

NATIONAL QUALITY FORUM

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

November 9, 2010

NATIONAL QUALITY FORUM

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

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 for Selecting the Focus for Measure Development.....	10
17	II. Recommendations on Sources of Evidence and Evidence Grading for the Present and the	
18	Future.....	11
19	III. Recommendations for the Evidence Needed to Justify the Focus of a Quality Measure....	13
20	IV. Recommendations for Evaluating Criterion 1c – Quantity, Quality, Consistency of Body of	
21	Evidence.....	16
22	V. Recommendations for Evaluating Importance to Measure and Report and the Other	
23	Subcriteria.....	21
24	Consequences of Measurement.....	24
25	VI. Recommendations for Modifications to the NQF Evaluation Criteria.....	25
26	VII. Recommendations for Modifications to the Measure Submission.....	27
27	VIII. Recommendations for Evidence Required for Practices Considered for NQF	
28	Endorsement.....	31
29	REFERENCES.....	31
30	APPENDIX A – EVALUATION CRITERIA.....	34
31	Current Measure Evaluation Criteria.....	34
32	Current Evaluation Criteria for Practices.....	39
33	APPENDIX B - TASK FORCE MEMBERS.....	40
34	APPENDIX C - US PREVENTIVE SERVICES TASK FORCE SYSTEM FOR GRADING	
35	EVIDENCE AND RECOMMENDATIONS.....	41

38 OVERVIEW AND PURPOSE

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

49
50 NQF's [evaluation criteria](#) require a variety of evidence as noted in the following table. Of these
51 criteria, some of the most rigorous evidence is required to justify what is being measured (1c)
52 and that is the primary focus of this report – *the evidence required to justify the measure focus*
53 (i.e., the specific process, structure, outcome, etc. that is being measured). Another task force
54 and subsequent report will address measure testing and the criterion of *Scientific Acceptability of*
55 *Measure Properties*.

56
57 Evidence refers to the information used to determine or demonstrate the truth of a hypothesis.
58 The highest quality evidence available should be used to support the focus of quality
59 performance measures. Evidence is not limited to quantitative studies and the best type of
60 evidence depends upon the question being studied (e.g., randomized controlled trials
61 appropriate for studying drug efficacy are not well suited for complex system changes). A body
62 of evidence includes all the evidence for a topic, which is systematically identified, based on
63 pre-established criteria for relevance and quality of evidence.

64
65 NQF endorses measures that are intended for use in public reporting as well as quality
66 improvement with the goal of improving the quality of healthcare. The evidence that supports
67 the focus for a quality measure is addressed under the must-pass criterion, *Importance to*
68 *Measure and Report* because if the measure focus is not supported by evidence that it can

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

72

73 Table 1. Measure Evaluation Criteria and Type of Evidence

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

74

75 **Task Force Charge**

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

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

84

85

86 **BACKGROUND**

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

89 However, much of healthcare has not been subjected to research studies, much less with
90 randomized controlled trials or comparative effectiveness studies. Lohr observed that “Perhaps
91 no more than half, or even one-third, of services are supported by compelling evidence that
92 benefits outweigh harms ¹.” For example, Tricoci, et al. ² reviewed recommendations in
93 American College of Cardiology/American Heart Association guidelines and found that only
94 314 of 2711 recommendations were classified as A-level evidence based on multiple
95 randomized trials with large numbers of patients. Many quality performance measures are
96 based on clinical practice guidelines, however not all guideline recommendations are
97 appropriate for performance measure development, which depends on the strength of the
98 evidence and relationship to meaningful outcomes ³.

99
100 Some aspects of healthcare (e.g., system change) may be more difficult to study with
101 quantitative methods, particularly with randomized controlled trials. Some clinical process
102 steps (i.e., assessing health status, diagnosing clinical conditions, recommending treatment,
103 teaching and counseling about conditions/treatment) may be unlikely to be subjected to
104 research. Even when research has been conducted, the body of evidence may not have been
105 systematically assessed and graded (e.g., care coordination, medication management). Lohr ¹
106 noted that absence of evidence about benefit is not the same as evidence of no benefit. Even
107 when available, evidence is rarely definitive. However, the level of confidence in a
108 recommendation (or measure) depends on the underlying research and synthesis of that
109 research.

110

111 Evidence Issues Identified with Measures Submitted to NQF

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

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

125
126 NQF consensus projects were not intended to undertake systematic evidence reviews for the
127 variety of measures that are submitted for consideration, nor is this feasible. Such detailed
128 evidence reviews have also not generally been viewed by developers as an integral part of the
129 measure development process. However, the responsibility for basing quality performance
130 measures on appropriate evidence does ultimately lie with measure developers. Measure
131 developers who do not have the expertise and resources to systematically assess the strength of
132 a body of evidence sometimes rely on other sources of evidence reviews and grading, such as
133 found in clinical practice guidelines or published systematic reviews. However, NQF wishes to
134 clearly signal, through this document and the measure submission form itself, that measure
135 developers are responsible for identifying, summarizing, and reporting the evidence that exists
136 to support the focus of measures submitted to NQF for potential endorsement.

137
138 **The Changing Environment**

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

151 conform to the highest measurement science principles. Recognizing the high stakes of
152 performance measurement in an increasingly transparent environment, some measure
153 developers have enhanced their requirements for the evidence base for performance measure
154 development ⁴.

155

156 **Clinical Practice Guidelines**

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

167

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

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

178

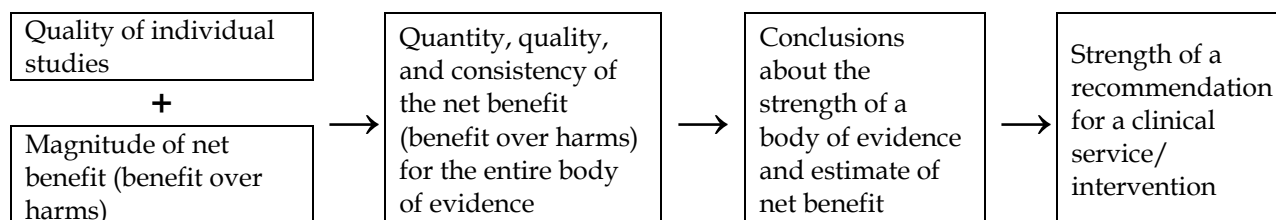
179 **Evidence Grading Systems**

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

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

185
186 There are commonalities among the various evidence grading systems. In general, the quality
187 and strength of the overall body of evidence is a function of the *quantity* and *quality* of
188 individual studies and the *consistency* among studies regarding judgments of net benefit (the
189 balance of benefits and harms). *Quality* of individual studies includes study design, sample size
190 and statistical power considerations, flaws such as selection bias, directness of the evidence
191 linking an intervention to health outcomes, and generalizability of findings. Of particular
192 interest for quality measures is how well the measure matches the population and intervention
193 in the evidence (e.g., cited studies). The general approach to determining the strength of
194 evidence and a recommendation for a particular intervention or service is depicted in Figure 1.

195
196 Figure 1. Approach to Determining Quality of Evidence and Strength of Recommendation



197
198 Differences in terminology and grading scales may inhibit understanding about the strength of
199 evidence. Differences can range from a rather minor but understandable difference in
200 terminology (e.g., strength, quality, or level of evidence) to pronounced differences in the
201 assignment of grades (e.g., a grade of A could indicate evidence based on consensus of opinion
202 in one system to evidence based on meta-analyses of randomized controlled trials in another
203 system). An international initiative to standardize grading evidence and recommendations,
204 [GRADE](#)⁹⁻¹⁵, is now supported by many [organizations](#) including the Cochrane Collaboration.
205 The Agency for Healthcare Research and Quality (AHRQ) supports two evidence grading
206 systems: one used by the US Preventive Services Task Force (USPSTF)^{16,17} and one used by the
207 Evidence-Based Practice Centers¹⁸ (consistent with GRADE). Table 2 provides examples of
208 terminology used by four evidence grading systems. It is important to note that grading
209 systems are tied to specific methods for reviewing and assessing the quality of evidence.

210

211 Table 2. Examples of Terminology in Selected Grading Scales

	<u>USPSTF</u>	<u>GRADE</u>	<u>AHRQ Evidence-Based Practice Centers</u>	<u>ACC/AHA</u>
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

212

213

214 Systematic reviews and meta-analyses are used to assess a body of evidence. PRISMA
 215 (Preferred Reporting Items for Systematic Reviews and Meta-Analyses) focuses on the
 216 transparent and full reporting of such reviews¹⁹. The Institute of Medicine (IOM) has two
 217 consensus projects underway that relate to grading the quality of evidence for clinical

218 interventions: [Standards for Developing Trustworthy Clinical Practice Guidelines](#) and
219 [Standards for Systematic Reviews of Clinical Effectiveness Research](#); however, reports will not
220 be ready until early 2011.

221

222

223 RECOMMENDATIONS

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

226

227 **Evidence** refers to the information used to determine or demonstrate the truth of a hypothesis.

228 The highest quality evidence available should be used to support the focus of quality
229 performance measures. Evidence is not limited to quantitative studies and the best type of
230 evidence depends upon the question being studied (e.g., randomized controlled trials
231 appropriate for studying drug efficacy are not well suited for complex system changes).

232

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

235

236 Principles

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

240

241 **Measures that will be used for public reporting should meet a high standard of evidence for**
242 **the focus of measurement.** NQF measures are intended to be useful for public reporting, as
243 well as to internal quality improvement activities. Measures used for public reporting often
244 impact large numbers of providers and entail investment of significant resources in
245 measurement and improvement. Consequently, measures that will be used for public reporting
246 should meet a high standard of evidence for the focus of measurement. The net benefit to
247 patients should outweigh any potential harm to patients, and be clinically or practically

248 meaningful to justify implementation. A lower standard of evidence may be deemed
249 appropriate by those selecting measures for use in smaller scale internal quality improvement
250 activities within a learning system that allows for rapid adjustments. Such measures, although
251 potentially of value, are not considered by NQF as they are not appropriate for public reporting.

252

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

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

263

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

269

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

274

275 I. Recommendations for Selecting the Focus for Measure Development

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

278

- 279 • There is a hierarchical preference for outcome measures (when possible) followed by
280 process measures, then structure measures. Outcome measures are preferred because
281 improving health outcomes is a central goal of healthcare. However, both outcome and
282 process measures have advantages and disadvantages²⁵ and both have a place in quality
283 assessment and the NQF portfolio.
- 284 • For process and structure measures, the focus of measurement should be on the aspect of
285 care with the most direct evidence of a strong relationship to the desired outcome. For
286 example, evidence about effective medication to control blood pressure is direct evidence
287 for the medication but only indirect evidence for the frequency of assessing blood pressure
288 (see Figure 2). Assessment of blood pressure, although necessary, is not sufficient to achieve
289 control. When there are multiple processes that affect a desired outcome, efforts should be
290 made to include measures for all processes that have a strong relationship to the desired
291 outcome.
- 292 • Specific drugs and devices included in quality performance measures should be FDA-
293 approved for the target condition.
- 294 • Structural measures are appropriate primarily when there are very well established
295 structure-process-outcome relationships; and when it is not feasible to directly measure the
296 outcome or processes.
- 297 • For any topic area, measures based on the best evidence should be considered over
298 measures based on lower quality evidence (e.g., expert opinion).
- 299

300 II. Recommendations on Sources of Evidence and Evidence Grading for the Present and the Future

- 301 • The preferred sources of evidence for quality performance measures are systematic reviews
302 and grading of a body of evidence conducted by independent organizations such as
303 [USPSTE](#), [AHRQ Evidence-based Practice Centers](#), and the [Cochrane Collaboration](#); or
304 guidelines that meet national standards for trustworthy guidelines (as being developed by
305 the IOM).
- 306 • Until such time when guidelines are certified as meeting a set standard, preferred guidelines
307 are those developed with balanced representation beyond one specialty group and with full
308 disclosure of biases and how they were addressed. Further, the evidence underlying a

309 guideline recommendation must be accessible in order to provide the information necessary
310 to meet the requirements set out in this report.

311 • An assigned evidence grade alone is not sufficient to evaluate whether the NQF criterion on
312 evidence for the focus of measurement (1c) is met, either now or in the future. The specific
313 information on the quantity, quality, and consistency of the body of evidence that was used
314 to determine an overall grade should be summarized in the measure submission.

315 • Explicit, transparent information on the quantity, quality, and consistency of the body of
316 evidence supporting a measure will facilitate identification of guideline recommendations
317 that do not have acceptable evidence as the basis for performance measurement. Explicit
318 information about the evidence also facilitates review by all stakeholders although TAPs
319 and Steering Committees will continue to include experts that possess knowledge about the
320 state of science for a particular topic.

321 • **Current Expectations**

322 ○ Most measure developers will rely on evidence reviews and grading conducted by
323 other organizations such as guideline developers or published systematic reviews.
324 However, it is the responsibility of the measure developer to understand the
325 strength of the evidence on which it is basing a measure and to provide a concise
326 summary of this evidence, not simply the end-result of the grading process.
327 Information on the evidence is useful to committees reviewing measures and the
328 public who use the measures.

329 ○ To promote transparency and standardization, NQF should require measure
330 developers to provide specific information about the quantity, quality, and
331 consistency of the body of evidence underlying a quality performance measure.
332 Information should include who graded the evidence, the evidence grading system
333 used and the grade assigned. If the developer fails to provide this information, NQF
334 should not review the proposed measure.

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

338

339

340 • **Future Expectations**

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

- 345 ○ Most measure developers will continue to rely on evidence reviews and grading
346 conducted by other organizations.
- 347 ○ Rather than identifying “preferred” grading systems as noted for the current
348 expectations, NQF should require that evidence used to support measures be graded
349 using one or two standardized evidence grading systems (e.g., the USPSTF, GRADE,
350 or possibly one adopted by the IOM).
- 351 ○ The evidence should be graded by identified credible sources, such as guideline
352 developers or review organizations, certified as meeting accepted standards.
- 353 ○ Even when basing measures on evidence graded with a standardized grading
354 system and potentially certified reviewers, explicit information on the quantity,
355 quality, and consistency of the specific evidence that led to the assignment of a grade
356 should be submitted for evaluation.

357

358 **III. Recommendations for the Evidence Needed to Justify the Focus of a Quality Measure**

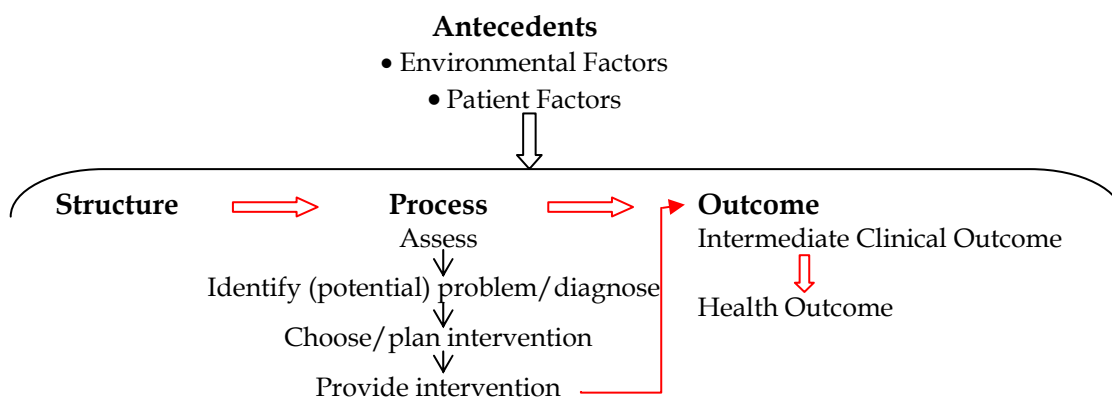
359 There has been widespread acceptance of Donabedian’s^{20, 21} structure-process-outcome model
360 for assessing healthcare quality. These three approaches to quality measurement can be used
361 with any topic of healthcare quality and the evidence required generally does not vary by topic.
362 The required evidence is for the links depicted by the red arrows in Figure 2. As depicted under
363 process, there may be multiple process steps prior to delivering an intervention; however, the
364 evidence is most often about the relationship between the intervention and outcome and
365 therefore, interventions are the preferred focus of process measures. Antecedents are depicted
366 in Figure 2. Although they influence structures, processes, and outcomes, patient factors that
367 influence outcomes are important to consider for risk adjustment for outcome measures.

368

369

370

371 Figure 2. Structure-Process-Outcome Model



372

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

377

378 As noted by the Task Force and articulated by NQF's Board of Directors, there is a preference
379 for measures of health outcomes. Achieving or improving health outcomes is a central goal of
380 healthcare treatments and services (e.g., health, function, survival, symptom control). Outcomes
381 also are viewed as useful quality indicators because they are integrative of the influence of
382 multiple care processes and disciplines involved in the care. . Further, once outcomes are
383 measured and reported, many outcomes that were not thought to be modifiable tend to be
384 improved and stimulate identification and adoption of effective practices. Because multiple
385 processes may influence a health outcome, several bodies of evidence could be relevant. For the
386 reasons noted above, measures of health outcomes are considered an exception to the
387 requirement of submitting a review of an empirical body of evidence. Instead, a rationale that
388 supports the relationship of the measured health outcome to processes of care and/or the
389 importance of measuring the outcome is considered acceptable.

390

391

392

393

394

395 Table 3. Evidence to Support the Focus of Measurement

Type of Measure	Evidence	Example of Measure Type & Evidence to be Addressed
<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 may be considered a proxy for a health state (e.g., hospitalization may represent a deterioration in health status)</p>	<p>A rationale supports the relationship of the health outcome to processes of care and/or the importance of measuring the outcome. See Table 5</p>	<p>#0230 Acute Myocardial Infarction 30-day Mortality Survival is a goal of seeking and providing treatment for AMI Rationale linking healthcare processes/ interventions (aspirin, reperfusion) to 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. Rationale linking healthcare processes (e.g., medication reconciliation, care coordination) to 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. Rationale linking healthcare processes (e.g., ventilator bundle) to ventilator acquired pneumonia</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 a desired health outcome See Table 4</p>	<p>#0059 Hemoglobin A1c Management [A1c>9] Evidence that hemoglobin A1c level leads to health outcomes (e.g., prevention of renal disease, heart disease, amputation, mortality)</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 with benefits that outweigh harms to patients</p> <p>Specific drugs and devices should have FDA approval for the target condition</p> <p>If the measure focus is on inappropriate use: Quantity, quality, and consistency of a body of evidence that the measured healthcare process does <u>not</u> lead to desired health outcomes in the target population See Table 4</p>	<p>#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 results in 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>Structure Structure of care is a feature of a health care organization or</p>	<p>Quantity, quality, and consistency of a body of evidence that the measured healthcare structure leads to desired health outcomes</p>	<p>#0190 Nurse Staffing Hours Evidence that higher nursing hours results in lower mortality, morbidity ; or</p>

Type of Measure	Evidence	Example of Measure Type & Evidence to be Addressed
clinician related to its capacity to provide high quality health care	with benefits that outweigh harms (including evidence for the link to effective care processes and the link from the care processes to desired health outcomes) See Table 4	leads to provision of effective care processes (e.g., lower medication errors) that lead to better outcomes
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 (often acquired through qualitative studies) OR Evidence that patient experience with care is correlated with desired outcomes	#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

396

397 **IV. Recommendations for Evaluating Criterion 1c – Quantity, Quality, Consistency of Body of**
398 **Evidence**

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

405

- 406 • Evidence should be evaluated on the *quantity* of studies, *quality* of studies, and *consistency* in
407 direction and magnitude of net benefit (clinically or practically meaningful benefits over
408 harms to patients) of a ***body of evidence*** on a scale of High, Moderate, or Low.
- 409 • The dimensions of *quantity*, *quality*, and *consistency* of a body of evidence apply to measures
410 based on guidelines as well as those for which guidelines may not exist (e.g., measures of

- 411 care coordination or team functioning may not be based on guidelines, but often have
412 bodies of evidence including non-clinical literature that should be systematically assessed)
- 413 • Measures without a clear description of the *quantity, quality, and consistency* of the
414 supporting body of evidence or without any evidence should not pass criterion 1c and the
415 threshold criterion of *Importance to Measure and Report*.
 - 416 • Use of only selected studies rather than an entire body of evidence that meets pre-
417 established criteria is not adequate to evaluate the evidence and should not pass criterion 1c
418 and the threshold criterion of *Importance to Measure and Report*.
 - 419 • Inconsistent and conflicting evidence should result in measures not passing both criterion 1c
420 and the threshold criterion of *Importance to measure and report*.
 - 421 • Outcome measures are considered an exception to the evidence requirement. A rationale
422 should support the relationship of the outcome to processes of care and/or the importance
423 of measuring the outcome.
 - 424 • Expert opinion is not considered empirical evidence and will only be considered in
425 exceptional circumstances when all of the following conditions are met.
 - 426 ○ No evidence is available.
 - 427 ○ Expert opinion is systematically assessed. That is, identified experts explicitly
428 address the certainty or confidence that benefits to patients from the specific process
429 or structure greatly outweigh potential harms, using a specified process that is
430 transparent and open to peer review (e.g., modified Delphi, formal consensus
431 process, [RAND Appropriateness Method](#)²²). The methods and results are reported
432 for review.
 - 433 ○ There is a strong rationale for why the specific structure or process should be the
434 focus of a quality performance measure.

435
436 Table 4 provides definitions and guidance on how to evaluate each of the dimensions of
437 *quantity, quality, and consistency* for a body of quantitative evidence. Each dimension is rated on
438 a scale of high, moderate, low, or inadequate to evaluate. A body of evidence could have
439 different ratings for each dimension, e.g., high on quantity, low on quality, and moderate on
440 consistency. Table 5 provides recommended decision rules for using the ratings for all three
441 dimensions to make a decision on whether a measure should pass criterion 1c, the evidence to

442 support the measure focus. Strong evidence usually requires multiple studies each with
443 sufficient numbers of patients to give precise estimates, but occasionally a large and
444 representative study can provide adequate evidence. For example, one study (low quantity) that
445 is a randomized controlled trial with a large representative sample of patients (high quality)
446 and substantial estimates of net benefit would pass the criterion, whereas, a body of evidence
447 with low consistency of estimates of net benefits indicates a measure should not pass the
448 criterion regardless of the ratings for quantity and quality of studies.

449
450 There are various ways to categorize research [study designs](#). However, for purposes of the
451 rating schema, the type of evidence for the structure-process-outcome linkages is grouped into
452 two categories as follows.

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

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

458
459 Although RCTs remain the gold standard for evidence of efficacy of treatment, there are many
460 areas where RCTs may not currently exist and are unlikely to be conducted. Furthermore, the
461 strict eligibility and exclusion criteria for randomized trials may sometimes result in findings
462 that are not fully generalizable in real world applications. NQF recognizes the evidentiary value
463 of well-conducted observational studies, particularly those that attempt to balance measured
464 covariates (e.g., using propensity scores) and account for other sources of bias as articulated in
465 the [GRACE principles](#) [Good Research for Comparative Effectiveness] ²³. This is particularly
466 true when there are multiple observational studies that arrive at similar conclusions.

467
468 Qualitative studies often are used to gain understanding of people’s attitudes, behaviors, and
469 values and may be suited to evidence regarding patient experience with care. The descriptions
470 of quality and consistency of the evidence in Table 4 do not apply to qualitative evidence. When
471 qualitative studies are used, appropriate qualitative research criteria should be used to judge
472 the strength of the evidence ²⁴.

473

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

478

479

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

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

482 *Study designs that affect certainty of confidence in estimates of effect include: Randomized controlled
 483 trials (RCT), which control for both observed and unobserved confounders, and non-RCTs (observational
 484 studies) with various levels of control for confounders.

485 Study flaws that may bias estimates of effect include: lack of allocation concealment; lack of blinding;
 486 large losses to follow-up; failure to adhere to intention to treat analysis; stopping early for benefit; failure
 487 to report important outcomes.
 488 Imprecision with wide confidence intervals around estimates of effects can occur in studies involving few
 489 patients and few events.
 490 Indirectness of evidence includes: indirect comparisons (e.g., two drugs compared to placebos rather than
 491 head-to head); differences between the population, intervention, comparator interventions, and outcome
 492 of interest and those included in the relevant studies. ¹⁴
 493 ** The suggested number of studies for rating levels of quantity is considered a general guideline.
 494

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

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

496

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

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

503

504 Generally, in measure submissions, high impact is easily demonstrated by alignment with a
505 specific NPP goal or epidemiologic or resource use data (incidence, prevalence, resource use,
506 consequences of quality problems). However, data on opportunity for improvement may be
507 lacking (e.g., submitter states that performance is unknown, or it may not be specific to the
508 focus of measurement, or only based on a sample from measure development and testing).
509 Reviewers sometimes question whether there is enough variation to justify importance to
510 measure and report, or how to judge overall poor performance. When a measure undergoes
511 review for continued endorsement, an issue that sometimes arises is whether a measure is
512 “topped out” meaning there are high levels of performance with little variation and therefore,
513 little room for further improvement.

514

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

522

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

532

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

539

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

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

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

548 o number and representativeness of the entities included in the measure performance
549 data; and

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

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

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

561

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

Pass Criterion, Importance to Measure and Report?			
All 3 subcriteria (1a,1b,1c) must be met to pass the threshold criterion, <i>Importance to Measure and Report</i>			
Subcriterion	Evidence	Example	Pass the subcriterion?
High impact (1a)	Addresses a <u>specific national health goal/priority</u> identified by the Secretary of DHHS or the NPP; OR Epidemiologic or resource use data; health services research - affects large numbers of patients and/or has a very substantial impact for smaller populations; leading cause of morbidity/mortality; high resource use (current and/or future); severity of illness; and patient/societal consequences of poor quality	#0140 Ventilator-associated pneumonia for ICU and high-risk nursery (HRN) patients 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)	Initial Endorsement Epidemiologic or resource use data; health services research - data demonstrating considerable variation, or overall less than optimal performance, for the focus of measurement across providers and/or population groups (disparities in care) Review for Endorsement Maintenance <u>Data for the measure as specified and endorsed</u> demonstrating considerable variation, or overall less than optimal performance	#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
All 3 subcriteria (1a,1b,1c) must be met to pass the threshold criterion, <i>Importance to Measure and Report</i>			

563

564 Consequences of Measurement

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

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

567 topic of measurement. Currently, unintended consequences of measurement are addressed
 568 under feasibility.

569 **4d.** Susceptibility to inaccuracies, errors, or unintended consequences of measurement and the
 570 ability to audit the data items to detect such problems are identified.

571
 572 The Task Force identified that actual vs. theoretical consequences to measurement are most
 573 likely to arise after implementation and should be addressed at the time of review for
 574 endorsement maintenance. For example, a measure of timing of antibiotic administration in
 575 patients with pneumonia may result in some patients receiving antibiotics before the diagnosis
 576 of pneumonia is confirmed by x-ray. The Task Force did not recommend moving subcriterion
 577 4d under *Importance to Measure and Report*.

578
 579 **VI. Recommendations for Modifications to the NQF Evaluation Criteria**

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

582
 583 Table 7. Current and Modified Measure Evaluation Criteria

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

Current Measure Evaluation Criteria	Modified Measure Evaluation Criteria
<p>considerable variation, or overall poor performance, in the quality of care across providers and/or population groups (disparities in care).</p> <p>1c. The measure focus is:</p> <ul style="list-style-type: none"> • an outcome (e.g., morbidity, mortality, function, health-related quality of life) that is relevant to, or associated with, a national health goal/priority, the condition, population, and/or care being addressed (2); OR • if an intermediate outcome, process, structure, etc., there is evidence (3) that supports the specific measure focus as follows: <ul style="list-style-type: none"> o <u>Intermediate outcome</u> – evidence that the measured intermediate outcome (e.g., blood pressure, Hba1c) leads to improved health/avoidance of harm or cost/benefit. o <u>Process</u> – evidence that the measured clinical or administrative process leads to improved health/avoidance of harm and if the measure focus is on one step in a multi-step care process (4), it measures the step that has the greatest effect on improving the specified desired outcome(s). o <u>Structure</u> – evidence that the measured structure supports the consistent delivery of effective processes or access that lead to improved health/avoidance of harm or cost/benefit. o <u>Patient experience</u> – evidence that an association exists between the measure of patient experience of health care and the outcomes, values and preferences of individuals/ the public. o <u>Access</u> – evidence that an association exists between access to a health service and the outcomes of, or experience with, care. o <u>Efficiency</u> (5) – demonstration of an association between the measured resource use and level of performance with respect to one or more of the other five IOM aims of quality. <p><i>If not important to measure and report, STOP.</i></p> <p>Footnotes 1 Examples of data on opportunity for improvement include, but are not limited to: prior studies, epidemiologic data, measure data from pilot testing or implementation. If data are not available, the measure focus is systematically assessed (e.g., expert panel rating) and judged to be a quality problem. 2 Generally, rare event outcomes do not provide adequate information for improvement or discrimination; however, “never events” that are compared to zero are appropriate outcomes for public reporting and quality improvement.</p>	<p>for improvement, i.e., data (footnote 1) demonstrating considerable variation, or overall less than optimal performance, in the quality of care across providers and/or population groups (disparities in care).</p> <p>AND</p> <p>1c. The measure focus is:</p> <ul style="list-style-type: none"> • a <u>health outcome</u> (footnote 2): with a rationale that supports the relationship of the health outcome to processes of care and/or the importance of measuring the health outcome; • OR • Is evidence-based as demonstrated by a systematic assessment and grading of the quantity, quality, and consistency of the body of evidence (footnote 3). • <u>Intermediate clinical outcome</u>: evidence that the measured intermediate clinical outcome leads to a desired health outcome • <u>Process</u> (footnote 4): evidence that the measured healthcare process leads to desired outcomes in the target population. • <u>Structure</u>: evidence that the measured structure leads to desired health outcomes (including evidence for the link to effective care processes and the link from the care processes to desired health outcomes). • Special Considerations by Topic of Measurement <ul style="list-style-type: none"> o <u>Patient experience with care</u>: evidence that the measured aspects of care are those valued by patients and for which the patient is the best and/or only source of information OR that patient experience with care is correlated with desired outcomes. o <u>Efficiency</u> (footnote 5): evidence for the quality component as noted above. <p>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, serious reportable events that are compared to zero are appropriate outcomes for public reporting and quality improvement. 3 The preferred systems for grading the evidence are the USPSTF grading definitions and methods, or GRADE. 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</p>

Current Measure Evaluation Criteria	Modified Measure Evaluation Criteria
<p>3 The strength of the body of evidence for the specific measure focus should be systematically assessed and rated (e.g., USPSTF grading system – 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.</p> <p>4 Clinical care processes typically include multiple steps: assess → identify problem/potential problem → choose/plan intervention (with patient input) → provide intervention → evaluate impact on health status. If the measure focus is one step in such a multi-step process, the step with the greatest effect on the desired outcome should be selected as the focus of measurement. For example, although assessment of immunization status and recommending immunization are necessary steps, they are not sufficient to achieve the desired impact on health status – patients must be vaccinated to achieve immunity. This does not preclude consideration of measures of preventive screening interventions where there is a strong link with desired outcomes (e.g., mammography) or measures for multiple care processes that affect a single outcome.</p> <p>5 Efficiency of care is a measurement construct of cost of care or resource utilization associated with a specified level of quality of care. It is a measure of the relationship of the cost of care associated with a specific level of performance measured with respect to the other five IOM aims of quality. Efficiency might be thought of as a ratio, with quality as the numerator and cost as the denominator. As such, efficiency is directly proportional to quality, and inversely proportional to cost. (NQF’s Measurement Framework: Evaluating Efficiency Across Episodes of Care; based on AQA Principles of Efficiency Measures).</p>	<p>step in such a multi-step process, the step with the strongest evidence for the link to the desired outcome should be selected as the focus of measurement.</p> <p>⁵ Measures of efficiency combine the concepts of resource use and quality (NQF’s Measurement Framework: Evaluating Efficiency Across Episodes of Care; AQA Principles of Efficiency Measures).</p>

584

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

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

591

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

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

599

600 Table 8. Current and Modified Measure Submission Items

Current Measure Submission (4.1) Items	Modified Measure Submission Items
	<p>Add to 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>High Impact (Measure evaluation criterion 1a) (for NQF staff use) <u>Specific NPP goal</u>: 1a.1. Demonstrated High Impact Aspect of Healthcare Affects large numbers Leading cause of morbidity/mortality Severity of illness Patient/societal consequences of poor quality Frequently performed procedure High resource use Other: 1a.3. Summary of Evidence of High Impact 1a.4. Citations for Evidence of High Impact Opportunity for Improvement (Measure evaluation criterion 1b) 1b.1. Briefly explain the benefits (improvements in quality) envisioned by use of this measure 1b.2. Summary of Data Demonstrating Performance Gap (<i>Variation or overall poor performance across providers</i>) 1b.3. Citations for Data on Performance Gap 1b.4. Summary of Data on Disparities by Population Group</p>	<p>High Impact (Measure evaluation criterion 1a) (for NQF staff use) Specific priority goal: 1a.1. Demonstrated High Impact Aspect of Healthcare Affects large numbers Leading cause of morbidity/mortality Severity of illness Patient/societal consequences of poor quality Frequently performed procedure High resource use Other: 1a.3. Summary of Evidence of High Impact (<i>provide epidemiologic or resource use data</i>) 1a.4. Citations for Evidence of High Impact Opportunity for Improvement (Measure evaluation criterion 1b) 1b.1. Briefly explain the benefits (improvements in quality) envisioned by use of this measure 1b.2. Summary of Data Demonstrating Performance Gap (<i>Variation or overall poor performance across providers</i>) 1b.3. Citations for Data on Performance Gap 1b.4. Summary of Data on Disparities by Population Group</p>

Current Measure Submission (4.1) Items	Modified Measure Submission Items
<p>1b.5. Citations for Data on Disparities</p> <p>1c.1. Relationship to Outcomes <i>(For non-outcome measures, briefly describe the relationship to desired outcome. For outcomes, describe why it is relevant to the target population.)</i></p> <p>1c.2. Type of Evidence <i>(Check all that apply)</i> Cohort study Observational study Evidence-based guideline Randomized controlled trial Expert opinion Systematic synthesis of research Meta-analysis Other: 1c.3.</p> <p>1c.4. Summary of Evidence <i>(For non-outcome measures, provide evidence of relationship to desired outcome. For outcomes, summarize any evidence that healthcare services/care processes influence the outcome.)</i></p> <p>1c.5. Rating of Strength/Quality of Evidence <i>(Also provide narrative description of the rating and by whom)</i></p> <p>1c.6. Method for Rating Evidence</p> <p>1c.7. Summary of Controversy/Contradictory Evidence</p>	<p>1b.5. Citations for Data on Disparities</p> <p>1c.1. Structure-Process-Outcome Relationship <i>(Briefly state the measure focus, e.g., structure, process, or outcome <u>and</u> identify the links and direction between: a) the measured health outcome and processes that influence the outcome; b) the measured process or intermediate clinical outcome and desired health outcome; or c) the measured structure and effective processes and desired outcome.)</i></p> <p>For health outcome measures, provide a rationale that supports the relationship of the health outcome to processes of care and/or the importance of measuring the outcome <i>(Provide references if applicable)</i></p> <p>For health outcome measures, items 1c.2 through 1c.15 may be skipped.</p> <p>c.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 (1c.3).</p> <p>1c.4. Summary of Body of Evidence Directness to focus of measurement & target population in proposed measure <i>(State the central topic, population, and outcomes addressed in the body of evidence and identify any differences from the measure focus and measure target population)</i> Quantity of Studies in Body of Evidence <i>(total number of studies, not articles):</i> Quality of Body of Evidence <i>(Summarize the certainty or confidence in the estimates of benefits and harms to patients <u>across studies</u> in the body of evidence resulting from <u>study factors</u> including: study design/ flaws; directness/indirectness regarding the specific process/structure being measured, outcomes assessed, target population, comparisons; imprecision (wide confidence intervals due to few patients or events):</i></p> <p>Consistency of Results across Studies <i>(Summarize the consistency of the magnitude and direction of the effect) :</i> Net Benefit <i>(Benefits over harms)</i> 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:</p>

Current Measure Submission (4.1) Items	Modified Measure Submission Items
<p>1c.8. Citations for Evidence (<i>Other than guidelines</i>)</p> <p>1c.9. Quote 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</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>By whom (<i>describe the entity that graded the evidence, including balance of representation and any disclosures regarding bias</i>) Grade Assigned to the Evidence:</p> <p>1c.6. System Used for Grading the Body of Evidence described above: 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 described above (<i>Other than guidelines</i>)</p> <p>If the measure is based on a clinical practice guideline, complete 1c.9-1c.14; otherwise complete 1c.15.</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 Guideline 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 Guideline 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>1c.15 Based on the NQF descriptions for rating the body of evidence, what was your assessment of the quantity, quality, and consistency of the body of evidence? (rate each as High, Moderate, or Low) Quantity: Quality: Consistency:</p>

602 VIII. Recommendations for Evidence Required for Practices Considered for NQF Endorsement
603 NQF also endorses practices such as [safe practices](#), care coordination practices, and substance
604 use treatment practices. The [criteria](#) for practices include evidence of effectiveness.

605
606 The Task Force recommends that the same evidence requirements as indicated for process
607 measures (Tables 3, 4, 5) be applied to practices considered for NQF endorsement.

608
609 Table 9. Evidence to Support a Practice

Evidence to Support a Practice	Example of Practice & Evidence to be Addressed
Quantity, quality, and consistency of a body of evidence that the measured healthcare process leads to desired health outcomes in the target population with benefits that outweigh harms to patients	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

611
612 **Modifications to Practice Evaluation Criteria**
613 **Evidence of Effectiveness.** A practice is evidence-based as demonstrated by a systematic
614 assessment of the quantity, quality, and consistency of the body of evidence and standardized
615 grading of the body of evidence. The preferred systems for grading the evidence are the
616 USPSTF [grading definitions](#) and [methods](#), or [GRADE](#). Evidence from non-healthcare industries
617 that should be substantially transferable to healthcare (e.g., safety practices of repeat-back of
618 verbal orders or standardizing abbreviations) also may be considered.

619
620
621 **REFERENCES**

622 1. Lohr KN. Rating the strength of scientific evidence: relevance for quality improvement
623 programs. *Int J Qual Health Care*. 2004;16(1):9-18.
624 2. Tricoci P, Allen JM, Kramer JM et al. Scientific evidence underlying the ACC/AHA clinical
625 practice guidelines. *JAMA*. 2009;301(8):831-841.
626 3. Spertus JA, Eagle KA, Krumholz HM et al. American College of Cardiology and American
627 Heart Association methodology for the selection and creation of performance measures
628 for quantifying the quality of cardiovascular care. *Circulation*. 2005;111(13):1703-1712.
629 4. Physician Consortium for Performance Improvement. Physician Consortium for
630 Performance Improvement® (PCPI) Position Statement - The Evidence Base Required for

- 631 Measures Development. *American Medical Association* 6-26-2009;1-18. Last accessed March
632 2010.
- 633 5. Grilli R, Magrini N, Penna A et al. Practice guidelines developed by specialty societies: the
634 need for a critical appraisal. *Lancet*. 2000;355(9198):103-106.
- 635 6. Shiffman RN, Shekelle P, Overhage JM et al. Standardized reporting of clinical practice
636 guidelines: a proposal from the Conference on Guideline Standardization. *Ann Intern Med*.
637 2003;139(6):493-498.
- 638 7. The AGREE Collaboration. *Appraisal of Guidelines for Research and Evaluation AGREE*
639 *Instrument*. 2001. Available at <http://www.agreecollaboration.org/instrument/>. Last
640 accessed February 2010.
- 641 8. The AGREE Collaboration. Development and validation of an international appraisal
642 instrument for assessing the quality of clinical practice guidelines: the AGREE project.
643 *Quality & safety in health care*. 2003;12(1):18-23.
- 644 9. Atkins D, Eccles M, Flottorp S et al. Systems for grading the quality of evidence and the
645 strength of recommendations I: critical appraisal of existing approaches The GRADE
646 Working Group. *BMC Health Serv Res*. 2004;4(1):38.
- 647 10. Atkins D, Best D, Briss PA et al. Grading quality of evidence and strength of
648 recommendations. *BMJ*. 2004;328(7454):1490-1494.
- 649 11. Guyatt GH, Oxman AD, Vist GE et al. GRADE: an emerging consensus on rating quality of
650 evidence and strength of recommendations. *BMJ*. 2008;336(7650):924-926.
- 651 12. Guyatt GH, Oxman AD, Kunz R et al. Incorporating considerations of resources use into
652 grading recommendations. *BMJ*. 2008;336(7654):1170-1173.
- 653 13. Guyatt GH, Oxman AD, Kunz R et al. Going from evidence to recommendations. *BMJ*.
654 2008;336(7652):1049-1051.
- 655 14. Guyatt GH, Oxman AD, Kunz R et al. What is "quality of evidence" and why is it important
656 to clinicians? *BMJ*. 2008;336(7651):995-998.
- 657 15. Guyatt GH, Oxman AD, Vist GE et al. GRADE: an emerging consensus on rating quality of
658 evidence and strength of recommendations. *BMJ*. 2008;336(7650):924-926.
- 659 16. Harris RP, Helfand M, Woolf SH et al. Current methods of the US Preventive Services Task
660 Force: a review of the process. *Am J Prev Med*. 2001;20(3 Suppl):21-35.
- 661 17. Sawaya GF, Guirguis-Blake J, LeFevre M et al. Update on the methods of the U.S. Preventive
662 Services Task Force: estimating certainty and magnitude of net benefit. *Ann Intern Med*.
663 2007;147(12):871-875.
- 664 18. Owens DK, Lohr KN, Atkins D et al. Grading the strength of a body of evidence when
665 comparing medical interventions-Agency for Healthcare Research and Quality and the
666 Effective Health Care Program. *J Clin Epidemiol*. 2009.
- 667 19. Liberati A, Altman DG, Tetzlaff J et al. The PRISMA statement for reporting systematic
668 reviews and meta-analyses of studies that evaluate health care interventions: explanation
669 and elaboration. *Ann Intern Med*. 2009;151(4):W65-W94.
- 670 20. Donabedian A. *An Introduction to Quality Assurance in Health Care*. New York, NY: Oxford
671 University Press; 2003.
- 672 21. Donabedian A. The role of outcomes in quality assessment and assurance. *Quality Review*
673 *Bulletin*. 1992;18(11):356-360.
- 674 22. Fitch K, Bernstein SJ, Aguilar MS et al. *The RAND/UCLA Appropriateness Method User's*
675 *Manual*. Santa Monica, CA: RAND Health; 2000. Available at
676 http://www.rand.org/pubs/monograph_reports/MR1269/.

- 677 23. Dreyer NA, Schneeweiss S, McNeil BJ et al. GRACE principles: recognizing high-quality
678 observational studies of comparative effectiveness. *Am J Manag Care*. 2010;16(6):467-471.
679 24. Cohen DJ, Crabtree BF. Evaluative criteria for qualitative research in health care:
680 controversies and recommendations. *Ann Fam Med*. 2008;6(4):331-339.
681 25. Donabedian A. The role of outcomes in quality assessment and assurance. *Quality Review*
682 *Bulletin*. 1992;18(11):356-360.
683
684

685
686
687
688
689

APPENDIX A – EVALUATION CRITERIA

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 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).

690
691

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.

692 **Current Evaluation Criteria for Practices**

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

698
699 **Benefit.** If the practice were more widely utilized, it would save lives endangered by healthcare
700 delivery, reduce disability or other morbidity, or reduce the likelihood of a serious reportable
701 event (e.g., an effective practice already in near universal use would lead to little new benefit to
702 patients by being designated a safe practice).

703
704 **Evidence of Effectiveness.** There must be clear evidence that the practice would be effective in
705 reducing patient safety events. Such evidence may take various forms, including the following:

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

715
716 **Generalizability.** The safe practice must be able to be utilized in multiple applicable clinical
717 care settings (e.g., a variety of inpatient and/or outpatient settings) and/or for multiple types of
718 patients.

719
720 **Readiness.** The necessary technology and appropriately skilled staff must be available to most
721 healthcare organizations.

722

723 **APPENDIX B - TASK FORCE MEMBERS**

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

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

728

729 **Kristine Martin Anderson, MBA**

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

732

733 **David Atkins MD, MPH**

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

736

737 **Arthur Levin, MPH**

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

740

741 **Mary Naylor, PhD, RN**

742 Marian S. Ware Professor in Gerontology
743 University of Pennsylvania School of Nursing
744 Board member

745

746 **Greg Pawlson, MD, MPH**

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

748

749 **Eric Schneider, MD, MSc, FACP**

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

755

756 APPENDIX C - US PREVENTIVE SERVICES TASK FORCE SYSTEM FOR GRADING EVIDENCE AND
 757 RECOMMENDATIONS

758 The following information was obtained from AHRQ websites describing the [United States](#)
 759 [Preventive Services Task Force \(USPSTF\) grade definitions](#) and [methods](#).

760
 761 **Table 1. U.S. Preventive Services Task Force Recommendation Grid***
 762

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			

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

766 **[What the Grades Mean and Suggestions for Practice](#)**

767 The USPSTF updated its definitions of the grades it assigns to recommendations and now includes "suggestions for
 768 practice" associated with each grade. The USPSTF has also defined levels of certainty regarding net benefit. These
 769 definitions apply to USPSTF recommendations voted on after May 2007.
 770
 771

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 Statement	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.

772
 773
 774 **Table 2. Questions Considered by the U.S. Preventive Services Task Force for Evaluating Evidence Related**
 775 **Both to Key Questions and to the Overall Certainty of the Evidence of Net Benefit for the Preventive Service**

- | |
|---|
| <ol style="list-style-type: none"> 1. Do the studies have the appropriate research design to answer the key question(s)? 2. To what extent are the existing studies of high quality? (i.e., what is the internal validity?) 3. To what extent are the results of the studies generalizable to the general U.S. primary care population and situation? (i.e., what is the external validity?) 4. How many studies have been conducted that address the key question(s)? How large are the studies? (i.e., what is the precision of the evidence?) 5. How consistent are the results of the studies? 6. Are there additional factors that assist us in drawing conclusions (e.g., presence or absence of dose-response) |
|---|

effects, fit within a biologic model)?

776
777
778

Table 3. U.S. Preventive Services Task Force Levels of Certainty Regarding Net Benefit

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 that are not generalizable to routine primary care practice a lack of information on important health outcomes. More information may allow an estimation of effects on health outcomes.

779 *The U.S. Preventive Services Task Force (USPSTF) defines *certainty* as "likelihood that the USPSTF assessment of
780 the net benefit of a preventive service is correct." The net benefit is defined as benefit minus harm of the preventive
781 service as implemented in a general primary care population. The USPSTF assigns a certainty level based on the
782 nature of the overall evidence available to assess the net benefit of a preventive service.
783

784
785
786
787

Table 4. U.S. Preventive Services Task Force Terminology to Describe the Critical Assessment of Evidence at 3 Levels: Individual Studies, Key Questions, and Overall Certainty of Net Benefit of the Preventive Service

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

788 *This terminology is not reflected in the carotid artery stenosis screening recommendation statement in this issue,¹
789 but it will appear in future recommendation statements.
790
791