

NATIONAL QUALITY FORUM

TO: NQF Members

FR: Karen Pace, PhD, RN

RE: Review of *Guidance for Evaluating the Evidence Related to the Focus of Quality Measurement*

DA: May 20, 2010

Background

NQF's [evaluation criteria](#) require a variety of evidence including the clinical evidence for the focus of a quality measure (criterion 1c). Steering committees have diverse backgrounds and expertise and could benefit from more guidance and support to consistently apply the NQF measure evaluation criteria. Both evidence and expert judgment play a role in evaluating measures against criteria, however, judgment can best be applied when Steering Committees have a thorough understanding of the evidence that does or does not exist.

In January, NQF convened a task force of seven members to assist with developing guidance on evaluating the evidence that supports the focus of a quality performance measure (1c), as well as the other subcriteria under *Importance to Measure and Report*. The task force was asked to address the following tasks.

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

The Task Force's recommendations are included in the draft document, *Guidance for Evaluating the Evidence Related to the Focus of Quality Measurement*. The draft report is posted on the NQF web site for review and comment only – not voting.

You may post your comments and view the comments of others on the NQF website. NQF is now using a program that facilitates electronic submission of comments on this draft report. **All comments must be submitted using the online submission process.**

NQF Member comments must be submitted no later than 6:00 PM ET, June 18, 2010; public comments are due 6:00 PM ET, June 11, 2010.

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

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

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due by June 18, 2010, 6:00 PM ET; PUBLIC comments due by June 11, 2010 by 6:00 PM ET

NATIONAL QUALITY FORUM

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

Draft Report for Review and Comment

May 20, 2010

NATIONAL QUALITY FORUM

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

Draft Report May 20, 2010

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	11
19	III. Recommendations for Evaluating Criterion 1c – Quantity, Quality, Consistency of Body of	
20	Evidence	14
21	Table 4. Evaluation of Quantity, Quality, and Consistency of Body of Evidence for Criterion	
22	1c – evidence for the measure focus.....	16
23	Table 5. Evaluation of Criterion 1c based on the quantity, quality and consistency of the	
24	body of evidence	17
25	IV. Recommendations for Selecting the Focus for Measure Development	17
26	V. Recommendations for Evaluating Importance to Measure and Report and the Other	
27	Subcriteria	18
28	Table 6. Evidence for Evaluating Importance to Measure and Report	20
29	VI. Recommendations for Modifications to the NQF Evaluation Criteria.....	20
30	VII. Recommendations for Modifications to the Measure Submission.....	22
31	VIII. Recommendations for Evidence Required for Practices Considered for NQF	
32	Endorsement.....	25
33	Table 8. Evidence to Support a Practice.....	25
34	Consequences of Measurement.....	26
35	REFERENCES.....	26
36	TASK FORCE MEMBERS.....	29
37	APPENDIX.....	30
38	Current Measure Evaluation Criteria.....	30
39	Current Evaluation Criteria for Practices	35
40	US Preventive Services Task Force System for Grading Evidence and Recommendations.....	36

44 **OVERVIEW AND PURPOSE**

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

55
56 NQF's [evaluation criteria](#) require a variety of evidence as noted in the following table. Of these
57 criteria, some of the most rigorous evidence is required to justify what is being measured (1c)
58 and that is the primary focus of this report – *the evidence required to justify the measure focus*
59 (i.e., the specific process, structure, outcome, etc. that is being measured). Another task force
60 and subsequent report will address measure testing and the criterion of *Scientific Acceptability of*
61 *Measure Properties*.

62
63 NQF endorses measures that are intended for use in public reporting as well as quality
64 improvement with the goal of improving the quality of healthcare. The evidence that supports
65 the focus for a quality measure is addressed under the must-pass criterion, *Importance to*
66 *Measure and Report* because if the measure focus is not supported by evidence that it can
67 facilitate gains in quality and health, then the use of limited resources for measuring and
68 reporting on it would be questionable. For most healthcare quality measures, the evidence will
69 be that of clinical effectiveness and the link to desired outcomes.

70
71
72
73
74
75
76

77 Table 1. Measure Evaluation Criteria and Type of Evidence
78

Evaluation Criteria	Type of Evidence
1. Importance to measure and report 1a . High impact	Epidemiologic data
1b . Opportunity for improvement	Epidemiologic data Health services research
1c . Evidence that supports the focus of measurement	Clinical research; Health services research
2. Scientific acceptability of measure properties 2a-h	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 Understanding what? (ECS)
4. Feasibility 4e . Demonstration the measure can be implemented	Data and/or qualitative information demonstrating the measure can be implemented

79

80 **Task Force Charge**

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

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

89

90

91 **BACKGROUND**

92 Ideally, quality performance measures are based on high quality evidence regarding the types
93 of interventions and services that will achieve desired outcomes and reflect high quality care.
94 However, much of healthcare has not been subjected to research studies, much less with
95 randomized controlled trials or comparative effectiveness studies. Lohr observed that “Perhaps

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

104

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

115

116 **Evidence Issues Identified with Measures Submitted to NQF**

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

127 immunization status rather than measures of administering the vaccine). Some measure
128 submitters question whether the suggested USPSTF evidence grading system is only applicable
129 to preventive services.

130
131 NQF consensus projects were not intended to undertake systematic evidence reviews for the
132 variety of measures that are submitted for consideration, nor is this feasible. The responsibility
133 for such reviews lies with measure developers. However, in the past, measure developers have
134 varied substantially in the expertise and resources they have to systematically assess the
135 strength of a body of evidence or crosswalk a different grading system to the USPSTF system.
136 Such detailed evidence reviews have frequently not been viewed by developers as an integral
137 part of the measure development process. NQF wishes to clearly signal, through this document
138 and the measure submission form itself, that evidence reviews need to be completed by
139 measure developers or their designates prior to measure submission for endorsement.

140

141 **The Changing Environment**

142 As guidelines and quality metrics are increasingly used not only for internal quality
143 improvement but also for public reporting, the necessity for a strong evidence base has become
144 more urgent and compelling. This need is further substantiated by the development of
145 reimbursement programs that utilize such publicly reported metrics. Although public reporting
146 and pay for performance have the potential to inform consumers, focus quality improvement
147 activities, and reward high performance; there are potential unintended negative consequences
148 if measures do not meet all the aspects of the importance criterion. Potential negative
149 consequences include confusion about the importance of particular care processes to quality,
150 the unnecessary resources to measure elements of care that may not impact quality, and
151 diversion of scarce resources to marginally effective activities. To achieve the intended positive
152 effects of quality measurement and minimize the unintended potential negative consequences,
153 measures should be based on the best evidence for the focus of measurement and also should
154 conform to the highest measurement science principles. Recognizing the high stakes of
155 performance measurement in an increasingly transparent environment, some measure
156 developers have enhanced their requirements for the evidence base for performance measure
157 development ⁴.

158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188

Clinical Practice Guidelines

Although they are not the only evidence base for performance measures, many measure developers rely on clinical practice guidelines to support the focus of measurement ^{2,4}. There has been a proliferation of such guidelines, some overlapping or even contradictory. There also is substantial variability in the methodological rigor of review and grading of the evidence and recommendations. In 2000, Grilli ⁵ and colleagues reported that of 431 specialty society guidelines reviewed, 82% did not apply explicit criteria to grade the scientific evidence used as a basis for recommendations, 87% did not report whether a systematic literature search was conducted, and 67% did not describe the professional involved. Some tools to assess clinical practice guidelines ⁶⁻⁸ are available and developing trustworthy guidelines is also the subject of a current IOM study.

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

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

Evidence Grading Systems

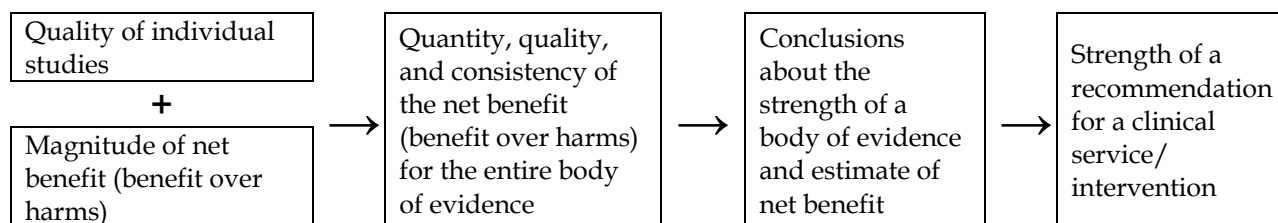
A variety of evidence grading systems currently are in use to achieve this enhanced degree of evidence review and assessment. These systems generally include methods for selection and review of the evidence, and rules or hierarchies related to grading the quality of evidence and the strength of a recommendation. These evidence grading systems are applicable to guidelines as well as other sources of evidence for performance measures.

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

198

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

200



201

202

203 Differences in terminology and grading scales may inhibit understanding about the strength of
204 evidence. Differences can range from a rather minor but understandable difference in
205 terminology (e.g., strength, quality, or level of evidence) to pronounced differences in the
206 assignment of grades (e.g., a grade of A could indicate evidence based on consensus of opinion
207 in one system to evidence based on meta-analyses of randomized controlled trials in another
208 system). An international initiative to standardize grading evidence and recommendations,
209 [GRADE](#)⁹⁻¹⁵, is now supported by many [organizations](#) including the Cochrane Collaboration.
210 The Agency for Healthcare Research and Quality (AHRQ) supports two evidence grading
211 systems: one used by the US Preventive Services Task Force (USPSTF)^{16,17} and one used by the
212 Evidence-Based Practice Centers¹⁸ (consistent with GRADE). Table 2 provides examples of
213 terminology used by four evidence grading systems. It is important to note that grading
214 systems are tied to specific methods for reviewing and assessing the quality of evidence.

215

216

217

218 Table 2. Examples of Terminology in Selected Grading Scales
 219

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

220
 221
 222 Systematic reviews and meta-analyses are used to assess a body of evidence. PRISMA
 223 (Preferred Reporting Items for Systematic Reviews and Meta-Analyses) focuses on the
 224 transparent and full reporting of such reviews¹⁹. The Institute of Medicine (IOM) has two
 225 consensus projects underway that relate to grading the quality of evidence for clinical
 226 interventions: [Standards for Developing Trustworthy Clinical Practice Guidelines](#) and

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

229

230

231 **RECOMMENDATIONS**

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

234

235 **Principles**

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

239

240 **Measures that will be used for public reporting should meet a high standard of evidence for**
241 **the focus of measurement.** NQF measures are intended to be useful for public reporting, as
242 well as to internal quality improvement activities. Measures used for public reporting often
243 impact large numbers of providers and entail investment of significant resources in
244 measurement and improvement. Consequently, measures that will be used for public reporting
245 should meet a high standard of evidence for the focus of measurement. A lower standard of
246 evidence may be deemed appropriate by those selecting measures for use in smaller scale
247 internal quality improvement activities within a learning system that allows for rapid
248 adjustments.

249

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

253 Much of healthcare has not been subjected to research studies and thus, does not have a strong
254 evidence base. In the absence of strong evidence, clinical interventions and services that are the
255 focus of quality performance measures should be judged to have benefits to patients that clearly
256 outweigh any potential risk. In the absence of strong evidence, administrative, management, or
257 system structures and processes that are the focus of quality performance measures should be

258 judged to have benefits to patients that clearly outweigh the system costs and resources to
259 implement those structures and processes.

260
261 **Standards for evidence grading are evolving and expectations for both the present and future**
262 **should be stated.** Standards for evidence review and grading and clinical practice guideline
263 development are evolving, as are expectations for measures endorsed by NQF. Explicit
264 information about the evidence supporting a measure and how (or if) it was graded is essential
265 for evaluating the evidence both now and in the future.

266
267 **Consistency with prior terminology, whenever possible, minimizes confusion.** Terminology
268 used in prior NQF documents should be changed only if incorrect or leads to increased
269 understanding. Whenever possible, narrative descriptions should be used instead of technical
270 terminology.

271
272 **I. Recommendations on Sources of Evidence and Evidence Grading for the Present and the Future**

- 273 • The preferred sources of evidence are systematic reviews and grading of evidence
274 conducted by independent organizations such as USPSTF, AHRQ Evidence-based Practice
275 Centers, and the Cochrane Collaborative; or guidelines that meet national standards for
276 trustworthy guidelines.
- 277 • Until such time when guidelines are certified to meet a set standard, preferred guidelines
278 are those developed with balanced representation beyond one specialty group and with full
279 disclosure of biases.
- 280 • An assigned evidence grade alone is not sufficient to evaluate whether the NQF criterion on
281 evidence for the focus of measurement (1c) is met, either now or in the future. The specific
282 information on the quantity, quality, and consistency of the body of evidence that was used
283 to determine an overall grade should be provided in the measure submission.
- 284 • Explicit, transparent information on the quantity, quality, and consistency of the body of
285 evidence supporting a measure will facilitate identification of guideline recommendations
286 that do not have acceptable evidence as the basis for performance measurement. Explicit
287 information about the evidence also facilitates review by all stakeholders although TAPs

288 and Steering Committees will continue to include experts that possess knowledge about the
289 state of science for a particular topic.

290 • **Current Expectations -**

291 ○ NQF should require measure developers to provide specific information about the
292 quantity, quality, and consistency of evidence. Information should include how the
293 evidence was graded and the grade assigned. If the developer fails to provide this
294 information, NQF should not review the proposed measure.

295 ○ NQF prefers that evidence be graded based on the systems of either the [USPSTF](#) or
296 [GRADE](#).

297 • **Future Expectations -**

298 ○ Rather than identifying “preferred” grading systems as noted for the current
299 expectations, NQF should require the use of one or two standardized evidence
300 grading systems (e.g., the USPSTF, GRADE, or possibly one adopted by the IOM).

301 ○ The evidence should be graded by identified credible sources, such as guideline
302 developers or review organizations certified as meeting accepted standards.

303 ○ Even with standardized grading systems and potentially certified reviewers, explicit
304 information on the quantity, quality, and consistency of the specific evidence that led
305 to the assignment of a grade should be submitted for evaluation. In other words,
306 NQF expects not simply the end-result of the grading process, but also a concise
307 summary of the evidence.

308

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

310 There has been widespread acceptance of Donabedian’s ^{20, 21} structure-process-outcome model
311 for assessing healthcare quality. These three approaches to quality measurement can be used
312 with any topic of healthcare quality and the evidence required generally does not vary by topic.

313 The required evidence is for the links depicted by the red arrows in Figure 2. As depicted under
314 process, there may be multiple process steps prior to delivering an intervention; however, the
315 evidence is most often about the relationship between the intervention and outcome.

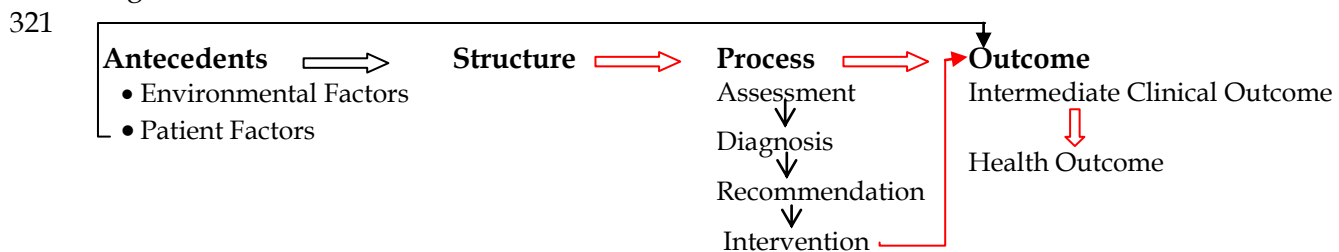
316

317

318

319

320 Figure 2. Structure-Process-Outcome Model



322

323 Table 3 outlines the evidence required to justify the structure, process, or outcome that is the

324 focus of measurement (i.e., what is being measured). It also identifies special considerations

325 related to certain quality topics. Subsequent tables lay out the approach for evaluating the

326 evidence and using it to determine if the NQF criterion is met.

327

328 Outcomes as a representation of quality also are based on the process-outcome link. Outcomes

329 are viewed as useful quality indicators because they are integrative of the influence of multiple

330 care processes and disciplines involved in the care. However, that also presents some challenges

331 related to presenting evidence to support the focus of measurement. Optimally, there will be a

332 body of evidence for the link between the outcome and at least one care process. However, the

333 lack of such evidence should not necessarily be reason to automatically dismiss the value of

334 measurement, particularly when the outcome represents the reason for seeking and providing

335 healthcare (e.g., health, function, survival, symptom control) or harm resulting from healthcare

336 provided or omitted. Once outcomes are measured and reported, many outcomes that were not

337 thought to be modifiable tend to be improved and stimulate identification and adoption of

338 effective practices.

339 Table 3. Evidence to Support the Focus of Measurement

340

Type of Measure	Evidence	Example of Measure Type & Evidence to be Addressed
<p>Structure 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(including evidence for the link to effective care processes and the link from the care processes to desired health outcomes)</p>	<p>#0190 Nurse Staffing Hours Evidence that higher nursing hours are associated with lower mortality, morbidity ; or associated with effective care processes (e.g., lower medication errors) that lead to better outcomes</p>
<p>Process 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</p> <p>If the measure focus is on inappropriate use: Quantity, quality, and consistency for a body of evidence that the measured healthcare process does <u>not</u> lead to desired health outcomes in the target population</p>	<p>#0551 ACE Inhibitor / Angiotensin Receptor Blocker(ARB) Use and Persistence Among Members with Coronary Artery Disease at High Risk for Coronary Events Evidence that use of ACE-I and ARB are associated with lower mortality and/or cardiac events</p> <p>#0058 Inappropriate antibiotic treatment for adults with acute bronchitis Evidence that antibiotics are not effective for acute bronchitis</p>
<p>Intermediate Clinical Outcome 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 the measured intermediate clinical outcome leads to desired health outcomes in the target population</p>	<p>#0059 Hemoglobin A1c Management Evidence that hemoglobin A1c > 9 is associated with more complications</p>
<p>Health Outcome 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 measures may be considered proxies for a health state (e.g., hospitalization may represent a deterioration in health status)</p>	<p>Optimally, quantity, quality, and consistency for a body of evidence that the measured outcome (desirable or adverse) is influenced by at least one healthcare process or service. However, outcomes do not necessarily require evidence.</p>	<p>#0230 Acute Myocardial Infarction 30-day Mortality Survival is a goal of seeking and providing treatment for AMI Evidence that healthcare processes/ interventions (aspirin, reperfusion) affect mortality/ survival</p> <p>#0171 Acute care hospitalization (risk-adjusted) [of home care patients] Improvement or stabilization of condition to remain at home is a goal of seeking and providing home care services. Evidence that healthcare processes (e.g., medication reconciliation, care coordination) affect hospitalization of patients receiving home care services</p> <p>#0140 Ventilator-associated pneumonia for ICU and high-risk nursery (HRN) patients Avoiding harm from treatment is a goal of when seeking and providing healthcare. Evidence that ventilator acquired pneumonia is affected by healthcare processes (e.g., ventilator</p>

Type of Measure	Evidence	Example of Measure Type & Evidence to be Addressed
		bundle)
Special Considerations by Topic		
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	#0166 HCAHPS Evidence 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)
Efficiency Measures of efficiency combine the concepts of resource use <u>and</u> quality	Efficiency Measured with combination of Quality measures and Resource Use measures Quality measure component Evidence for the selected quality measure(s) as described in this table Resource use measure component 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 Evidence for diabetes quality measure(s) as described in this table

341

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

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

350

- 351 • Evidence should be evaluated on *quantity, quality* of studies, *consistency* in direction, and
352 magnitude of net benefit (benefits over harms) of a ***body of evidence*** on a scale of High,
353 Moderate, or Low.
- 354 • The dimensions of *quantity, quality, and consistency* for a body of evidence apply to measures
355 based on guidelines as well as those for which guidelines may not exist (e.g., care
356 coordination or team functioning may not be based on guidelines, but often have bodies of
357 evidence including non-clinical literature that should be systematically assessed)

- 358 • Measures without a clear description of the *quantity, quality, and consistency* of the
359 supporting body of evidence or without any evidence should not pass criterion 1c and the
360 threshold criterion of *Importance to Measure and Report*.
- 361 • Use of only selected individual studies from a body of evidence is not adequate to evaluate
362 the evidence and should not pass criterion 1c and the threshold criterion of *Importance to*
363 *Measure and Report*. This should be flagged in the measure submission form.
- 364 • Inconsistent and conflicting evidence should result in measures not passing both criterion 1c
365 and the threshold criterion of *Importance to measure and report*.
- 366 • Expert opinion is acceptable evidence; it should be systematically assessed and fully
367 described and will be evaluated as outlined in Table 4.

368

369 Table 4 provides guidance on how to evaluate each of the dimensions of *quantity, quality, and*
370 *consistency* for a body of evidence. Table 5 provides recommended decision rules for using the
371 ratings for all three dimensions to make a decision on whether a measure should pass the
372 criterion 1c, the evidence to support the measure focus. High quality evidence usually requires
373 multiple studies each with sufficient numbers of patients to give precise estimates, but
374 occasionally a large and representative study can give high quality evidence. For example, one
375 study (low quantity) that is a RCT with a large representative sample of patients (high quality)
376 and substantial estimates of net benefit would pass the criterion, whereas, a body of evidence
377 with low consistency of estimates of net benefits indicates a measure should not pass the
378 criterion regardless of the ratings for quantity and quality of studies.

379
380

381 Table 4. Evaluation of Quantity, Quality, and Consistency of Body of Evidence for Criterion 1c –
 382 evidence for the measure focus
 383

	Quantity of Body of Evidence	Quality of Body of Evidence	Consistency of Results of Body of Evidence
Definition	Total number of studies (not articles or papers)	Certainty or confidence in the estimates of benefits and harms to patients across studies in the body of evidence resulting from study factors* including: study design or flaws; directness/indirectness regarding the specific process or structure being measured, outcomes assessed, target population, comparisons; imprecision (wide confidence intervals due to few patients or events)	Stability in both the magnitude and direction of benefits and harms to patients (benefits over harms) across studies in the body of evidence
High	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 benefits and harms to patients are consistent in direction and similar in magnitude across studies in the body of evidence
Moderate	2-4 studies**	<ul style="list-style-type: none"> • Non-RCTs with control for confounders that could account for other plausible explanations, with large, precise estimate of effect; or • RCTs without serious flaws that introduce bias, but with either indirect evidence, or imprecise estimate of effect 	Estimates of benefits and harms to patients are consistent in direction, but differ in magnitude across studies in the body of evidence; or If only one study, the estimate of benefits greatly outweighs the estimate of potential harms OR For expert opinion that is systematically assessed, agreement that benefits to patients clearly outweigh potential harms
Low	0-1 studies**	<ul style="list-style-type: none"> • Expert opinion that is systematically assessed; or • RCTs with flaws that introduce bias; or • Non-RCTs with small or imprecise estimate of effect, or without control for confounders that could account for other plausible explanations 	Differences in both magnitude and direction of benefits and harms to patients across studies in the body of evidence; or wide confidence intervals prevent estimating net benefit OR For expert opinion evidence that is systematically assessed: <u>lack of agreement</u> that benefits to patients clearly outweigh potential harms
Inadequate to Evaluate	No empirical evidence; OR only selected individual studies from a larger body of evidence	Expert opinion only and it was not systematically assessed	No assessment of magnitude and direction of benefits and harms to patients

384 *Study designs that affect certainty of confidence in estimates of effect include: Randomized controlled
 385 trials (RCT), which control for both observed and unobserved confounders, and non-RCTs (observational
 386 studies) with various levels of control for confounders.
 387 Study flaws that may bias estimates of effect include: lack of allocation concealment; lack of blinding;
 388 large losses to follow-up; failure to adhere to intention to treat analysis; stopping early for benefit; failure
 389 to report important outcomes.
 390 Imprecision with wide confidence intervals around estimates of effects can occur in studies involving few
 391 patients and few events.
 392 Indirectness of evidence includes: indirect comparisons (e.g., two drugs compared to placebos rather than
 393 head-to head), differences between the population, intervention, outcome of interest, or comparator
 394 interventions and those included in the relevant studies.¹⁴
 395 ** The suggested number of studies for rating levels of quantity is considered a general guideline.
 396

397 Table 5. Evaluation of Criterion 1c based on the quantity, quality and consistency of the body of
 398 evidence
 399

Quantity of Body of Evidence	Quality of Body of Evidence	Consistency of Body of Evidence	Pass Criterion 1c
Moderate-High	Moderate-High	Moderate-High	Yes
Low	Moderate-High	Moderate-High	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	Low	No
Low	Low	Low	No

400

401 IV. Recommendations for Selecting the Focus for Measure Development

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

404

- 405 • For any topic area, measures based on the best evidence should be considered over
 406 measures based on lower quality evidence (e.g., expert opinion).
- 407 • There is a hierarchical preference for outcome measures (when possible) followed by
 408 process measures. Structural measures are appropriate primarily when there are very well
 409 established structure-process-outcome relationships; and when it is not feasible to directly
 410 measure the outcomes or processes. For process and structure measures, the focus of
 411 measurement should be on the aspect of care with the most direct evidence of a strong

412 relationship to the desired outcome. For example, evidence about effective medication to
413 control blood pressure is direct evidence for the medication but only indirect evidence for
414 the frequency of assessing blood pressure (see Figure 2). Assessment of blood pressure,
415 although necessary, is not sufficient to achieving control. When there are multiple processes
416 that affect a desired outcome, efforts should be made to include measures for all processes
417 that have a strong relationship to the desired outcome.

418

419 **V. Recommendations for Evaluating Importance to Measure and Report and the Other Subcriteria**

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

425

426 Generally, in measure submissions, high impact is easily demonstrated by alignment with a
427 specific NPP goal or epidemiologic data (incidence, prevalence, resource use, consequences of
428 quality problems). However, data on opportunity for improvement may be lacking (e.g.,
429 submitter states that performance is unknown), or it may not be specific to the focus of
430 measurement, or only based on a sample from measure development and testing. When data
431 are presented, reviewers sometimes question whether there is enough variation to justify
432 importance to measure and report or how to judge overall poor performance. When a measure
433 undergoes review for continued endorsement, one issue that sometimes arises is whether a
434 measure is “topped out” meaning there are high levels of performance with little variation and
435 therefore, little room for further improvement. The Task Force did not recommend specific
436 quantitative thresholds for identifying conformance with the subcriteria of high impact (1a) and
437 opportunity for improvement (1b).

438

- 439 • Threshold values for opportunity for improvement would be difficult to standardize. It
440 depends on the size of the population at risk, effectiveness, and the consequences of the
441 quality problem. For example, even modest variation would be sufficient justification for

442 some highly effective, potentially life-saving treatments (e.g., certain vaccinations) that are
443 critical to the public health.

- 444 • At the time of review for continued endorsement, being “topped out” is not a reason in itself
445 to remove endorsement for a measure; however, it may be a signal of some other problem
446 with the measure (e.g. imprecise specification, overly broad exclusions). If a measure is an
447 important and valid indicator of quality, it may still be justified to retain endorsement, as
448 overall performance could deteriorate if not monitored. However, a “topped out” process
449 measure might have endorsement withdrawn if there is an associated outcome measure.

450

451

452 Table 6. Evidence for Evaluating Importance to Measure and Report
 453

Subcriterion	Evidence	Example	Pass the subcriterion?
High impact (1a)	Addresses a specific national health goal/priority identified by the Secretary of DHHS or the NPP OR Epidemiologic data – 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 NPP goal: . . . focus relentlessly on continually reducing and seeking to eliminate all healthcare-associated infections (HAIs) Evidence 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)	Subcriterion 1a Yes – Demonstrated at least one of the aspects of high impact No – Did not demonstrate at least one of the aspects of high impact
Opportunity for improvement (1b)	Epidemiologic data; health services research – data demonstrating considerable variation, or overall poor performance, in the quality of care across providers and/or population groups (disparities in care)	#0432 Influenza Vaccination of Nursing Home/ Skilled Nursing Facility Residents NPP goal: All Americans will receive the most effective preventive services recommended by the U.S. Preventive Services Task Force Evidence that vaccination rates vary (e.g., 39% fail to reach the Healthy People 2010 objective of vaccinating at least 90% of nursing home residents)	Subcriterion 1b Yes – Demonstrated either variation or overall less than optimal performance No – Did not demonstrate either variation or overall less than optimal performance
Evidence for the focus of measurement (1c)	See Table 3	See Table 3	Subcriterion 1c See Table 4 and Table 5
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>			

454
 455 **VI. Recommendations for Modifications to the NQF Evaluation Criteria**
 456 As noted previously, the Task Force recommended that all three subcriteria be met to pass the
 457 threshold criterion of *Importance to Measure and Report*. The following redlined modifications to
 458 the criteria are based on the Task Force recommendations as reported above, as well as a few
 459 editorial changes.
 460
 461

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

468 |
469 | **1a.** The measure focus addresses:

- 470 | • a specific national health goal/priority identified by NQF's DHHS or the National
471 | Priorities Partnership convened by NQF;
- 472 | **OR**
- 473 | • a demonstrated high impact aspect of healthcare (e.g., affects large numbers of patients
474 | and/or has a substantial impact for a smaller population; leading cause of
475 | morbidity/mortality; high resource use (current and/or future); severity of illness; and
476 | severity of patient/societal consequences of poor quality).

477 | **AND**

478 | **1b.** Demonstration of quality problems and opportunity for improvement, i.e., data (footnote 1)
479 | demonstrating considerable variation, or overall poor performance, in the quality of care across
480 | providers and/or population groups (disparities in care).

481 | **AND**

482 | **1c.** The measure focus is evidence-based as demonstrated by a systematic assessment of the
483 | quantity, quality, and consistency of the body of evidence (see Tables 3-5) and standardized
484 | grading of the body of evidence (footnote 3).

- 485 | • ~~an~~ Health outcome (footnote 2): optimally, evidence that the measured outcome (desirable or adverse)
486 | is influenced by at least one healthcare process or service. However, outcomes do not necessarily
487 | require evidence. (e.g., morbidity, mortality, function, health related quality of life) that is relevant to,
488 | or associated with, a national health goal/priority, the condition, population, and/or care being
489 | addressed (2);
- 490 | **OR**
- 491 | • ~~if an intermediate outcome, process, structure, etc., there is~~ **evidence (3)** that supports the specific
492 | ~~measure focus as follows:~~
- 493 | • Intermediate clinical outcome: evidence that the measured intermediate outcome (e.g., blood pressure,
494 | Hba1c) leads to desired outcomes in the target population-improved health/avoidance of harm or
495 | cost/benefit.
- 496 | • Process (footnote 4): evidence that the measured clinical or administrative healthcare process leads to
497 | desired outcomes in the target population-improved health/avoidance of harm and
498 | if the measure focus is on one step in a multi-step care process (4), it measures the step that has the
499 | greatest effect on improving the specified desired outcome(s).
- 500 | • Structure: evidence that the measured structure leads to desired health outcomes (including evidence
501 | for the link to effective care processes and the link from the care processes to desired health outcomes)
502 | ~~supports the consistent delivery of effective processes or access that lead to improved~~
503 | health/avoidance of harm or cost/benefit.
- 504 | • Special Considerations by Topic of Measurement
 - 505 | Patient experience with care: – evidence that the measured aspects of care are those valued
506 | by patients and for which the patient is the best and/or only source of information~~an~~
507 | association exists between the measure of patient experience of health care and the outcomes,
508 | values and preferences of individuals/ the public.
 - 509 | Access—evidence that an association exists between access to a health service and the outcomes
510 | of, or experience with, care.

○—Efficiency ([footnote 5](#)): ~~evidence for the quality component as noted above 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.~~

Footnotes

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, ~~“never events”~~ serious reportable events that are compared to zero are appropriate outcomes for public reporting and quality improvement.

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

⁴ Clinical care processes typically include multiple steps: assess → identify problem/potential problem → choose/plan intervention (with patient input) → provide intervention → evaluate impact on health status. If the measure focus is one step in such a multi-step process, the step with the greatest effect on strongest evidence for the link to the desired outcome should be selected as the focus of measurement. ~~For example, although assessment of immunization status and recommending immunization are necessary steps, they are not sufficient to achieve the desired impact on health status—patients must be vaccinated to achieve immunity. This does not preclude consideration of measures of preventive screening interventions where there is a strong link with desired outcomes (e.g., mammography) or measures for multiple care processes that affect a single outcome.~~

⁵ Measures of efficiency combine the concepts of resource use and quality. 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](#)).

VII. Recommendations for Modifications to the Measure Submission

The information requested on NQF’s measure submission form is consistent with those identified in a 2009 collaborative effort undertaken with CMS, The Joint Commission, NCQA, and PCPI to identify common data fields. AHRQ participated, but maintained its own data items for the [National Quality Measures Clearinghouse](#).

The Task Force suggested the following modifications to the information requested on the NQF [measure submission form](#). The intent is full transparency about the supporting evidence for the submitted measure. This will facilitate understanding of the adequacy of the evidence presented (selected evidence vs. a body of evidence) and the developer’s representation of the quality of the evidence. Currently, evidence graded using the USPSTF or GRADE systems may not be

557 available, however, an accurate description of the evidence and any grading system used
 558 should still be expected. The following items pertain to the recommendations related to
 559 evidence (1c) under *Importance to Measure and Report*.

560

561 Table 7. Current and Proposed Information Requested on Measure Submission

562

Current Measure Submission (4.1) Items	Proposed Measure Submission Items
	<p>Introduction <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>1c.1. Relationship to Outcomes (For non-outcome measures, briefly describe the relationship to desired outcome. For outcomes, describe why it is relevant to the target population.)</p> <p>1c.2. Type of Evidence (Check all that apply) Cohort study Observational study Evidence-based guideline Randomized controlled trial Expert opinion Systematic synthesis of research Meta-analysis Other:</p> <p>1c.4. Summary of Evidence (For non-outcome measures, provide evidence of relationship to desired outcome. For outcomes, summarize any evidence that healthcare services/care processes influence the outcome.)</p> <p>1c.5. Rating of Strength/Quality of Evidence (Also provide narrative description of the rating and by whom)</p> <p>1c.6. Method for Rating Evidence</p> <p>1c.7. Summary of Controversy/Contradictory Evidence</p> <p>1c.8. Citations for Evidence (Other than guidelines)</p> <p>1c.9. Quote the Specific Guideline Recommendation (Including guideline number and/or page number)</p>	<p>1c.1. Structure-Process-Outcome Relationship (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.)</p> <p>1c.2. Source of Evidence Clinical practice guideline Systematic review of body of evidence (other than within guideline development) Selected individual studies (rather than entire body of evidence) Other</p> <p>1c.4. Summary of Body of Evidence Quantity of Studies in Body of Evidence (total number of studies, not articles): Quality of Body of Evidence (Certainty or confidence in the estimates of benefits and harms to patients across studies 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): Directness to focus of measurement & target population in proposed measure: Consistency of Results across Studies: Net Benefit (Benefits over harms) Benefit/outcome – estimate of effect Harms addressed – estimate of effect</p> <p>1c.5. Grading of Strength/Quality of Body of Evidence Has the body of evidence been graded? Yes No If graded: By whom (describe the entity that graded the evidence,</p>

Current Measure Submission (4.1) Items	Proposed Measure Submission Items
<p>1c.10. Clinical Practice Guideline Citation</p> <p>1c.11. National Guideline Clearinghouse or Other URL</p> <p>1c.12. Rating Strength of Recommendation <i>(Also provide narrative description of the rating and by whom)</i></p> <p>1c.13. Method for Rating Strength of Recommendation <i>(If different from USPSTF system, also describe rating and how it relates to USPSTF)</i></p> <p>1c.14. Rationale for Using This Guideline Over Others</p>	<p><i>including balance of representation and any disclosures regarding bias)</i></p> <p>Grade Assigned to the Evidence: Based on the NQF descriptions for rating the evidence, what was your assessment of the body of evidence (rate each as High, Moderate, or Low)</p> <p>Quantity: Quality: Consistency:</p> <p>1c.6. System for Grading Evidence: USPSTF GRADE Other <i>(provide description of grading scale with definitions)</i></p> <p>1c.7. Summary of Controversy/Contradictory Evidence</p> <p>1c.8. Citations for Evidence <i>(Other than guidelines)</i></p> <p>1c.9. Quote Verbatim the Specific Guideline Recommendation <i>(Including guideline number and/or page number)</i></p> <p>1c.10. Clinical Practice Guideline Citation</p> <p>1c.11. National Guideline Clearinghouse or Other URL for the cited guideline</p> <p>1c.12. Grading of Strength of Recommendation Has the recommendation 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>1c.13. System for Grading Strength of Recommendation: USPSTF GRADE Other <i>(provide description of grading scale with definitions)</i></p> <p>1c.14. Rationale for Using This Guideline Over Others</p>
<p>Descriptive Information</p> <p>De.4. National Priority Partnership priority area <i>(Select the most relevant)</i></p> <ul style="list-style-type: none"> Patient and family engagement Population health Safety Care coordination Palliative and end of life care Overuse <p>De.5. IOM Quality Domain <i>(Select the most relevant)</i></p> <ul style="list-style-type: none"> Effectiveness Efficiency 	<p>Descriptive Information - no change</p> <p>De.4. National Priority Partnership priority area <i>(Select the most relevant)</i></p> <ul style="list-style-type: none"> Patient and family engagement Population health Safety Care coordination Palliative and end of life care Overuse <p>De.5. IOM Quality Domain <i>(Select the most relevant)</i></p> <ul style="list-style-type: none"> Effectiveness Efficiency

Current Measure Submission (4.1) Items	Proposed Measure Submission Items
Equity Patient-centered Safety Timeliness De.6. Consumer Care Need (Select the most relevant) Getting better Living with illness Staying healthy	Equity Patient-centered Safety Timeliness De.6. Consumer Care Need (Select the most relevant) Getting better Living with illness Staying healthy

563
564 **VIII. Recommendations for Evidence Required for Practices Considered for NQF Endorsement**

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

567
568 The Task Force recommends that the same evidence requirements as indicated for process
569 measures (Tables 3, 4, 5) be applied to practices considered for NQF endorsement.

570
571 Table 8. Evidence to Support a Practice
572

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	Safe Practice 16 Safe Adoption of Computerized Prescriber Order Entry Evidence that computerized order entry systems are associated with lower medication errors and adverse events

573
574 **Modifications to Practice Evaluation Criteria**

575 **Evidence of Effectiveness.** ~~A practice is evidence-based as demonstrated by a systematic~~
576 ~~assessment of the quantity, quality, and consistency of the body of evidence (see Tables 3-5) and~~
577 ~~standardized grading of the body of evidence (footnote). There must be clear evidence that the~~
578 ~~practice would be effective in reducing patient safety events. Such evidence may take various~~
579 ~~forms, including the following:~~
580 ~~o Research studies showing a direct connection between improved clinical outcomes (e.g.,~~
581 ~~reduced mortality or morbidity) and the practice;~~
582 ~~o experiential data (including broad expert agreement, widespread opinion, or professional~~
583 ~~consensus) showing the practice is "obviously beneficial" or self-evident (i.e., the practice~~
584 ~~absolutely constrains a potential problem or forces an improvement to occur, reduces reliance~~
585 ~~on memory, standardizes equipment or process steps, or promotes teamwork); or~~
586 ~~o Research findings or experiential data.~~ Evidence from non-healthcare industries that should be
587 substantially transferable to healthcare (e.g., [safety practices of](#) repeat-back of verbal orders or
588 standardizing abbreviations) also may be considered.

589 **Footnote:**
590 The preferred systems for grading the evidence are the USPSTF ([grading definitions and methods](#)) or [GRADE](#).
591

592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620

Consequences of Measurement

Consequences of measurement are not the same as the consequences of the measured structure or process, i.e., the benefits or harms to the patient related to the specific topic of measurement. Currently, unintended consequences of measurement are addressed under feasibility.

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

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

REFERENCES

1. Lohr KN. Rating the strength of scientific evidence: relevance for quality improvement programs. *Int J Qual Health Care*. 2004;16(1):9-18.
2. Spertus JA, Eagle KA, Krumholz HM et al. American College of Cardiology and American Heart Association methodology for the selection and creation of performance measures for quantifying the quality of cardiovascular care. *Circulation*. 2005;111(13):1703-1712.
3. Tricoci P, Allen JM, Kramer JM et al. Scientific evidence underlying the ACC/AHA clinical practice guidelines. *JAMA*. 2009;301(8):831-841.
4. Physician Consortium for Performance Improvement. Physician Consortium for Performance Improvement® (PCPI) Position Statement - The Evidence Base Required for Measures Development. *American Medical Association* 6-26-2009;1-18. Last accessed March 2010.
5. Grilli R, Magrini N, Penna A et al. Practice guidelines developed by specialty societies: the need for a critical appraisal. *Lancet*. 2000;355(9198):103-106.

- 621 6. Shiffman RN, Shekelle P, Overhage JM et al. Standardized reporting of clinical practice
622 guidelines: a proposal from the Conference on Guideline Standardization. *Ann Intern Med.*
623 2003;139(6):493-498.
- 624 7. The AGREE Collaboration. *Appraisal of Guidelines for Research and Evaluation AGREE*
625 *Instrument*. 2001. Available at <http://www.agreecollaboration.org/instrument/>. Last
626 accessed February 2010.
- 627 8. The AGREE Collaboration. Development and validation of an international appraisal
628 instrument for assessing the quality of clinical practice guidelines: the AGREE project.
629 *Quality & safety in health care*. 2003;12(1):18-23.
- 630 9. Atkins D, Eccles M, Flottorp S et al. Systems for grading the quality of evidence and the
631 strength of recommendations I: critical appraisal of existing approaches The GRADE
632 Working Group. *BMC Health Serv Res*. 2004;4(1):38.
- 633 10. Atkins D, Best D, Briss PA et al. Grading quality of evidence and strength of
634 recommendations. *BMJ*. 2004;328(7454):1490-1494.
- 635 11. Guyatt GH, Oxman AD, Vist GE et al. GRADE: an emerging consensus on rating quality of
636 evidence and strength of recommendations. *BMJ*. 2008;336(7650):924-926.
- 637 12. Guyatt GH, Oxman AD, Kunz R et al. Incorporating considerations of resources use into
638 grading recommendations. *BMJ*. 2008;336(7654):1170-1173.
- 639 13. Guyatt GH, Oxman AD, Kunz R et al. Going from evidence to recommendations. *BMJ*.
640 2008;336(7652):1049-1051.
- 641 14. Guyatt GH, Oxman AD, Kunz R et al. What is "quality of evidence" and why is it important
642 to clinicians? *BMJ*. 2008;336(7651):995-998.
- 643 15. Guyatt GH, Oxman AD, Vist GE et al. GRADE: an emerging consensus on rating quality of
644 evidence and strength of recommendations. *BMJ*. 2008;336(7650):924-926.
- 645 16. Harris RP, Helfand M, Woolf SH et al. Current methods of the US Preventive Services Task
646 Force: a review of the process. *Am J Prev Med*. 2001;20(3 Suppl):21-35.
- 647 17. Sawaya GF, Guirguis-Blake J, LeFevre M et al. Update on the methods of the U.S. Preventive
648 Services Task Force: estimating certainty and magnitude of net benefit. *Ann Intern Med.*
649 2007;147(12):871-875.

650 18. Owens DK, Lohr KN, Atkins D et al. Grading the strength of a body of evidence when
651 comparing medical interventions-Agency for Healthcare Research and Quality and the
652 Effective Health Care Program. *J Clin Epidemiol.* 2009.

653 19. Liberati A, Altman DG, Tetzlaff J et al. The PRISMA statement for reporting systematic
654 reviews and meta-analyses of studies that evaluate health care interventions: explanation
655 and elaboration. *Ann Intern Med.* 2009;151(4):W65-W94.

656 20. Donabedian A. *An Introduction to Quality Assurance in Health Care.* New York, NY: Oxford
657 University Press; 2003.

658 21. Donabedian A. The role of outcomes in quality assessment and assurance. *Quality Review*
659 *Bulletin.* 1992;18(11):356-360.

660

661

662

663

664 **TASK FORCE MEMBERS**

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

666 Center for Quality and Safety and Department of Surgery,
667 Massachusetts General Hospital
668 Professor of Surgery, Harvard Medical School

669

670 **Kristine Martin Anderson, MBA**

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

673

674 **David Atkins MD, MPH**

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

677

678 **Arthur Levin, MPH**

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

681

682 **Mary Naylor, PhD, RN**

683 Marian S. Ware Professor in Gerontology
684 University of Pennsylvania School of Nursing
685 Board member

686

687 **Greg Pawlson, MD, MPH**

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

689

690 **Eric Schneider, MD, MSc, FACP**

691 Senior Scientist and Director, RAND Boston
692 Associate Professor, Division of General Medicine and Primary Care
693 Brigham and Women's Hospital and
694 Department of Health Policy and Management
695 Harvard School of Public Health

696

NATIONAL QUALITY FORUM

Current Measure Evaluation Criteria

December 2009

Conditions for Consideration

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

- A. The measure is in the public domain or an intellectual property agreement is signed.
- B. The measure owner/steward verifies there is an identified responsible entity and process to maintain and update the measure on a schedule that is commensurate with the rate of clinical innovation, but at least every 3 years.
- C. The intended use of the measure includes both public reporting and quality improvement.
- D. The requested measure submission information is complete. Generally, measures should be fully developed and tested so that all the evaluation criteria have been addressed and information needed to evaluate the measure is provided. Measures that have not been tested are only potentially eligible for a time-limited endorsement and in that case, measure owners must verify that testing will be completed within 24-12 months of endorsement.

Criteria for Evaluation

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

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

1a. The measure focus addresses:

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

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

1c. The measure focus is:

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

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.

associated with, a national health goal/priority, the condition, population, and/or care being addressed²;

OR

- if an intermediate outcome, process, structure, etc., there is **evidence**³ that supports the specific measure focus as follows:
 - o Intermediate outcome – evidence that the measured intermediate outcome (e.g., blood pressure, HbA1c) leads to improved health/avoidance of harm or cost/benefit.
 - o Process – evidence that the measured clinical or administrative process leads to improved health/avoidance of harm and
if the measure focus is on one step in a multi-step care process⁴, it measures the step that has the greatest effect on improving the specified desired outcome(s).
 - o Structure – evidence that the measured structure supports the consistent delivery of effective processes or access that lead to improved health/avoidance of harm or cost/benefit.
 - o Patient experience – evidence that an association exists between the measure of patient experience of health care and the outcomes, values and preferences of individuals/ the public.
 - o Access – evidence that an association exists between access to a health service and the outcomes of, or experience with, care.
 - o Efficiency⁵ – demonstration of an association between the measured resource use and level of performance with respect to one or more of the other five IOM aims of quality.

If not important to measure and report, STOP.

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

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

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

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

⁴ Clinical care processes typically include multiple steps: assess → identify problem/potential problem → choose/plan intervention (with patient input) → provide intervention → evaluate impact on health status. If the measure focus is one step in such a multi-step process, the step with the greatest effect on the desired outcome should be selected as the focus of measurement. For example, although assessment of immunization status and recommending immunization are necessary steps, they are not sufficient to achieve the desired impact on health status – patients must be vaccinated to achieve immunity. This does not preclude consideration of measures of preventive screening interventions where there is a strong link with desired outcomes (e.g., mammography) or measures for multiple care processes that affect a single outcome.

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

⁶ Measure specifications include the target population (e.g., denominator) to whom the measure applies, identification of those from the target population who achieved the specific measure focus (e.g., numerator),

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

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

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

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

AND

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

AND

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

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

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

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

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

⁸ Examples of reliability testing include, but are not limited to: inter-rater/abstractor or intra-rater/abstractor studies; internal consistency for multi-item scales; test-retest for survey items. Reliability testing may address the data items or final measure score.

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

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

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

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

¹³ Risk models should not obscure disparities in care for populations by including factors that are associated with differences/inequalities in care such as race, socioeconomic status, gender (e.g., poorer treatment outcomes of

OR

- rationale/data support no risk adjustment.

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

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

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

OR

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

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

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

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

3c. Review of existing endorsed measures and measure sets demonstrates that the measure provides a 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.¹⁴ With large enough sample sizes, small differences that are statistically significant may or may not be practically or clinically meaningful. The substantive question may be, for example, whether a statistically significant difference of one percentage point in the percentage of patients who received smoking cessation counseling (e.g., 74% v. 75%) is clinically meaningful; or whether a statistically significant difference of \$25 in cost for an episode of care (e.g., \$5,000 v. \$5,025) is practically meaningful. Measures with overall poor performance may not demonstrate much variability across providers.

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

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

byproduct of care processes during care delivery.

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

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

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

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

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

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

702
703

¹⁷ All data collection must conform to laws regarding protected health information. Patient confidentiality is of particular concern with measures based on patient surveys and when there are small numbers of patients.

704 **Current Evaluation Criteria for Practices**

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

710
711 **Benefit.** If the practice were more widely utilized, it would save lives endangered by healthcare
712 delivery, reduce disability or other morbidity, or reduce the likelihood of a serious reportable
713 event (e.g., an effective practice already in near universal use would lead to little new benefit to
714 patients by being designated a safe practice).

715
716 **Evidence of Effectiveness.** There must be clear evidence that the practice would be effective in
717 reducing patient safety events. Such evidence may take various forms, including the following:
718 o Research studies showing a direct connection between improved clinical outcomes (e.g.,
719 reduced mortality or morbidity) and the practice;
720 o experiential data (including broad expert agreement, widespread opinion, or professional
721 consensus) showing the practice is "obviously beneficial" or self-evident (i.e., the practice
722 absolutely constrains a potential problem or forces an improvement to occur, reduces reliance
723 on memory, standardizes equipment or process steps, or promotes teamwork); or
724 o Research findings or experiential data from non-healthcare industries that should be
725 substantially transferable to healthcare (e.g., repeat-back of verbal orders or standardizing
726 abbreviations).

727
728 **Generalizability.** The safe practice must be able to be utilized in multiple applicable clinical
729 care settings (e.g., a variety of inpatient and/or outpatient settings) and/or for multiple types of
730 patients.

731
732 **Readiness.** The necessary technology and appropriately skilled staff must be available to most
733 healthcare organizations.

734

735 **US Preventive Services Task Force System for Grading Evidence and Recommendations**

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

738
 739 **What the Grades Mean and Suggestions for Practice**

740 The USPSTF updated its definitions of the grades it assigns to recommendations and now includes "suggestions for
 741 practice" associated with each grade. The USPSTF has also defined levels of certainty regarding net benefit. These
 742 definitions apply to USPSTF recommendations voted on after May 2007.
 743

Grade	Definition	Suggestions for Practice
A	The USPSTF recommends the service. There is high certainty that the net benefit is substantial.	Offer or provide this service.
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.
C	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.
D	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.
I State ment	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.

744
 745 **Levels of Certainty Regarding Net Benefit**
 746

Level of Certainty*	Description
High	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.
Moderate	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"> • The number, size, or quality of individual studies. • Inconsistency of findings across individual studies. • Limited generalizability of findings to routine primary care practice. • Lack of coherence in the chain of evidence. 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.
Low	The available evidence is insufficient to assess effects on health outcomes. Evidence is insufficient because of: <ul style="list-style-type: none"> • The limited number or size of studies. • Important flaws in study design or methods. • Inconsistency of findings across individual studies. • Gaps in the chain of evidence. • Findings not generalizable to routine primary care practice. • Lack of information on important health outcomes. More information may allow estimation of effects on health outcomes.

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

751
 752
 753
 754

U.S. Preventive Services Task Force Recommendation Grid*

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			

755 *A, B, C, D, and *Insufficient* represent the letter grades of recommendation or statement of insufficient evidence
 756 assigned by the U.S. Preventive Services Task Force after assessing certainty and magnitude of net benefit of the
 757 service.
 758

759 **U.S. Preventive Services Task Force Terminology to Describe the Critical Assessment of Evidence at 3**
 760 **Levels: Individual Studies, Key Questions, and Overall Certainty of Net Benefit of the Preventive Service**
 761

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 Table 2 ; judgment
Overall certainty of net benefit of the preventive service	High, moderate, low (certainty)	6 questions in Table 2 ; judgment

762 *This terminology is not reflected in the carotid artery stenosis screening recommendation statement in this issue,¹
 763 but it will appear in future recommendation statements.
 764