source for Computing Score	Aspect of Reliability/Test
Retrospective chart abstraction (including registry data abstracted retrospectively from medical records)	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
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;
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)	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
Instrument/scale	Internal consistency reliability (Cronbach’s alpha) Analysis of the extent to which item responses obtained at the same time correlate highly with each other
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)
Other data source	Rationale should be provided for method chosen to demonstrate reliability
Reliability Testing—Measure Score	
Data	Aspect of Reliability/Test
Computed measure scores and individual patient-level data for a sample of measured entities Reliability testing of the computed <u>measure scores</u> does not vary by data source or type of measure	Statistical reliability (precision) of sample average as an estimate of the underlying population average 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) or Intraclass Correlation Coefficient (ICC) Generalizability analysis based on generalizability theory on the sources of variation Other methods may be appropriate and rationale for method chosen should be provided

875 Table A-2. Common Approaches to Testing Validity Applied to Quality Measures

Validity Testing – Data elements	
Data Source	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 source of the same data (“gold standard”)</p> <p>Analysis of agreement using appropriate statistical tests (e.g., kappa, ICC) 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>
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 source of the same data (“gold standard”)</p> <p>Analysis of agreement using appropriate statistical tests (e.g., kappa, ICC) with manual abstraction from the full medical record as the “gold standard”</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 source of the same data (“gold standard”)</p> <p>Analysis of agreement using appropriate statistical tests (e.g., kappa, ICC) with “expert” assessor (conducted at approximately the same time)</p>
EHR clinical record information	<p>Validity of data elements extracted from specified fields in EHRs as compared to some criterion source of the same data (“gold standard”)</p> <p>Analysis of agreement using appropriate statistical tests (e.g., kappa, ICC) 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 and the known values for the critical data elements and computed measure score</p>
Survey – single items	<p>Validity of data elements from survey as compared to some criterion source of the same data (“gold standard”)</p> <p>Analysis of agreement using appropriate statistical tests (e.g., kappa, ICC) with some other source of the same information considered to be valid (e.g., medical record, vital statistics)</p>
Instrument/scale	<p>Validity of the content of the items in an instrument or scale</p> <p>Systematic assessment by subject matter experts that the content of the instrument/scale is representative of the domain being measured</p> <p>Confirmatory factor analysis</p>
Other	Rationale should be provided for method chosen to demonstrate reliability

876

877

878

879 Table A-2 Cont. Common Approaches to Testing Validity Applied to Quality Measures

Validity Testing—Measure Score	
Data	Aspect of Validity/Test
<p>Computed measure scores for a sample of measured entities and other data as necessary for the chosen validity study</p> <p>Validity testing of the computed <u>measure scores</u> does not vary by data source or type of measure</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 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 modified Delphi or other process that is open to peer review, e.g., RAND Appropriateness Method¹⁴, ACC/AHA method)¹⁵</p> <p>Studies to check the performance of the computed measure score against some criterion. These studies are based on a prediction about how the operationalization of the measure will <i>perform</i> based on the theory of the construct.</p> <p><i>Predictive</i> - assess the ability of the measure score to predict something it should theoretically be able to predict</p> <p><i>Concurrent</i> - assess the ability of the measure score to distinguish between groups that it should theoretically be able to distinguish between</p> <p><i>Convergent</i> - 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><i>Discriminative</i> - examine the degree to which the measure score is not similar to (diverges from) other measures that it theoretically should be not be similar to</p>
Threats to Validity—Measure Score	
Data Needed	Aspect of Validity/Test
<p>Patient level data divided into development and validation samples</p>	<p>Threat that differences in measure scores are due to differences in severity of patients served rather than differences in quality</p> <p>For outcome and resource use measures, empirical evidence for the adequacy of controlling for patient factors (analysis of risk factors, discrimination and calibration of risk models);</p> <p>OR evidence that risk adjustment/ stratification is not necessary for fair comparisons (patient outcomes do not vary by patient characteristics)</p>
<p>Computed measure scores for each specified data source</p>	<p>Threat of bias from differences in data source and/or differences in data collection practices</p> <p>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)</p>
<p>Patient level data</p>	<p>Threat of bias from missing or “incorrect” data</p> <p>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)</p>
<p>Patient level data and computed measure scores</p>	<p>Threat of misclassification or inability to distinguish comparative quality from measure scores related to statistical methods</p> <p>Analysis of computed measure scores demonstrates the ability to identify statistically significant and practically or clinically meaningful differences in performance <u>or overall poor performance</u></p>

880

881

882 Table A-3. Interpretation of Statistical Results

Test	Interpretation												
<p>Kappa ^{16, 17} Measure of agreement between two raters that adjusts for chance agreements for categorical data (nominal, ordinal)</p>	<p>Kappa values range between 0 and 1. 0 and are interpreted as degree of agreement beyond chance</p> <table border="0"> <tr> <td>0</td> <td>No better than chance</td> </tr> <tr> <td>0.01-0.20</td> <td>Slight</td> </tr> <tr> <td>0.21-0.40</td> <td>Fair</td> </tr> <tr> <td>0.41-0.60</td> <td>Moderate</td> </tr> <tr> <td>0.61-0.80</td> <td>Substantial</td> </tr> <tr> <td>0.81-1.0</td> <td>Almost perfect</td> </tr> </table>	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
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												
<p>ICC Alternative measure of agreement when more than two raters or quantitative data (interval, ratio)</p>	<p>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</p>												
<p>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)</p>	<p>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 ¹⁸ demonstrated that value of F should be 10 or greater</p>												
<p>Cronbach’s alpha Measure of the average correlation of the items comprising a scale or subscale</p>	<p>A widely-accepted cut-off is .70 or higher for a set of items to be considered a scale, but 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.</p>												
<p>Pearson Correlation Measure of the degree of association (not agreement) between two quantitative variables</p>	<p>Interpretation depends on significance, size, and context. 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).</p>												

883

APPENDIX B – TASK FORCE MEMBERS

884	Timothy G. Ferris, MD, Mphil, MPH	916	Rebecca S. Lipner, PhD
885	(Chair)	917	Vice President of Psychometrics and
886	Associate Professor of Medicine and	918	Research Analysis
887	Pediatrics	919	American Board of Internal Medicine
888	Massachusetts General Hospital/Institute	920	
889	for Health Policy	921	
890	CSAC	922	Jerod Loeb, PhD
891		923	Executive Vice President for Research
892	Andy Amster, MSPH	924	The Joint Commission
893	Director, Integrated Analytics	925	
894	Kaiser Permanente	926	Sean O'Brien, PhD
895		927	Assistant Prof., Dept. of Biostatistics and
896	Nancy Dunton, PhD	928	Bioinformatics
897	Research Professor	929	Duke University Medical Center
898	University of Kansas School of Nursing	930	
899		931	Patrick Romano, MD, MPH
900	Steven Findlay, MPH	932	Professor of Medicine and Pediatrics
901	Senior Health Policy Analyst	933	UC Davis School of Medicine
902	Consumers Union	934	
903		935	Amy K. Rosen, PhD
904	David S. P. Hopkins, MS, PhD	936	VA Research Career Scientist
905	Director of Quality Measurement	937	VA Boston Healthcare System
906	Pacific Business Group on Health	938	
907	CSAC	939	Jed Weissberg, MD
908		940	Senior Vice President, Quality and Care
909	Karen Kmetik, PhD	941	Delivery Excellence
910	Vice President for Performance	942	Kaiser Permanente
911	Improvement		
912	American Medical Association-Physician		
913	Consortium for Performance Improvement		
914			
915			
944			
945			

946 APPENDIX C – GLOSSARY

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

953
954 **Data element, quality:** A quality data element is a single piece of information that is used in
955 quality measures to describe part of the clinical care process, including both a clinical entity and
956 its context of use (e.g., diagnosis, active) ⁹

957
958 **Electronic Health Record (EHR):** (also electronic patient record, electronic medical record, or
959 computerized patient record) As [defined by Healthcare Information Management and Systems
960 Society \(HIMSS\)](#), the Electronic Health Record (EHR) is a secure, real-time, point-of-care,
961 patient-centric information resource for clinicians.

962
963 **EHR measure:** An EHR measure is specified for use with electronic health records and
964 composed of data elements from the quality data set (see below), measure logic, and machine
965 readable specifications.

966
967 **eMeasure:** As [defined by Health Level Seven \(HL7\)](#), an eMeasure is a health quality measure
968 encoded in the Health Quality Measures Format (HQMF) format is referred to as an
969 "eMeasure." The HQMF is a standard for representing a health quality measure as an electronic
970 document. Through standardization of a measure's structure, metadata, definitions, and logic,
971 the HQMF provides for quality measure consistency and unambiguous interpretation.

972
973 **Empirical evidence:** Analyses of data for the measure as specified, unpublished or published

974
975 **Measure Testing:** Empirical analysis to demonstrate the reliability (2b) and validity (2c) of the
976 measure as specified including analysis of issues that pose threats to the validity of conclusions

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

980
981 **Quality Data Set (QDS):** Clinical data necessary to measure quality performance. The QDS
982 framework contains three levels of information: standard elements, quality data elements, and
983 data flow attributes. Standard elements (e.g., diagnosis) represent the atomic unit of data
984 identified by a data element name, a code set, and a code list composed of one or more
985 enumerated values. The quality data element includes the standard element plus quality data
986 type or context (e.g., diagnosis active). Data flow attributes include source (originator), recorder,
987 setting, and health record field. ⁹

988
989 **Reliability:** Measure results are repeatable, producing the same results a high proportion of the
990 time when assessed in the same population in the same time period

991
992 **Reliability testing:** Empirical analysis of the measure as specified including the data elements
993 and/or the computed measure scores that demonstrate repeatability and reproducibility of the
994 data elements and computed measure scores in the same population in the same time period.
995 Reliability testing focuses on random error in measurement and generally involves testing the
996 agreement between repeated measurements of data elements (often referred to as inter-rater or
997 inter-observer, which also applies to abstractors and coders); and the amount of error associated
998 with the computed measure scores.

999
1000 **Reliability, threats:** Some aspects of the measure specifications or the specific topic of
1001 measurement can affect reliability. Ambiguous measure specifications can result in unreliable
1002 measures. Statistical issues, such as small case volume or sample size or rare events can affect
1003 the precision (reliability) of the measure score.

1004
1005 **Untested Measure:** Measure without empirical evidence of reliability and validity. Untested
1006 measures are only eligible for time-limited endorsement if the conditions for considering time-
1007 limited endorsement are met.

1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037

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

Validity: Measure results correctly reflect the quality of care provided, adequately distinguishing good and poor quality. A measure cannot be valid without being reliable.

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

Validity, threats to of conclusions about quality: In addition to unreliability, some aspects of measure specifications and data can affect the validity of conclusions about quality. Potential threats include patients excluded from measurement; differences in patient mix for outcome and resource use measures; measure scores generated with multiple data sources/methods; specifications that obscure disparities in care; systematic missing or “incorrect” data (unintentional or intentional); statistical methods used to produce estimates or identify differences in performance.

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 24 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).

1044
1045
1046

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.