

# NATIONAL QUALITY FORUM

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

## BOARD ACTION

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

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

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

## BACKGROUND

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

## Process

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

## Overview

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

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

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

Following are the key features of the guidance.

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

#### Comments Received

The key issues raised in the comments included the following.

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

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

## GUIDANCE FOR MEASURE TESTING AND SCIENTIFIC ACCEPTABILITY OF MEASURE PROPERTIES

Following are the key features of the guidance.

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

### Comments Received

The key issues raised in the comments included the following.

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

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

# NATIONAL QUALITY FORUM

## **Guidance for Evaluating the Evidence Related to the Focus of Quality Measurement And Importance to Measure and Report**

September 13, 2010

# NATIONAL QUALITY FORUM

## Guidance for Evaluating the Evidence Related to the Focus of Quality Measurement

### CONTENTS

8	OVERVIEW AND PURPOSE.....	2
9	BACKGROUND.....	3
10	Evidence Issues Identified with Measures Submitted to NQF.....	4
11	The Changing Environment.....	5
12	Clinical Practice Guidelines.....	6
13	Evidence Grading Systems.....	6
14	RECOMMENDATIONS.....	9
15	Principles.....	9
16	I. Recommendations on Sources of Evidence and Evidence Grading for the Present and the	
17	Future.....	10
18	II. Recommendations for the Evidence Needed to Justify the Focus of a Quality Measure .....	12
19	III. Recommendations for Evaluating Criterion 1c – Quantity, Quality, Consistency of Body of	
20	Evidence .....	15
21	Table 4. Evaluation of Quantity, Quality, and Consistency of Body of Evidence for Criterion	
22	1c – evidence for the measure focus.....	19
23	Table 5. Evaluation of Subcriterion 1c based on the quantity, quality and consistency of the	
24	body of evidence .....	20
25	IV. Recommendations for Selecting the Focus for Measure Development .....	20
26	V. Recommendations for Evaluating Importance to Measure and Report and the Other	
27	Subcriteria .....	21
28	Consequences of Measurement.....	24
29	VI. Recommendations for Modifications to the NQF Evaluation Criteria.....	25
30	VII. Recommendations for Modifications to the Measure Submission.....	27
31	VIII. Recommendations for Evidence Required for Practices Considered for NQF	
32	Endorsement.....	31
33	Table 9. Evidence to Support a Practice.....	31
34	REFERENCES.....	31
35	APPENDIX A – EVALUATION CRITERIA.....	34
36	APPENDIX A – EVALUATION CRITERIA.....	34
37	Current Measure Evaluation Criteria.....	34
38	Current Evaluation Criteria for Practices .....	39
39	APPENDIX B - TASK FORCE MEMBERS.....	40
40	APPENDIX C - US PREVENTIVE SERVICES TASK FORCE SYSTEM FOR GRADING	
41	EVIDENCE AND RECOMMENDATIONS .....	41

45 OVERVIEW AND PURPOSE

46 Steering committees have diverse backgrounds and expertise and could benefit from more  
47 guidance and support to consistently apply NQF measure evaluation criteria. Both evidence  
48 and expert judgment play a role in evaluating measures against criteria. However, judgment  
49 can best be applied when Steering Committees have a thorough understanding of the evidence  
50 that does or does not exist. Evidence comes in many different forms (e.g., peer reviewed  
51 publications; practice guidelines from authoritative sources; expert assessments); there are often  
52 inconsistencies and gaps; and it can be difficult to interpret and reach conclusions. In  
53 October 2009, the Board directed that NQF should take steps to strengthen its processes to  
54 evaluate the synthesis and scoring of evidence and to present this information in ways that will  
55 be best understood and useful to Steering Committees.

56  
57 NQF's [evaluation criteria](#) require a variety of evidence as noted in the following table. Of these  
58 criteria, some of the most rigorous evidence is required to justify what is being measured (1c)  
59 and that is the primary focus of this report – *the evidence required to justify the measure focus*  
60 (i.e., the specific process, structure, outcome, etc. that is being measured). Another task force  
61 and subsequent report will address measure testing and the criterion of *Scientific Acceptability of*  
62 *Measure Properties*.

63  
64 Evidence refers to the information used to determine or demonstrate the truth of a hypothesis.  
65 The highest quality evidence available should be used to support the focus of quality  
66 performance measures. Evidence is not limited to quantitative studies and the best type of  
67 evidence depends upon the question being studied (e.g., randomized controlled trials  
68 appropriate for studying drug efficacy are not well suited for complex system changes). A body  
69 of evidence includes all the evidence for a topic, which is systematically identified, based on  
70 pre-established criteria for relevance and quality of evidence.

71  
72 NQF endorses measures that are intended for use in public reporting as well as quality  
73 improvement with the goal of improving the quality of healthcare. The evidence that supports  
74 the focus for a quality measure is addressed under the must-pass criterion, *Importance to*  
75 *Measure and Report* because if the measure focus is not supported by evidence that it can

76 facilitate gains in quality and health, then the use of limited resources for measuring and  
 77 reporting on it would be questionable. For most healthcare quality measures, the evidence will  
 78 be that of clinical effectiveness and the link to desired outcomes.

79

80 Table 1. Measure Evaluation Criteria and Type of Evidence

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

81

82 **Task Force Charge**

83 The task force was asked to address the following tasks.

- 84 • Identify the type of evidence needed to justify the focus of a quality measure (1c) (i.e.,  
 85 what is being measured).
- 86 • Identify the evidence needed to demonstrate high impact (1a) and opportunity for  
 87 improvement (1b).
- 88 • Develop guidance on how technical advisors and steering committees use the evidence  
 89 provided to evaluate submitted measures for possible endorsement.
- 90 • Make recommendations for potential enhancements to the evaluation criteria.

91

92

93 **BACKGROUND**

94 Ideally, quality performance measures are based on high quality evidence regarding the types  
 95 of interventions and services that will achieve desired outcomes and reflect high quality care.

96 However, much of healthcare has not been subjected to research studies, much less with  
97 randomized controlled trials or comparative effectiveness studies. Lohr observed that “Perhaps  
98 no more than half, or even one-third, of services are supported by compelling evidence that  
99 benefits outweigh harms <sup>1</sup>.” For example, Tricoci, et al. <sup>2</sup> reviewed recommendations in  
100 American College of Cardiology/American Heart Association guidelines and found that only  
101 314 of 2711 recommendations were classified as A-level evidence based on multiple  
102 randomized trials with large numbers of patients. Many quality performance measures are  
103 based on clinical practice guidelines, however not all guideline recommendations are  
104 appropriate for performance measure development, which depends on the strength of the  
105 evidence and relationship to meaningful outcomes <sup>3</sup>.

106  
107 Some aspects of healthcare (e.g., system change) may be more difficult to study with  
108 quantitative methods, particularly with randomized controlled trials. Some clinical process  
109 steps (i.e., assessing health status, diagnosing clinical conditions, recommending treatment,  
110 teaching and counseling about conditions/treatment) may be unlikely to be subjected to  
111 research. Even when research has been conducted, the body of evidence may not have been  
112 systematically assessed and graded (e.g., care coordination, medication management). Lohr <sup>1</sup>  
113 noted that absence of evidence about benefit is not the same as evidence of no benefit. Even  
114 when available, evidence is rarely definitive. However, the level of confidence in a  
115 recommendation (or measure) depends on the underlying research and synthesis of that  
116 research.

#### 117 118 **Evidence Issues Identified with Measures Submitted to NQF**

119 The NQF evaluation criteria ([1c](#), Footnotes [3](#) & [4](#)) and submission questions may not provide  
120 enough direction to reviewers or measure developers. Measure submissions often have  
121 insufficient information on the strength of the evidence or strength of a guideline  
122 recommendation. Measures have been submitted with no evidence; no systematic grading or  
123 incorrect grading of the evidence or guideline recommendation; use of a different grading  
124 system than the recommended USPSTF system with no explanation; or low quality evidence. In  
125 some cases, a grade might be assigned without using the associated methods to assess the body  
126 of evidence. Some submitted measures are focused on process steps far removed from the



127 desired outcome, even when there is evidence for a particular intervention or intermediate  
128 outcome that is more directly linked to the desired outcome (e.g., measures to assess  
129 immunization status rather than measures of administering the vaccine). Some measure  
130 submitters question whether the suggested USPSTF evidence grading system is only applicable  
131 to preventive services.

132  
133 NQF consensus projects were not intended to undertake systematic evidence reviews for the  
134 variety of measures that are submitted for consideration, nor is this feasible. Such detailed  
135 evidence reviews have also not generally been viewed by developers as an integral part of the  
136 measure development process. However, the responsibility for basing quality performance  
137 measures on appropriate evidence does ultimately lie with measure developers. Measure  
138 developers who do not have the expertise and resources to systematically assess the strength of  
139 a body of evidence sometimes rely on other sources of evidence reviews and grading, such as  
140 found in clinical practice guidelines or published systematic reviews. However, NQF wishes to  
141 clearly signal, through this document and the measure submission form itself, that measure  
142 developers are responsible for identifying, summarizing, and reporting the evidence that exists  
143 to support the focus of measures submitted to NQF for potential endorsement.

## 144 145 **The Changing Environment**

146 As guidelines and quality metrics are increasingly used not only for internal quality  
147 improvement but also for public reporting, the necessity for a strong evidence base has become  
148 more urgent and compelling. This need is further substantiated by the development of  
149 reimbursement programs that utilize such publicly reported metrics. Although public reporting  
150 and pay for performance have the potential to inform consumers, focus quality improvement  
151 activities, and reward high performance; there are potential unintended negative consequences  
152 if measures do not meet all the aspects of the importance criterion. Potential negative  
153 consequences include confusion about the importance of particular care processes to quality,  
154 the unnecessary resources to measure elements of care that may not impact quality, and  
155 diversion of scarce resources to marginally effective activities. To achieve the intended positive  
156 effects of quality measurement and minimize the unintended potential negative consequences,  
157 measures should be based on the best evidence for the focus of measurement and also should

158 conform to the highest measurement science principles. Recognizing the high stakes of  
159 performance measurement in an increasingly transparent environment, some measure  
160 developers have enhanced their requirements for the evidence base for performance measure  
161 development <sup>4</sup>.

162

### 163 **Clinical Practice Guidelines**

164 Although they are not the only evidence base for performance measures, many measure  
165 developers rely on clinical practice guidelines to support the focus of measurement <sup>3,4</sup>. There  
166 has been a proliferation of such guidelines, some overlapping or even contradictory. There also  
167 is substantial variability in the methodological rigor of review and grading of the evidence and  
168 recommendations. In 2000, Grilli <sup>5</sup> and colleagues reported that of 431 specialty society  
169 guidelines reviewed, 82% did not apply explicit criteria to grade the scientific evidence used as  
170 a basis for recommendations, 87% did not report whether a systematic literature search was  
171 conducted, and 67% did not describe the professional involved. Some tools to assess clinical  
172 practice guidelines <sup>6-8</sup> are available and developing trustworthy guidelines is also the subject of  
173 a current IOM study.

174

175 At the January 11, 2010 IOM meeting on developing trustworthy guidelines, Vivian Coates  
176 [presented](#) the following information about the [National Guidelines Clearinghouse](#) (NGC):

- 177 • Currently, NGC contains more than 2500 guidelines from more than 200 developers.
- 178 • Most of the developers whose guidelines are represented in NGC (158 of 204; 77%) use  
179 some sort of rating scheme to grade the underlying evidence and/or strength of the  
180 recommendations. Of these:
  - 181 ○ Ten developers report using GRADE or modified GRADE.
  - 182 ○ Six report using the USPSTF approach, either as is, or modified.
  - 183 ○ The great majority (142 developers) does not identify the origin of their rating  
184 schemes, and appear to be using schemes unique to their organizations.

185

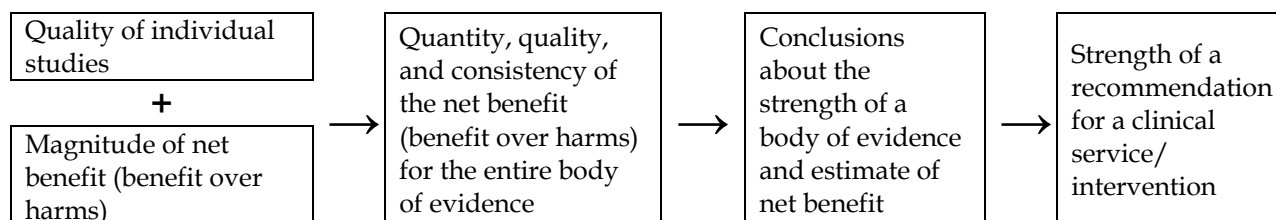
### 186 **Evidence Grading Systems**

187 A variety of evidence grading systems currently are in use to achieve this enhanced degree of  
188 evidence review and assessment. These systems generally include methods for selection and

189 review of the evidence, and rules or hierarchies related to grading the quality of evidence and  
190 the strength of a recommendation. These evidence grading systems are applicable to guidelines  
191 as well as other sources of evidence for performance measures.

192  
193 There are commonalities among the various evidence grading systems. In general, the quality  
194 and strength of the overall body of evidence is a function of the *quantity* and *quality* of  
195 individual studies and the *consistency* among studies regarding judgments of net benefit (the  
196 balance of benefits and harms). *Quality* of individual studies includes study design, sample size  
197 and statistical power considerations, flaws such as selection bias, directness of the evidence  
198 linking an intervention to health outcomes, and generalizability of findings. Of particular  
199 interest for quality measures is how well the measure matches the population and intervention  
200 in the evidence (e.g., cited studies). The general approach to determining the strength of  
201 evidence and a recommendation for a particular intervention or service is depicted in Figure 1.

202  
203 Figure 1. Approach to Determining Quality of Evidence and Strength of Recommendation



204  
205 Differences in terminology and grading scales may inhibit understanding about the strength of  
206 evidence. Differences can range from a rather minor but understandable difference in  
207 terminology (e.g., strength, quality, or level of evidence) to pronounced differences in the  
208 assignment of grades (e.g., a grade of A could indicate evidence based on consensus of opinion  
209 in one system to evidence based on meta-analyses of randomized controlled trials in another  
210 system). An international initiative to standardize grading evidence and recommendations,  
211 [GRADE](#)<sup>9-15</sup>, is now supported by many [organizations](#) including the Cochrane Collaboration.  
212 The Agency for Healthcare Research and Quality (AHRQ) supports two evidence grading  
213 systems: one used by the US Preventive Services Task Force (USPSTF)<sup>16,17</sup> and one used by the  
214 Evidence-Based Practice Centers<sup>18</sup> (consistent with GRADE). Table 2 provides examples of  
215 terminology used by four evidence grading systems. It is important to note that grading  
216 systems are tied to specific methods for reviewing and assessing the quality of evidence.

217

218 Table 2. Examples of Terminology in Selected Grading Scales

	<u>USPSTF</u>	<u>GRADE</u>	<u>AHRQ Evidence-Based Practice Centers</u>	<u>ACC/AHA</u>
<b>Evidence</b>	Certainty of Net Benefit: <ul style="list-style-type: none"> <li>• High</li> <li>• Moderate</li> <li>• Low</li> </ul> Magnitude of Net Benefit: <ul style="list-style-type: none"> <li>• Substantial</li> <li>• Moderate</li> <li>• Small</li> <li>• Zero/Negative</li> </ul>	Quality of Evidence: (confidence in estimate of effect to support recommendation) <ul style="list-style-type: none"> <li>• High</li> <li>• Moderate</li> <li>• Low</li> <li>• Very Low</li> </ul>	Strength of Evidence: (confidence that estimate of effect is correct) <ul style="list-style-type: none"> <li>• High</li> <li>• Moderate</li> <li>• Low</li> <li>• Insufficient</li> </ul>	Estimate of certainty of treatment effect <ul style="list-style-type: none"> <li>• A: multiple pop, RCT, meta-analysis</li> <li>• B: limited pop, single RCT or non-RCT</li> <li>• C: very limited pop, consensus expert opinion, case studies</li> </ul> Size of treatment effect <ul style="list-style-type: none"> <li>• Class I: Benefit &gt;&gt;&gt;Risk</li> <li>• Class IIa: Benefit &gt;&gt;Risk</li> <li>• Class IIb: Benefit &gt; or = Risk</li> <li>• Class III: Risk &gt; or = Benefit</li> </ul>
<b>Recommendation</b>	Grade of Recommendation: Certainty/Magnitude <ul style="list-style-type: none"> <li>• <b>A - Recommend:</b> High/Substantial</li> <li>• <b>B - Recommend:</b> High/Moderate; Moderate/Substantial; Moderate/Moderate</li> <li>• <b>C - Recommend against routine use:</b> High or Mod/Small</li> <li>• <b>D - Recommend against:</b> High or Mod/Zero-Neg</li> <li>• <b>I-Insufficient evidence:</b> Low/any magnitude</li> </ul>	Strength of Recommendation: <ul style="list-style-type: none"> <li>• Strong</li> <li>• Weak</li> </ul>	Does not make recommendation	<ul style="list-style-type: none"> <li>• Should be performed: Class 1-A, B, C</li> <li>• Reasonable to perform: Class IIa-A,B,C</li> <li>• May be considered: Class IIb-A,B,C</li> <li>• Not helpful/may be harmful: Class III-A,B,C</li> </ul>

219

220

221 Systematic reviews and meta-analyses are used to assess a body of evidence. PRISMA  
 222 (Preferred Reporting Items for Systematic Reviews and Meta-Analyses) focuses on the  
 223 transparent and full reporting of such reviews<sup>19</sup>. The Institute of Medicine (IOM) has two  
 224 consensus projects underway that relate to grading the quality of evidence for clinical  
 225 interventions: [Standards for Developing Trustworthy Clinical Practice Guidelines](#) and

226 [Standards for Systematic Reviews of Clinical Effectiveness Research](#); however, reports will not  
227 be ready until early 2011.

228

229

## 230 RECOMMENDATIONS

231 The Task force identified some definitions and principles that guided its discussion and the  
232 recommendations that follow.

233

234 **Evidence** refers to the information used to determine or demonstrate the truth of a hypothesis.  
235 The highest quality evidence available should be used to support the focus of quality  
236 performance measures. Evidence is not limited to quantitative studies and the best type of  
237 evidence depends upon the question being studied (e.g., randomized controlled trials  
238 appropriate for studying drug efficacy are not well suited for complex system changes).

239

240 A **body of evidence** includes all the evidence for a topic, which is systematically identified,  
241 based on pre-established criteria for relevance and quality of evidence.

242

### 243 Principles

244 **Transparency is a primary goal.** All stakeholders need to have a clear understanding of the  
245 evidence supporting a performance measure in order to make informed decisions about the  
246 importance of measuring and reporting on the topic.

247

248 **Measures that will be used for public reporting should meet a high standard of evidence for**  
249 **the focus of measurement.** NQF measures are intended to be useful for public reporting, as  
250 well as to internal quality improvement activities. Measures used for public reporting often  
251 impact large numbers of providers and entail investment of significant resources in  
252 measurement and improvement. Consequently, measures that will be used for public reporting  
253 should meet a high standard of evidence for the focus of measurement. The net benefit to  
254 patients should outweigh any potential harm to patients, and be clinically or practically  
255 meaningful to justify implementation. A lower standard of evidence may be deemed

256 appropriate by those selecting measures for use in smaller scale internal quality improvement  
257 activities within a learning system that allows for rapid adjustments. Such measures, although  
258 potentially of value, are not considered by NQF as they are not appropriate for public reporting.

259

260 **In the absence of strong evidence of certainty of net benefit for a structure or process being**  
261 **measured, expert judgment must conclude that potential benefits to patients clearly**  
262 **outweigh potential harms to patients from the specific structure, intervention or service.**

263 Much of healthcare has not been subjected to research studies and thus, does not have a strong  
264 evidence base. In the absence of strong evidence, clinical interventions and services that are the  
265 focus of quality performance measures should be judged to have benefits to patients that clearly  
266 outweigh any potential risk. In the absence of strong evidence, administrative, management, or  
267 system structures and processes that are the focus of quality performance measures should be  
268 judged to have benefits to patients that clearly outweigh the system costs and resources to  
269 implement those structures and processes.

270

271 **Standards for evidence grading are evolving and expectations for both the present and future**  
272 **should be stated.** Standards for evidence review and grading and clinical practice guideline  
273 development are evolving, as are expectations for measures endorsed by NQF. Explicit  
274 information about the evidence supporting a measure and how (or if) it was graded is essential  
275 for evaluating the evidence both now and in the future.

276

277 **Consistency with prior terminology, whenever possible, minimizes confusion.** Terminology  
278 used in prior NQF documents should be changed only if incorrect or leads to increased  
279 understanding. Whenever possible, narrative descriptions should be used instead of technical  
280 terminology.

281

## 282 I. Recommendations on Sources of Evidence and Evidence Grading for the Present and the Future

- 283 • The preferred sources of evidence are systematic reviews and grading of a body of evidence  
284 conducted by independent organizations such as [USPSTF](#), [AHRQ Evidence-based Practice](#)  
285 [Centers](#), and the [Cochrane Collaboration](#); or guidelines that meet national standards for  
286 trustworthy guidelines (as being developed by the IOM).

- 287 • Until such time when guidelines are certified as meeting a set standard, preferred guidelines  
288 are those developed with balanced representation beyond one specialty group and with full  
289 disclosure of biases and how they were addressed. Further, the evidence underlying a  
290 guideline recommendation must be accessible in order to meet the requirements set out in  
291 this report.
- 292 • An assigned evidence grade alone is not sufficient to evaluate whether the NQF criterion on  
293 evidence for the focus of measurement (1c) is met, either now or in the future. The specific  
294 information on the quantity, quality, and consistency of the body of evidence that was used  
295 to determine an overall grade should be summarized in the measure submission.
- 296 • Explicit, transparent information on the quantity, quality, and consistency of the body of  
297 evidence supporting a measure will facilitate identification of guideline recommendations  
298 that do not have acceptable evidence as the basis for performance measurement. Explicit  
299 information about the evidence also facilitates review by all stakeholders although TAPs  
300 and Steering Committees will continue to include experts that possess knowledge about the  
301 state of science for a particular topic.
- 302 • **Current Expectations**
- 303     ○ Most measure developers will rely on evidence reviews and grading conducted by  
304     other organizations such as guideline developers or published systematic reviews.  
305     However, it is the responsibility of the measure developer to understand the  
306     strength of the evidence on which it is basing a measure and to provide a concise  
307     summary of this evidence, not simply the end-result of the grading process.  
308     Information on the evidence is useful to committees reviewing measures and the  
309     public who use the measures.
- 310     ○ To promote transparency and standardization, NQF should require measure  
311     developers to provide specific information about the quantity, quality, and  
312     consistency of the body of evidence underlying a quality performance measure.  
313     Information should include who graded the evidence, the evidence grading system  
314     used and the grade assigned. If the developer fails to provide this information, NQF  
315     should not review the proposed measure.

316           ○ NQF prefers (but does not require) that submitted evidence be graded based on the  
317           systems of either the [USPSTF](#) or [GRADE](#) because such standardization facilitates  
318           broader understanding of the strength of the evidence.

319   • **Future Expectations**

320   The Task Force identified the following future expectations to signal support for  
321   standardized evidence grading and methods for guideline development. However, even  
322   with standardized grading, reporting the quantity, quality, and consistency of the body of  
323   evidence will be required for transparency and NQF measure evaluation.

- 324           ○ Most measure developers will continue to rely on evidence reviews and grading  
325           conducted by other organizations.
- 326           ○ Rather than identifying “preferred” grading systems as noted for the current  
327           expectations, NQF should require that evidence used to support measures be graded  
328           using one or two standardized evidence grading systems (e.g., the USPSTF, GRADE,  
329           or possibly one adopted by the IOM).
- 330           ○ The evidence should be graded by identified credible sources, such as guideline  
331           developers or review organizations, certified as meeting accepted standards.
- 332           ○ Even when basing measures on evidence graded with a standardized grading  
333           system and potentially certified reviewers, explicit information on the quantity,  
334           quality, and consistency of the specific evidence that led to the assignment of a grade  
335           should be submitted for evaluation.

336

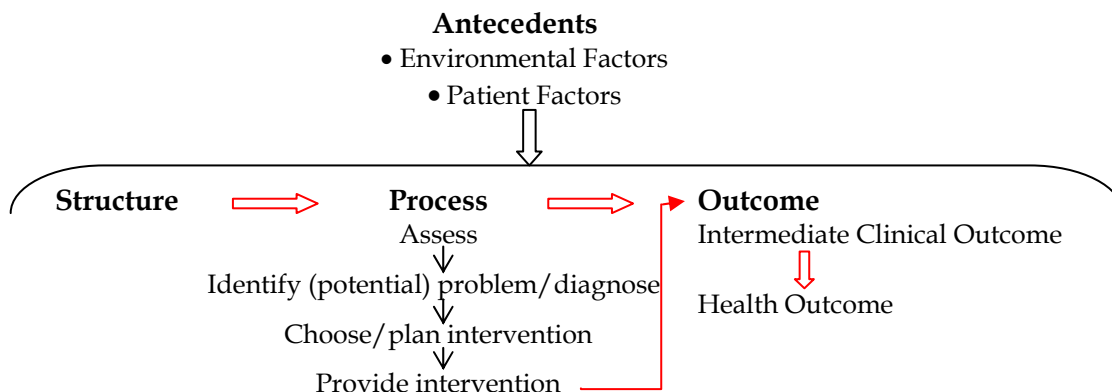
337   **II. Recommendations for the Evidence Needed to Justify the Focus of a Quality Measure**

338   There has been widespread acceptance of Donabedian’s<sup>20,21</sup> structure-process-outcome model  
339   for assessing healthcare quality. These three approaches to quality measurement can be used  
340   with any topic of healthcare quality and the evidence required generally does not vary by topic.  
341   The required evidence is for the links depicted by the red arrows in Figure 2. As depicted under  
342   process, there may be multiple process steps prior to delivering an intervention; however, the  
343   evidence is most often about the relationship between the intervention and outcome and  
344   therefore, interventions are the preferred focus of process measures. Antecedents are depicted  
345   in Figure 2. Although they influence structures, processes, and outcomes, patient factors that  
346   influence outcomes are important to consider for risk adjustment for outcome measures.



347

348 Figure 2. Structure-Process-Outcome Model



349

350 Table 3 outlines the evidence required to justify the structure, process, or outcome that is the  
351 focus of measurement (i.e., what is being measured). It also identifies special considerations  
352 related to certain quality topics. Subsequent tables lay out the approach for evaluating the  
353 evidence and using it to determine if the NQF criterion is met.

354

355 Outcomes as a representation of quality also are based on the process-outcome link. Outcomes  
356 are viewed as useful quality indicators because they are integrative of the influence of multiple  
357 care processes and disciplines involved in the care. However, that also presents some challenges  
358 related to presenting evidence to support the focus of measurement. Optimally, there will be a  
359 body of evidence for the link between the outcome and at least one care process. However, the  
360 lack of such evidence should not necessarily be reason to automatically dismiss the value of  
361 measurement, particularly when the outcome represents a central goal of healthcare treatments  
362 and services (e.g., health, function, survival, symptom control) or harm resulting from  
363 healthcare provided or omitted. Once outcomes are measured and reported, many outcomes  
364 that were not thought to be modifiable tend to be improved and stimulate identification and  
365 adoption of effective practices. If an outcome does not have a body of evidence linking it to a  
366 healthcare process, it may be considered for an exception to the evidence subcriterion if there is  
367 a rationale for the relationship of the outcome to processes of care and/or the importance of  
368 measuring the outcome. Measuring outcomes is important and NQF will need to monitor  
369 whether this guidance on evidence presents a barrier to endorsing reliable and valid outcome  
370 measures.

371 Table 3. Evidence to Support the Focus of Measurement

Type of Measure	Evidence	Example of Measure Type & Evidence to be Addressed
<p><b>Structure</b> Structure of care is a feature of a health care organization or clinician related to its capacity to provide high quality health care</p>	<p>Quantity, quality, and consistency of a body of evidence that the measured healthcare structure leads to desired health outcomes with benefits that outweigh harms (including evidence for the link to effective care processes and the link from the care processes to desired health outcomes) See Table 4</p>	<p><b>#0190 Nurse Staffing Hours</b> <b>Evidence</b> that higher nursing hours are associated with lower mortality, morbidity ; or is associated with effective care processes (e.g., lower medication errors) that lead to better outcomes</p>
<p><b>Process</b> A process of care is a health care-related activity performed for, on behalf of, or by a patient</p>	<p>Quantity, quality, and consistency of a body of evidence that the measured healthcare process leads to desired health outcomes in the target population with benefits that outweigh harms to patients</p> <p>Specific drugs and devices should have FDA approval for the target condition</p> <p>If the measure focus is on inappropriate use: Quantity, quality, and consistency of a body of evidence that the measured healthcare process does <u>not</u> lead to desired health outcomes in the target population See Table 4</p>	<p><b>#0551 ACE Inhibitor / Angiotensin Receptor Blocker(ARB) Use and Persistence Among Members with Coronary Artery Disease at High Risk for Coronary Events</b> <b>Evidence</b> that use of ACE-I and ARB are associated with lower mortality and/or cardiac events</p> <p><b>#0058 Inappropriate antibiotic treatment for adults with acute bronchitis</b> <b>Evidence</b> that antibiotics are not effective for acute bronchitis</p>
<p><b>Intermediate Clinical Outcome</b> An intermediate outcome is a change in physiologic state that leads to a longer-term health outcome</p>	<p>Quantity, quality, and consistency of a body of evidence that at least one healthcare intervention influences the measured intermediate clinical outcome and leads to desired health outcomes See Table 4</p> <p>OR</p> <p>If such evidence does not exist, there is a rationale for the relationship of the intermediate outcome to processes of care and to the desired health outcome. See Table 5</p>	<p><b>#0059 Hemoglobin A1c Management [A1c&gt;9]</b> <b>Evidence</b> that hemoglobin A1c is influenced by interventions (e.g., medication, lifestyle) and is associated with health outcomes (e.g., renal disease, heart disease, amputation, mortality)</p>
<p><b>Health Outcome</b> An outcome of care is a health state of a patient (or change in health status) resulting from healthcare – desirable or adverse</p> <p>In some situations, resource use may be considered a proxy for a health state (e.g., hospitalization may represent a deterioration in health status)</p>	<p>Quantity, quality, and consistency of a body of evidence that the measured health outcome (desirable or adverse) is influenced by at least one healthcare intervention, process, or service. See Table 4</p> <p>OR</p> <p>If such evidence does not exist, there is a rationale for the relationship of the health outcome to processes of care and/or the importance of measuring the outcome. See Table 5</p>	<p><b>#0230 Acute Myocardial Infarction 30-day Mortality</b> Survival is a goal of seeking and providing treatment for AMI <b>Evidence</b> that healthcare processes/ interventions (aspirin, reperfusion) affect mortality/ survival</p> <p><b>#0171 Acute care hospitalization (risk-adjusted) [of home care patients]</b> Improvement or stabilization of condition to remain at home is a goal of seeking and providing home care services. <b>Evidence</b> that healthcare processes (e.g., medication reconciliation, care</p>

Type of Measure	Evidence	Example of Measure Type & Evidence to be Addressed
		coordination) affect hospitalization of patients receiving home care services  <b>#0140</b> Ventilator-associated pneumonia for ICU and high-risk nursery (HRN) patients Avoiding harm from treatment is a goal of when seeking and providing healthcare. <b>Evidence</b> that ventilator acquired pneumonia is affected by healthcare processes (e.g., ventilator bundle)
<b>Special Considerations by Topic</b>		
<b>Patient experience with care</b>	Evidence that the measured aspects of care are those valued by patients and for which the patient is the best and/or only source of information (often acquired through qualitative studies) OR Evidence that patient experience with care is correlated with desired outcomes	<b>#0166</b> HCAHPS <b>Evidence</b> that patients/consumers value the aspects of care being measured (e.g., communication with doctors and nurses, responsiveness of hospital staff, pain control, communication about medicines, cleanliness and quiet of the hospital environment, and discharge information)
<b>Efficiency</b> Measures of efficiency combine the concepts of resource use <u>and</u> quality	<b>Efficiency Measured with combination of Quality measures and Resource Use measures</b>  <b>Quality measure component</b> Evidence for the selected quality measure(s) as described in this table <b>Resource use measure component</b> Does not require clinical evidence as described in this table	Currently, there are no NQF-endorsed efficiency measures that combine quality and resource use  Potential Measure: Diabetes quality measure(s) or composite used in conjunction with a measure of resource use per episode <b>Evidence</b> for diabetes quality measure(s) as described in this table

372

373 **III. Recommendations for Evaluating Criterion 1c – Quantity, Quality, Consistency of Body of**  
374 **Evidence**

375 The following recommendations and decision rules apply to evaluating evidence whether for  
376 initial endorsement, endorsement maintenance, or ad hoc review. The state of science may  
377 change over time, therefore at the time of review for endorsement maintenance, it also is  
378 appropriate to reexamine the evidence to assess whether new and innovative ways of  
379 organizing and providing care have evolved which achieve the same or better outcomes  
380 potentially at less cost.

381

- 382 • Evidence should be evaluated on *quantity* of studies, *quality* of studies, and *consistency* in  
383 direction and magnitude of net benefit (clinically or practically meaningful benefits over  
384 harms to patients) of a ***body of evidence*** on a scale of High, Moderate, or Low.

- 385 • The dimensions of *quantity*, *quality*, and *consistency* of a body of evidence apply to measures  
386 based on guidelines as well as those for which guidelines may not exist (e.g., care  
387 coordination or team functioning may not be based on guidelines, but often have bodies of  
388 evidence including non-clinical literature that should be systematically assessed)
- 389 • Measures without a clear description of the *quantity*, *quality*, and *consistency* of the  
390 supporting body of evidence or without any evidence should not pass criterion 1c and the  
391 threshold criterion of *Importance to Measure and Report*.
- 392 • Use of only selected studies rather than an entire body of evidence that meets pre-  
393 established criteria is not adequate to evaluate the evidence and should not pass criterion 1c  
394 and the threshold criterion of *Importance to Measure and Report*.
- 395 • Inconsistent and conflicting evidence should result in measures not passing both criterion 1c  
396 and the threshold criterion of *Importance to measure and report*.
- 397 • Outcome measures may be considered for an exception to the evidence subcriterion if a  
398 body of evidence linking the outcome to at least one healthcare process does not exist, but  
399 there is a rationale for the relationship of the outcome to processes of care and/or the  
400 importance of measuring the outcome.
- 401 • Expert opinion is not considered empirical evidence and will only be considered in  
402 exceptional circumstances when all of the following conditions are met.
- 403     o No evidence is available.
- 404     o Expert opinion is systematically assessed. That is, identified experts explicitly  
405 address the certainty or confidence that benefits to patients from the specific process  
406 or structure greatly outweigh potential harms, using a specified process that is  
407 transparent and open to peer review (e.g., modified Delphi, formal consensus  
408 process, [RAND Appropriateness Method](#)<sup>22</sup>). The methods and results are reported  
409 for review.
- 410     o There is a strong rationale for why the specific structure or process should be the  
411 focus of a quality performance measure.

412

413 Table 4 provides definitions and guidance on how to evaluate each of the dimensions of  
414 *quantity*, *quality*, and *consistency* for a body of quantitative evidence. Each dimension is rated on  
415 a scale of high, moderate, low, or inadequate to evaluate. A body of evidence could have

416 different ratings for each dimension, e.g., high on quantity, low on quality, and moderate on  
417 consistency. Table 5 provides recommended decision rules for using the ratings for all three  
418 dimensions to make a decision on whether a measure should pass the criterion 1c, the evidence  
419 to support the measure focus. Strong evidence usually requires multiple studies each with  
420 sufficient numbers of patients to give precise estimates, but occasionally a large and  
421 representative study can provide adequate evidence. For example, one study (low quantity) that  
422 is a RCT with a large representative sample of patients (high quality) and substantial estimates  
423 of net benefit would pass the criterion, whereas, a body of evidence with low consistency of  
424 estimates of net benefits indicates a measure should not pass the criterion regardless of the  
425 ratings for quantity and quality of studies.

426

427 There are various ways to categorize research [study designs](#). However, for purposes of the  
428 rating schema, the type of evidence for the structure-process-outcome linkages is categorized  
429 into two categories as follows.

430 **Randomized Controlled Trial (RCT):** Research study design in which subjects are randomly  
431 assigned to various interventions.

432 **Non-RCT:** Research study designs without random assignment to intervention groups,  
433 including quasi-experimental studies, observational studies (e.g., cohort, case-control, cross-  
434 sectional, epidemiologic studies), and qualitative studies.

435

436 Although RCTs remain the gold standard for evidence of efficacy of treatment, there are many  
437 areas where RCTs may not currently exist and are unlikely to be conducted. Furthermore, the  
438 strict eligibility and exclusion criteria for randomized trials may sometimes result in findings  
439 that are not fully generalizable in real world applications. NQF recognizes the evidentiary value  
440 of well-conducted observational studies, particularly those that attempt to balance measured  
441 covariates (e.g., using propensity scores) and account for other sources of bias as articulated in  
442 the [GRACE principles](#) [Good Research for Comparative Effectiveness] <sup>23</sup>. This is particularly  
443 true when there are multiple such studies that arrive at similar conclusions.

444

445 Qualitative studies often are used to gain understanding of people’s attitudes, behaviors, and  
446 values and may be suited to evidence regarding patient experience with care. Table 4 does not

447 apply to qualitative evidence. When qualitative studies are used, appropriate qualitative  
448 research criteria should be used to judge the strength of the evidence <sup>24</sup>.

449

450 Quality improvement studies are not among the types of study designs listed above, but quality  
451 improvement may be a topic of study. Quality improvement studies may include a variety of  
452 study designs from RCTs to qualitative studies. They could be included in a body of evidence  
453 and the assessment of the strength of evidence would not differ from that of other studies.

454

455

456 Table 4. Evaluation of Quantity, Quality, and Consistency of Body of Evidence for Criterion 1c – evidence for  
 457 the measure focus

Definition/ Rating	Quantity of Body of Evidence	Quality of Body of Evidence	Consistency of Results of Body of Evidence
<b>Definition</b>	Total number of studies (not articles or papers)	Certainty or confidence in the estimates of benefits and harms to patients across studies in the body of evidence related to <a href="#">study factors*</a> including: study design or flaws; directness/indirectness (regarding: the specific process or structure that is the measure focus, outcomes assessed, target population, comparisons); imprecision (wide confidence intervals due to few patients or events)	Stability in both the direction and magnitude of clinically/practically meaningful benefits and harms to patients (benefit over harms) across studies in the body of evidence
<b>High</b>	5+ studies**	Randomized controlled trials (RCTs) of direct evidence, with adequate size to obtain precise estimates of effect, and without serious flaws that introduce bias	Estimates of clinically/practically meaningful benefits and harms to patients are consistent in direction, and similar in magnitude across the preponderance of studies in the body of evidence
<b>Moderate</b>	2-4 studies**	<ul style="list-style-type: none"> <li>• Non-RCTs with control for confounders that could account for other plausible explanations, with large, precise estimate of effect;</li> <li>OR</li> <li>• RCTs without serious flaws that introduce bias, but with either indirect evidence, or imprecise estimate of effect</li> </ul>	<p>Estimates of clinically/practically meaningful benefits and harms to patients are consistent in direction across the preponderance of studies in the body of evidence, but may differ in magnitude</p> <p>If only one study, the estimate of benefits greatly outweighs the estimate of potential harms to patients (1 study cannot achieve high consistency rating)</p>
<b>Low</b>	0-1 studies**	<ul style="list-style-type: none"> <li>• RCTs with flaws that introduce bias;</li> <li>OR</li> <li>• Non-RCTs with small or imprecise estimate of effect, or without control for confounders that could account for other plausible explanations</li> </ul>	<p>Estimates of clinically/practically meaningful benefits and harms to patients differ in both direction and magnitude across the preponderance of studies in the body of evidence; OR wide confidence intervals prevent estimating net benefit</p> <p>If only 1 study, estimate of benefits do not greatly outweigh harms to patients</p>
<b>Inadequate to Evaluate</b> <i>See Table 5 for exceptions</i>	No empirical evidence; OR only selected studies from a larger body of evidence	No empirical evidence; OR only selected studies from a larger body of evidence	No assessment of magnitude and direction of benefits and harms to patients

458 \*Study designs that affect certainty of confidence in estimates of effect include: Randomized controlled  
 459 trials (RCT), which control for both observed and unobserved confounders, and non-RCTs (observational  
 460 studies) with various levels of control for confounders.  
 461 Study flaws that may bias estimates of effect include: lack of allocation concealment; lack of blinding;  
 462 large losses to follow-up; failure to adhere to intention to treat analysis; stopping early for benefit; failure  
 463 to report important outcomes.  
 464 Imprecision with wide confidence intervals around estimates of effects can occur in studies involving few  
 465 patients and few events.  
 466 Indirectness of evidence includes: indirect comparisons (e.g., two drugs compared to placebos rather than  
 467 head-to-head), differences between the population, intervention, outcome of interest, or comparator  
 468 interventions and those included in the relevant studies.<sup>14</sup>  
 469 \*\* The suggested number of studies for rating levels of quantity is considered a general guideline.  
 470

471 Table 5. Evaluation of Subcriterion 1c based on the quantity, quality and consistency of the body of evidence

Quantity of Body of Evidence	Quality of Body of Evidence	Consistency of Body of Evidence	Pass Subcriterion 1c
Moderate-High	Moderate-High	Moderate-High	Yes
Low	Moderate-High	Moderate (if only 1 study high consistency not possible)	Yes, but only if judgment that additional research is unlikely to change conclusion that benefits to patients outweigh harms; otherwise, No
Moderate-High	Low	Moderate-High	Yes, but only if judgment that potential benefits to patients clearly outweigh potential harms; otherwise, No
Low-Mod-High	Low-Mod-High	<b>Low</b>	No
Low	Low	Low	No
<b>Potential Exceptions to Empirical Evidence</b> <ul style="list-style-type: none"> <li>• For a <i>structure or process measure</i>: there is no empirical evidence, <u>and</u> expert opinion is systematically assessed with agreement that the benefits to patients greatly outweigh potential harms and there is a strong rationale for the importance of measuring performance.</li> <li>• For a <i>health outcome measure</i>: a body of evidence linking the outcome to at least one healthcare process does not exist, <u>and</u> there is a rationale for the relationship of the outcome to processes of care and/or the importance of measuring the outcome.</li> </ul>			Yes, but only if judgment that potential benefits to patients clearly outweigh potential harms; otherwise, No

472

473 **IV. Recommendations for Selecting the Focus for Measure Development**

474 Based on its discussion and recommendations regarding evidence to support the measure focus,  
 475 the following recommendations address selecting a focus for measure development.  
 476



- 477 • For any topic area, measures based on the best evidence should be considered over  
478 measures based on lower quality evidence (e.g., expert opinion).
- 479 • There is a hierarchical preference for outcome measures (when possible) followed by  
480 process measures, then structure measures. Outcome measures are preferred because  
481 improving health outcomes is a central goal of healthcare. However, both outcome and  
482 process measures have advantages and disadvantages<sup>25</sup> and both have a place in quality  
483 assessment and the NQF portfolio.
- 484 • Specific drugs and devices included in quality performance measures should be FDA-  
485 approved for the target condition.
- 486 • Structural measures are appropriate primarily when there are very well established  
487 structure-process-outcome relationships; and when it is not feasible to directly measure the  
488 outcome or processes.
- 489 • For process and structure measures, the focus of measurement should be on the aspect of  
490 care with the most direct evidence of a strong relationship to the desired outcome. For  
491 example, evidence about effective medication to control blood pressure is direct evidence  
492 for the medication but only indirect evidence for the frequency of assessing blood pressure  
493 (see Figure 2). Assessment of blood pressure, although necessary, is not sufficient to  
494 achieving control. When there are multiple processes that affect a desired outcome, efforts  
495 should be made to include measures for all processes that have a strong relationship to the  
496 desired outcome.

497

## 498 V. Recommendations for Evaluating Importance to Measure and Report and the Other Subcriteria

499 Although the criterion *Importance to Measure and Report* has been a threshold, must-pass  
500 criterion, the weight of the individual subcriteria in making the determination of whether the  
501 criterion was met was not specified. The Task Force recommended that all three subcriteria  
502 must be met: High impact (1a), Opportunity for improvement (1b), and Evidence for the focus  
503 of measurement (1c) as noted above.

504

505 Generally, in measure submissions, high impact is easily demonstrated by alignment with a  
506 specific NPP goal or epidemiologic or resource use data (incidence, prevalence, resource use,  
507 consequences of quality problems). However, data on opportunity for improvement may be

508 lacking (e.g., submitter states that performance is unknown, or it may not be specific to the  
509 focus of measurement, or only based on a sample from measure development and testing).  
510 Reviewers sometimes question whether there is enough variation to justify importance to  
511 measure and report, or how to judge overall poor performance. When a measure undergoes  
512 review for continued endorsement, an issue that sometimes arises is whether a measure is  
513 “topped out” meaning there are high levels of performance with little variation and therefore,  
514 little room for further improvement.

515  
516 The Task Force did not recommend specific quantitative thresholds for identifying conformance  
517 with the subcriteria of high impact (1a) and opportunity for improvement (1b). Threshold  
518 values for opportunity for improvement would be difficult to standardize. It depends on the  
519 size of the population at risk, effectiveness of an intervention, and the consequences of the  
520 quality problem. For example, even modest variation would be sufficient justification for some  
521 highly effective, potentially life-saving treatments (e.g., certain vaccinations) that are critical to  
522 the public health.

523  
524 The Task Force noted that at the time of review for endorsement maintenance, measure  
525 performance data that indicates overall high performance with little variation would require  
526 justification to continue endorsement. The CSAC added that the default action should be  
527 removal of endorsement unless there is a strong justification to continue endorsement. Failing  
528 opportunity for improvement (subcriterion 1b) results in not passing the threshold criterion,  
529 *Importance to Measure and Report* and thus the measure is not suitable for endorsement. The  
530 CSAC noted that opportunity for improvement also could be considered at the time of review  
531 of measures with time-limited endorsement if there were enough data to make such a  
532 judgment.

533  
534 Measures with overall high performance and little variation might be considered for inclusion  
535 in composite measures; however that does not reduce measurement burden. Additionally, the  
536 measure would still require evaluation of the measure properties because sometimes overall  
537 high performance is a symptom of problems with the measure construction. Further, it would

538 require the analysis of the relationship and contribution of the component measures to the  
539 composite score called for in the composite measure evaluation criteria.

540

541 Recommendations related to opportunity for improvement (1b) include the following.

542 • At the time of initial endorsement, evidence for opportunity for improvement will be based  
543 on research studies, or epidemiologic or resource use data. However, at the time of review  
544 for endorsement maintenance, the primary interest is on the endorsed measure as specified,  
545 and the evidence for opportunity for improvement should be based on data on the specific  
546 endorsed measure.

547 • When assessing measure performance data for opportunity for improvement, the following  
548 factors should be considered:

549     ○ number and representativeness of the entities included in the measure performance  
550     data; and

551     ○ size of the population at risk, effectiveness of an intervention, likely occurrence of an  
552     outcome, and consequences of the quality problem.

553 • At the time of review for endorsement maintenance, an overall high level of performance  
554 with little variation in the endorsed measure scores should result in removal of  
555 endorsement. If other evidence (e.g., epidemiologic or research) is consistent with the  
556 measure performance data, it confirms the lack of opportunity for improvement. If other  
557 evidence is not consistent with the measure performance data, it is suggestive of potential  
558 problems with the measure as specified.

559 • In exceptional situations, a strong justification for continuing endorsement could be  
560 considered (e.g., evidence that overall performance will likely deteriorate if not monitored  
561 and the magnitude of potential harm if outcomes deteriorate when not monitored).

562

563

564 Table 6. Evidence for Evaluating Importance to Measure and Report

Pass Criterion, Importance to Measure and Report?			
All 3 subcriteria (1a,1b,1c) must be met to pass the threshold criterion, <i>Importance to Measure and Report</i>			
Subcriterion	Evidence	Example	Pass the subcriterion?
High impact (1a)	Addresses a <a href="#">specific national health goal/priority</a> identified by the Secretary of DHHS or the NPP; <b>OR</b> Epidemiologic or resource use data; health services research – affects large numbers of patients and/or has a very substantial impact for smaller populations; leading cause of morbidity/mortality; high resource use (current and/or future); severity of illness; and patient/societal consequences of poor quality	#0140 Ventilator-associated pneumonia for ICU and high-risk nursery (HRN) patients <b>NPP goal:</b> . . . focus relentlessly on continually reducing and seeking to eliminate all healthcare-associated infections (HAIs) <b>Evidence</b> related to numbers of patients (e.g., 250,205 VAPs reported 35,969 (14.4%) were fatal; cost (e.g., total annual cost of VAP \$2.5 billion)	<b>Subcriterion 1a</b>  <b>Yes</b> – Demonstrated at least one of the aspects of high impact  <b>No</b> – Did not demonstrate at least one of the aspects of high impact
Opportunity for improvement (1b)	<b>Initial Endorsement</b> Epidemiologic or resource use data; health services research – data demonstrating considerable variation, or overall less than optimal performance, for the focus of measurement across providers and/or population groups (disparities in care)  <b>Review for Endorsement Maintenance</b> Data for the measure as specified and endorsed demonstrating considerable variation, or overall less than optimal performance	#0432 Influenza Vaccination of Nursing Home/ Skilled Nursing Facility Residents <b>NPP goal:</b> All Americans will receive the most effective preventive services recommended by the U.S. Preventive Services Task Force <b>Evidence</b> that vaccination rates vary (e.g., 39% fail to reach the Healthy People 2010 objective of vaccinating at least 90% of nursing home residents)	<b>Subcriterion 1b</b>  <b>Yes</b> – Demonstrated either variation or overall less than optimal performance  <b>No</b> – Did not demonstrate either variation <b>or</b> overall less than optimal performance
Evidence for the focus of measurement (1c)	See Table 3	See Table 3	<b>Subcriterion 1c</b>  See Table 4 and Table 5
<b>All 3 subcriteria (1a,1b,1c) must be met</b> to pass the threshold criterion, <i>Importance to Measure and Report</i>			

565

566 Consequences of Measurement

567 Consequences of measurement are not the same as the consequences of implementing the

568 measured structure or process, i.e., the benefits or harms to the patient related to the specific

569 topic of measurement. Currently, unintended consequences of measurement are addressed  
 570 under feasibility.

571 **4d.** Susceptibility to inaccuracies, errors, or unintended consequences of measurement and the  
 572 ability to audit the data items to detect such problems are identified.

573  
 574 The Task Force identified that actual vs. theoretical consequences to measurement are most  
 575 likely to arise after implementation and should be addressed at the time of review for  
 576 endorsement maintenance. For example, a measure of timing of antibiotic administration in  
 577 patients with pneumonia may result in some patients receiving antibiotics before the diagnosis  
 578 of pneumonia is confirmed by x-ray. The Task Force did not recommend moving subcriterion  
 579 4d under *Importance to Measure and Report*.

580  
 581 **VI. Recommendations for Modifications to the NQF Evaluation Criteria**

582 The following criteria reflect changes to implement the recommendations including that all  
 583 three subcriteria be met to pass the threshold criterion of *Importance to Measure and Report*.

584  
 585 Table 7. Current and Modified Measure Evaluation Criteria

Current Measure Evaluation Criteria	Modified Measure Evaluation Criteria
<p><b>1. Importance to measure and report:</b> 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. <i>Candidate measures must be judged to be important to measure and report in order to be evaluated against the remaining criteria.</i></p> <p><b>1a.</b> The measure focus addresses:</p> <ul style="list-style-type: none"> <li>• a specific national health goal/priority identified by NQF’s National Priorities Partners;</li> <li>OR</li> <li>• 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).</li> </ul> <p><b>1b.</b> Demonstration of quality problems and opportunity for improvement, i.e., data (1) demonstrating</p>	<p><b>1. Importance to measure and report:</b> Extent to which the specific measure focus is evidence-based, important to making significant gains in health care quality and improving health outcomes for a specific high impact aspect of healthcare where there is variation in or overall poor performance. <i>Candidate measures must be judged to be important to measure and report in order to be evaluated against the remaining criteria.</i></p> <p><b>1a.</b> The measure focus addresses:</p> <ul style="list-style-type: none"> <li>• a specific national health goal/priority identified by DHHS or the <a href="#">National Priorities Partnership</a> convened by NQF;</li> <li>OR</li> <li>• a demonstrated high impact aspect of healthcare (e.g., affects large numbers of patients and/or has a substantial impact for a smaller population; leading cause of morbidity/mortality; high resource use (current and/or future); severity of illness; and severity of patient/societal consequences of poor quality).</li> </ul> <p><b>AND</b></p> <p><b>1b.</b> Demonstration of quality problems and opportunity</p>

Current Measure Evaluation Criteria	Modified Measure Evaluation Criteria
<p>considerable variation, or overall poor performance, in the quality of care across providers and/or population groups (disparities in care).</p> <p><b>1c.</b> The measure focus is:</p> <ul style="list-style-type: none"> <li>• 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 addressed (2); OR</li> <li>• if an intermediate outcome, process, structure, etc., there is evidence (3) that supports the specific measure focus as follows: <ul style="list-style-type: none"> <li>o <u>Intermediate outcome</u> – evidence that the measured intermediate outcome (e.g., blood pressure, Hba1c) leads to improved health/avoidance of harm or cost/benefit.</li> <li>o <u>Process</u> – 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 process (4), it measures the step that has the greatest effect on improving the specified desired outcome(s).</li> <li>o <u>Structure</u> – 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.</li> <li>o <u>Patient experience</u> – evidence that an association exists between the measure of patient experience of health care and the outcomes, values and preferences of individuals/ the public.</li> <li>o <u>Access</u> – evidence that an association exists between access to a health service and the outcomes of, or experience with, care.</li> <li>o <u>Efficiency</u> (5) – 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.</li> </ul> </li> </ul> <p><i>If not important to measure and report, STOP.</i></p> <p><b>Footnotes</b>  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.  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.</p>	<p>for improvement, i.e., data (<a href="#">footnote 1</a>) demonstrating considerable variation, or overall less than optimal performance, in the quality of care across providers and/or population groups (disparities in care).</p> <p><b>AND</b></p> <p><b>1c.</b> The measure focus is evidence-based as demonstrated by a systematic assessment and grading of the quantity, quality, and consistency of the body of evidence (<a href="#">footnote 3</a>).</p> <ul style="list-style-type: none"> <li>• Health outcome/intermediate clinical outcome (<a href="#">footnote 2</a>): evidence that the measured outcome (desirable or adverse) is influenced by at least one healthcare intervention, process, or service; OR there is a rationale for the relationship of the outcome to processes of care and/or the importance of measuring the outcome.</li> <li>• Process (<a href="#">footnote 4</a>): evidence that the measured healthcare process leads to desired outcomes in the target population.</li> <li>• Structure: evidence that the measured structure leads to desired health outcomes (including evidence for the link to effective care processes and the link from the care processes to desired health outcomes).</li> <li>• Special Considerations by Topic of Measurement <ul style="list-style-type: none"> <li>o Patient experience with care: evidence that the measured aspects of care are those valued by patients and for which the patient is the best and/or only source of information OR that patient experience with care is correlated with desired outcomes.</li> <li>o Efficiency (<a href="#">footnote 5</a>): evidence for the quality component as noted above.</li> </ul> </li> </ul> <p><b>Footnotes</b>  1 Examples of data on opportunity for improvement include, but are not limited to: prior studies, epidemiologic data, or data from pilot testing or implementation of the proposed measure. If data are not available, the measure focus is systematically assessed (e.g., expert panel rating) and judged to be a quality problem.  2 Generally, rare event outcomes do not provide adequate information for improvement or discrimination; however, serious reportable events that are compared to zero are appropriate outcomes for public reporting and quality improvement.  3 The preferred systems for grading the evidence are the USPSTF <a href="#">grading definitions</a> and <a href="#">methods</a>, or <a href="#">GRADE</a>.  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 strongest</p>

Current Measure Evaluation Criteria	Modified Measure Evaluation Criteria
<p>3 The strength of the body of evidence for the specific measure focus should be systematically assessed and rated (e.g., USPSTF grading system – <a href="#">grade definitions</a> and <a href="#">methods</a>). 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.</p> <p>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.</p> <p>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 <a href="#">Measurement Framework: Evaluating Efficiency Across Episodes of Care</a>; based on <a href="#">AQA Principles of Efficiency Measures</a>).</p>	<p>evidence for the link to the desired outcome should be selected as the focus of measurement.</p> <p><sup>5</sup> Measures of efficiency combine the concepts of resource use and quality (NQF’s <a href="#">Measurement Framework: Evaluating Efficiency Across Episodes of Care</a>; <a href="#">AQA Principles of Efficiency Measures</a>).</p>

586

587 **VII. Recommendations for Modifications to the Measure Submission**

588 The information requested on NQF’s measure submission form is consistent with those  
589 identified in a 2009 collaborative effort undertaken with AHRQ, CMS, The Joint Commission,  
590 NCQA, and PCPI to identify common data fields. The Task Force suggested modifications to  
591 the information requested on the NQF [measure submission form](#) to implement the above  
592 recommendations.

593

594 The intent is full transparency about the supporting evidence for the submitted measure. This  
595 will facilitate understanding of the adequacy of the evidence presented (selected evidence vs. a

596 body of evidence) and the developer’s representation of the quality of the evidence. Currently,  
 597 evidence graded using the USPSTF or GRADE systems may not be available, however, an  
 598 accurate description of the evidence and any grading system used should still be expected. The  
 599 following items pertain to the recommendations related to evidence (1c) under *Importance to*  
 600 *Measure and Report*.

601

602 Table 8. Current and Modified Measure Submission Items

Current Measure Submission (4.1) Items	Modified Measure Submission Items
	<p><b>Add to Introduction</b>  <i>Importance to Measure and Report</i> is a threshold criterion that must be met in order to recommend a measure for endorsement. All three subcriteria (1a, 1b, and 1c) must be met in order to pass this criterion. The following items request the information the committees will need to evaluate whether the criterion is met.</p>
<p><b>High Impact (Measure evaluation criterion 1a)</b>  <b>1a.1. Demonstrated High Impact Aspect of Healthcare</b>            Affects large numbers            Leading cause of morbidity/mortality            Severity of illness            Patient/societal consequences of poor quality            Frequently performed procedure            High resource use            Other:</p> <p><b>1a.3. Summary of Evidence of High Impact</b></p> <p><b>1a.4. Citations for Evidence of High Impact</b></p> <p><b>Opportunity for Improvement (Measure evaluation criterion 1b)</b>  <b>1b.1. Briefly explain the benefits (improvements in quality) envisioned by use of this measure</b></p> <p><b>1b.2. Summary of Data Demonstrating Performance Gap</b> (<i>Variation or overall poor performance across providers</i>)</p> <p><b>1b.3. Citations for Data on Performance Gap</b></p> <p><b>1b.4. Summary of Data on Disparities by Population Group</b></p>	<p><b>High Impact (Measure evaluation criterion 1a)</b>  <b>1a.1. Demonstrated High Impact Aspect of Healthcare</b>            Affects large numbers            Leading cause of morbidity/mortality            Severity of illness            Patient/societal consequences of poor quality            Frequently performed procedure            High resource use            Other:</p> <p><b>1a.3. Summary of Evidence of High Impact</b> (<i>provide epidemiologic or resource use data</i>)</p> <p><b>1a.4. Citations for Evidence of High Impact</b></p> <p><b>Opportunity for Improvement (Measure evaluation criterion 1b)</b>  <b>1b.1. Briefly explain the benefits (improvements in quality) envisioned by use of this measure</b></p> <p><b>1b.2. Summary of Data Demonstrating Performance Gap</b> (<i>Variation or overall poor performance across providers</i>)</p> <p><b>1b.3. Citations for Data on Performance Gap</b></p> <p><b>1b.4. Summary of Data on Disparities by Population Group</b></p> <p><b>1b.5. Citations for Data on Disparities</b></p>



Current Measure Submission (4.1) Items	Modified Measure Submission Items
<p><b>1b.5. Citations for Data on Disparities</b></p> <p><b>1c.1. Relationship to Outcomes</b> <i>(For non-outcome measures, briefly describe the relationship to desired outcome. For outcomes, describe why it is relevant to the target population.)</i></p> <p><b>1c.2. Type of Evidence</b> <i>(Check all that apply)</i>  Cohort study  Observational study  Evidence-based guideline  Randomized controlled trial  Expert opinion  Systematic synthesis of research  Meta-analysis  Other: 1c.3.</p> <p><b>1c.4. Summary of Evidence</b> <i>(For non-outcome measures, provide evidence of relationship to desired outcome. For outcomes, summarize any evidence that healthcare services/care processes influence the outcome.)</i></p> <p><b>1c.5. Rating of Strength/Quality of Evidence</b> <i>(Also provide narrative description of the rating and by whom)</i></p> <p><b>1c.6. Method for Rating Evidence</b></p> <p><b>1c.7. Summary of Controversy/Contradictory Evidence</b></p>	<p><b>1c.1. Structure-Process-Outcome Relationship</b> <i>(Briefly state the measured structure, process, or outcome and the links and direction between: a) the measured process and desired outcome; b) the measured outcome and processes that influence the outcome; or c) the measured structure and effective processes and desired outcome.)</i></p> <p><b>1c.2. Source of Evidence</b>  Clinical practice guideline  Systematic review of body of evidence (other than within guideline development)  Selected individual studies (rather than entire body of evidence)  Other 1c.3.</p> <p><b>1c.4. Summary of Body of Evidence</b>  Quantity of Studies in Body of Evidence <i>(total number of studies, not articles):</i>  Quality of Body of Evidence <i>(Certainty or confidence in the estimates of benefits and harms to patients <u>across studies</u> in the body of evidence resulting from <u>study factors</u> including: study design/ flaws; directness/indirectness regarding the specific process/structure being measured, outcomes assessed, target population, comparisons; imprecision (wide confidence intervals due to few patients or events):</i>  Directness to focus of measurement &amp; target population in proposed measure:  Consistency of Results across Studies:  Net Benefit <i>(Benefits over harms)</i>  Benefit/outcome – estimate of effect  Harms addressed – estimate of effect</p> <p><b>1c.5. Grading of Strength/Quality of Body of Evidence</b>  Has the <b>body of evidence</b> been graded? Yes No  If graded:  By whom <i>(describe the entity that graded the evidence, including balance of representation and any disclosures regarding bias)</i>  Grade Assigned to the Evidence:</p> <p><b>1c.6. System for Grading Evidence:</b> USPSTF GRADE  Other <i>(provide description of grading scale with definitions)</i></p> <p><b>1c.7. Summary of Controversy/Contradictory Evidence</b></p>

Current Measure Submission (4.1) Items	Modified Measure Submission Items
<p><b>1c.8. Citations for Evidence</b> (<i>Other than guidelines</i>)</p> <p><b>1c.9. Quote the Specific Guideline Recommendation</b> (<i>Including guideline number and/or page number</i>)</p> <p><b>1c.10. Clinical Practice Guideline Citation</b></p> <p><b>1c.11. National Guideline Clearinghouse or Other URL</b></p> <p><b>1c.12. Rating Strength of Recommendation</b> (<i>Also provide narrative description of the rating and by whom</i>)</p> <p><b>1c.13. Method for Rating Strength of Recommendation</b> (<i>If different from USPSTF system, also describe rating and how it relates to USPSTF</i>)</p> <p><b>1c.14. Rationale for Using This Guideline Over Others</b></p>	<p><b>1c.8. Citations for Evidence</b> (<i>Other than guidelines</i>)</p> <p><b>1c.9. Quote Verbatim the Specific Guideline Recommendation</b> (<i>Including guideline number and/or page number</i>)</p> <p><b>1c.10. Clinical Practice Guideline Citation</b></p> <p><b>1c.11. National Guideline Clearinghouse or Other URL for the cited guideline</b></p> <p><b>1c.12. Grading of Strength of Guideline Recommendation</b>  Has the <b>recommendation</b> been graded? Yes No  If graded:  By whom (<i>describe the entity that graded the evidence, including balance of representation and any disclosures regarding bias</i>)  Grade Assigned to the Recommendation:</p> <p><b>1c.13. System for Grading Strength of Guideline Recommendation:</b> USPSTF GRADE Other (<i>provide description of grading scale with definitions</i> )</p> <p><b>1c.14. Rationale for Using This Guideline Over Others</b></p> <p><b>1c.15 Based on the NQF descriptions for rating the evidence, what was your assessment of the quantity, quality, and consistency of the body of evidence?</b> (rate each as High, Moderate, or Low)  Quantity:  Quality:  Consistency:</p>
<p><b>Descriptive Information</b></p> <p><b>De.4. National Priority Partnership priority area</b> (<i>Select the most relevant</i>)  Patient and family engagement  Population health  Safety  Care coordination  Palliative and end of life care  Overuse</p> <p><b>De.5. IOM Quality Domain</b> (<i>Select the most relevant</i>)  Effectiveness  Efficiency  Equity  Patient-centered  Safety  Timeliness</p> <p><b>De.6. Consumer Care Need</b> (<i>Select the most relevant</i>)  Getting better</p>	<p><b>Descriptive Information - no change</b></p> <p><b>De.4. National Priority Area</b> (<i>Select the most relevant</i>)  [May change with DHHS priorities]  Patient and family engagement  Population health  Safety  Care coordination  Palliative and end of life care  Overuse</p> <p><b>De.5. IOM Quality Domain</b> (<i>Select the most relevant</i>)  Effectiveness  Efficiency  Equity  Patient-centered  Safety  Timeliness</p> <p><b>De.6. Consumer Care Need</b> (<i>Select the most relevant</i>)  Getting better</p>

Current Measure Submission (4.1) Items	Modified Measure Submission Items
Living with illness Staying healthy	Living with illness Staying healthy

603

604 **VIII. Recommendations for Evidence Required for Practices Considered for NQF Endorsement**

605 NQF also endorses practices such as [safe practices](#), care coordination practices, and substance  
606 use treatment practices. The [criteria](#) for practices include evidence of effectiveness.

607

608 The Task Force recommends that the same evidence requirements as indicated for process  
609 measures (Tables 3, 4, 5) be applied to practices considered for NQF endorsement.

610

611 **Table 9. Evidence to Support a Practice**

612

Evidence to Support a Practice	Example of Practice & Evidence to be Addressed
Quantity, quality, and consistency of a body of evidence that the measured healthcare process leads to desired health outcomes in the target population with benefits that outweigh harms to patients	<b>Safe Practice 16</b> Safe Adoption of Computerized Prescriber Order Entry <b>Evidence</b> that computerized order entry systems are associated with lower medication errors and adverse events

613

614 **Modifications to Practice Evaluation Criteria**

615 **Evidence of Effectiveness.** A practice is evidence-based as demonstrated by a systematic  
616 assessment of the quantity, quality, and consistency of the body of evidence and standardized  
617 grading of the body of evidence. The preferred systems for grading the evidence are the  
618 USPSTF [grading definitions](#) and [methods](#), or [GRADE](#). Evidence from non-healthcare industries  
619 that should be substantially transferable to healthcare (e.g., safety practices of repeat-back of  
620 verbal orders or standardizing abbreviations) also may be considered.

621

622

623 **REFERENCES**

- 624 1. Lohr KN. Rating the strength of scientific evidence: relevance for quality improvement  
625 programs. *Int J Qual Health Care*. 2004;16(1):9-18.  
626 2. Tricoci P, Allen JM, Kramer JM et al. Scientific evidence underlying the ACC/AHA clinical  
627 practice guidelines. *JAMA*. 2009;301(8):831-841.

- 628 3. Spertus JA, Eagle KA, Krumholz HM et al. American College of Cardiology and American  
629 Heart Association methodology for the selection and creation of performance measures  
630 for quantifying the quality of cardiovascular care. *Circulation*. 2005;111(13):1703-1712.
- 631 4. Physician Consortium for Performance Improvement. Physician Consortium for  
632 Performance Improvement® (PCPI) Position Statement - The Evidence Base Required for  
633 Measures Development. *American Medical Association* 6-26-2009;1-18. Last accessed March  
634 2010.
- 635 5. Grilli R, Magrini N, Penna A et al. Practice guidelines developed by specialty societies: the  
636 need for a critical appraisal. *Lancet*. 2000;355(9198):103-106.
- 637 6. Shiffman RN, Shekelle P, Overhage JM et al. Standardized reporting of clinical practice  
638 guidelines: a proposal from the Conference on Guideline Standardization. *Ann Intern Med*.  
639 2003;139(6):493-498.
- 640 7. The AGREE Collaboration. *Appraisal of Guidelines for Research and Evaluation AGREE*  
641 *Instrument*. 2001. Available at <http://www.agreecollaboration.org/instrument/>. Last  
642 accessed February 2010.
- 643 8. The AGREE Collaboration. Development and validation of an international appraisal  
644 instrument for assessing the quality of clinical practice guidelines: the AGREE project.  
645 *Quality & safety in health care*. 2003;12(1):18-23.
- 646 9. Atkins D, Eccles M, Flottorp S et al. Systems for grading the quality of evidence and the  
647 strength of recommendations I: critical appraisal of existing approaches The GRADE  
648 Working Group. *BMC Health Serv Res*. 2004;4(1):38.
- 649 10. Atkins D, Best D, Briss PA et al. Grading quality of evidence and strength of  
650 recommendations. *BMJ*. 2004;328(7454):1490-1494.
- 651 11. Guyatt GH, Oxman AD, Vist GE et al. GRADE: an emerging consensus on rating quality of  
652 evidence and strength of recommendations. *BMJ*. 2008;336(7650):924-926.
- 653 12. Guyatt GH, Oxman AD, Kunz R et al. Incorporating considerations of resources use into  
654 grading recommendations. *BMJ*. 2008;336(7654):1170-1173.
- 655 13. Guyatt GH, Oxman AD, Kunz R et al. Going from evidence to recommendations. *BMJ*.  
656 2008;336(7652):1049-1051.
- 657 14. Guyatt GH, Oxman AD, Kunz R et al. What is "quality of evidence" and why is it important  
658 to clinicians? *BMJ*. 2008;336(7651):995-998.
- 659 15. Guyatt GH, Oxman AD, Vist GE et al. GRADE: an emerging consensus on rating quality of  
660 evidence and strength of recommendations. *BMJ*. 2008;336(7650):924-926.
- 661 16. Harris RP, Helfand M, Woolf SH et al. Current methods of the US Preventive Services Task  
662 Force: a review of the process. *Am J Prev Med*. 2001;20(3 Suppl):21-35.
- 663 17. Sawaya GF, Guirguis-Blake J, LeFevre M et al. Update on the methods of the U.S. Preventive  
664 Services Task Force: estimating certainty and magnitude of net benefit. *Ann Intern Med*.  
665 2007;147(12):871-875.
- 666 18. Owens DK, Lohr KN, Atkins D et al. Grading the strength of a body of evidence when  
667 comparing medical interventions-Agency for Healthcare Research and Quality and the  
668 Effective Health Care Program. *J Clin Epidemiol*. 2009.
- 669 19. Liberati A, Altman DG, Tetzlaff J et al. The PRISMA statement for reporting systematic  
670 reviews and meta-analyses of studies that evaluate health care interventions: explanation  
671 and elaboration. *Ann Intern Med*. 2009;151(4):W65-W94.
- 672 20. Donabedian A. *An Introduction to Quality Assurance in Health Care*. New York, NY: Oxford  
673 University Press; 2003.

- 674 21. Donabedian A. The role of outcomes in quality assessment and assurance. *Quality Review*  
675 *Bulletin*. 1992;18(11):356-360.
- 676 22. Fitch K, Bernstein SJ, Aguilar MS et al. *The RAND/UCLA Appropriateness Method User's*  
677 *Manual*. Santa Monica, CA: RAND Health; 2000. Available at  
678 [http://www.rand.org/pubs/monograph\\_reports/MR1269/](http://www.rand.org/pubs/monograph_reports/MR1269/).
- 679 23. Dreyer NA, Schneeweiss S, McNeil BJ et al. GRACE principles: recognizing high-quality  
680 observational studies of comparative effectiveness. *Am J Manag Care*. 2010;16(6):467-471.
- 681 24. Cohen DJ, Crabtree BF. Evaluative criteria for qualitative research in health care:  
682 controversies and recommendations. *Ann Fam Med*. 2008;6(4):331-339.
- 683 25. Donabedian A. The role of outcomes in quality assessment and assurance. *Quality Review*  
684 *Bulletin*. 1992;18(11):356-360.
- 685
- 686

687  
688  
689  
690  
691

## APPENDIX A – EVALUATION CRITERIA

# NATIONAL QUALITY FORUM

### Measure Evaluation Criteria December 2009

#### Conditions for Consideration

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

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

#### Criteria for Evaluation

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

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

**1a.** The measure focus addresses:

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

**1b.** Demonstration of quality problems and opportunity for improvement, i.e., data<sup>1</sup> 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:

- an outcome (e.g., morbidity, mortality, function, health-related quality of life) that is relevant to, or

<sup>1</sup> 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.

associated with, a national health goal/priority, the condition, population, and/or care being addressed<sup>2</sup>;

OR

- if an intermediate outcome, process, structure, etc., there is **evidence**<sup>3</sup> 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 process<sup>4</sup>, 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 Efficiency<sup>5</sup> – 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 specified<sup>6</sup> 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)<sup>7</sup>.

<sup>2</sup> 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.

<sup>3</sup> 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.

<sup>4</sup> 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.

<sup>5</sup> 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](#)).

<sup>6</sup> 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),

**2b.** Reliability testing<sup>8</sup> 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 testing<sup>9</sup> 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 evidence<sup>10</sup> 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 focus<sup>11</sup>;

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 transparent<sup>12</sup> (e.g., numerator category computed separately, denominator exclusion category computed separately).

**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<sup>11,13</sup>

---

measurement time window, exclusions, risk adjustment, definitions, data elements, data source and instructions, sampling, scoring/computation.

<sup>7</sup> 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.

<sup>8</sup> 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.

<sup>9</sup> 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.

<sup>10</sup> 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.

<sup>11</sup> Risk factors that influence outcomes should not be specified as exclusions.

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

<sup>13</sup> 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



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 meaningful<sup>14</sup> 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)<sup>15</sup>. 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 harmonized<sup>16</sup> 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 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

---

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.<sup>14</sup> 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.

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

<sup>16</sup> 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.

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 confidentiality<sup>17</sup>, 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).

692

693

---

<sup>17</sup> 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.

694 **Current Evaluation Criteria for Practices**

695 **Specificity.** The practice must be a clearly and precisely defined process or manner of providing  
696 a healthcare service. All candidate safe practices were screened according to this threshold  
697 criterion. Candidate safe practices that met the threshold criterion of specificity were then rated  
698 against four additional criteria relating to the likelihood of the practice improving patient  
699 safety.

700  
701 **Benefit.** If the practice were more widely utilized, it would save lives endangered by healthcare  
702 delivery, reduce disability or other morbidity, or reduce the likelihood of a serious reportable  
703 event (e.g., an effective practice already in near universal use would lead to little new benefit to  
704 patients by being designated a safe practice).

705  
706 **Evidence of Effectiveness.** There must be clear evidence that the practice would be effective in  
707 reducing patient safety events. Such evidence may take various forms, including the following:

- 708 • Research studies showing a direct connection between improved clinical outcomes (e.g.,  
709 reduced mortality or morbidity) and the practice;
- 710 • experiential data (including broad expert agreement, widespread opinion, or professional  
711 consensus) showing the practice is "obviously beneficial" or self-evident (i.e., the practice  
712 absolutely constrains a potential problem or forces an improvement to occur, reduces  
713 reliance on memory, standardizes equipment or process steps, or promotes teamwork); or
- 714 • Research findings or experiential data from non-healthcare industries that should be  
715 substantially transferable to healthcare (e.g., repeat-back of verbal orders or standardizing  
716 abbreviations).

717  
718 **Generalizability.** The safe practice must be able to be utilized in multiple applicable clinical  
719 care settings (e.g., a variety of inpatient and/or outpatient settings) and/or for multiple types of  
720 patients.

721  
722 **Readiness.** The necessary technology and appropriately skilled staff must be available to most  
723 healthcare organizations.

724

725 **APPENDIX B - TASK FORCE MEMBERS**

726 **David Shahian, MD (chair)**

727 Center for Quality and Safety and Department of Surgery,  
728 Massachusetts General Hospital  
729 Professor of Surgery, Harvard Medical School

730

731 **Kristine Martin Anderson, MBA**

732 Senior Vice President, Booz Allen Hamilton, Rockville, MD  
733 Consensus Standards Approval Committee (CSAC) member

734

735 **David Atkins MD, MPH**

736 Director of Quality Enhancement Research Initiative (QUERI),  
737 Department of Veterans Affairs, Health Services Research & Development Service

738

739 **Arthur Levin, MPH**

740 Director, Center for Medical Consumers, New York, NY  
741 Consensus Standards Approval Committee (CSAC) member

742

743 **Mary Naylor, PhD, RN**

744 Marian S. Ware Professor in Gerontology  
745 University of Pennsylvania School of Nursing  
746 Board member

747

748 **Greg Pawlson, MD, MPH**

749 Executive Vice President, National Committee for Quality Assurance (NCQA)

750

751 **Eric Schneider, MD, MSc, FACP**

752 Senior Scientist and Director, RAND Boston  
753 Associate Professor, Division of General Medicine and Primary Care  
754 Brigham and Women's Hospital and  
755 Department of Health Policy and Management  
756 Harvard School of Public Health

757

758 **APPENDIX C - US PREVENTIVE SERVICES TASK FORCE SYSTEM FOR GRADING EVIDENCE AND**  
 759 **RECOMMENDATIONS**

760 The following information was obtained from AHRQ websites describing the [grade definitions](#)  
 761 and [methods](#).

762  
 763 **What the Grades Mean and Suggestions for Practice**

764 The USPSTF updated its definitions of the grades it assigns to recommendations and now includes "suggestions for  
 765 practice" associated with each grade. The USPSTF has also defined levels of certainty regarding net benefit. These  
 766 definitions apply to USPSTF recommendations voted on after May 2007.  
 767

Grade	Definition	Suggestions for Practice
<b>A</b>	The USPSTF recommends the service. There is high certainty that the net benefit is substantial.	Offer or provide this service.
<b>B</b>	The USPSTF recommends the service. There is high certainty that the net benefit is moderate or there is moderate certainty that the net benefit is moderate to substantial.	Offer or provide this service.
<b>C</b>	The USPSTF recommends against routinely providing the service. There may be considerations that support providing the service in an individual patient. There is at least moderate certainty that the net benefit is small.	Offer or provide this service only if other considerations support the offering or providing the service in an individual patient.
<b>D</b>	The USPSTF recommends against the service. There is moderate or high certainty that the service has no net benefit or that the harms outweigh the benefits.	Discourage the use of this service.
<b>I State ment</b>	The USPSTF concludes that the current evidence is insufficient to assess the balance of benefits and harms of the service. Evidence is lacking, of poor quality, or conflicting, and the balance of benefits and harms cannot be determined.	Read the clinical considerations section of USPSTF Recommendation Statement. If the service is offered, patients should understand the uncertainty about the balance of benefits and harms.

768  
 769 **Levels of Certainty Regarding Net Benefit**  
 770

Level of Certainty*	Description
<b>High</b>	The available evidence usually includes consistent results from well-designed, well-conducted studies in representative primary care populations. These studies assess the effects of the preventive service on health outcomes. This conclusion is therefore unlikely to be strongly affected by the results of future studies.
<b>Moderate</b>	The available evidence is sufficient to determine the effects of the preventive service on health outcomes, but confidence in the estimate is constrained by such factors as: <ul style="list-style-type: none"> <li>• The number, size, or quality of individual studies.</li> <li>• Inconsistency of findings across individual studies.</li> <li>• Limited generalizability of findings to routine primary care practice.</li> <li>• Lack of coherence in the chain of evidence.</li> </ul> As more information becomes available, the magnitude or direction of the observed effect could change, and this change may be large enough to alter the conclusion.
<b>Low</b>	The available evidence is insufficient to assess effects on health outcomes. Evidence is insufficient because of: <ul style="list-style-type: none"> <li>• The limited number or size of studies.</li> <li>• Important flaws in study design or methods.</li> <li>• Inconsistency of findings across individual studies.</li> <li>• Gaps in the chain of evidence.</li> <li>• Findings not generalizable to routine primary care practice.</li> <li>• Lack of information on important health outcomes.</li> </ul>

More information may allow estimation of effects on health outcomes.
--

771 \* The USPSTF defines certainty as "likelihood that the USPSTF assessment of the net benefit of a preventive service is correct."  
 772 The net benefit is defined as benefit minus harm of the preventive service as implemented in a general, primary care population.  
 773 The USPSTF assigns a certainty level based on the nature of the overall evidence available to assess the net benefit of a preventive  
 774 service.

775  
 776 **U.S. Preventive Services Task Force Recommendation Grid\***  
 777  
 778

Certainty of Net Benefit	Magnitude of Net Benefit			
	Substantial	Moderate	Small	Zero/Negative
High	A	B	C	D
Moderate	B	B	C	D
Low	Insufficient			

779 \*A, B, C, D, and *Insufficient* represent the letter grades of recommendation or statement of insufficient evidence  
 780 assigned by the U.S. Preventive Services Task Force after assessing certainty and magnitude of net benefit of the  
 781 service.  
 782

783 **U.S. Preventive Services Task Force Terminology to Describe the Critical Assessment of Evidence at 3**  
 784 **Levels: Individual Studies, Key Questions, and Overall Certainty of Net Benefit of the Preventive Service**  
 785

Level of Evidence Assessed	Terminology	Criteria Used to Select Terminology
Individual studies	Good, fair, poor (quality)	Critical appraisal; judgment
Key questions in analytic framework*	Convincing, adequate, inadequate (evidence)	6 questions in <a href="#">Table 2</a> ; judgment
Overall certainty of net benefit of the preventive service	High, moderate, low (certainty)	6 questions in <a href="#">Table 2</a> ; judgment

786 \*This terminology is not reflected in the carotid artery stenosis screening recommendation statement in this issue,<sup>1</sup>  
 787 but it will appear in future recommendation statements.  
 788