

NATIONAL QUALITY FORUM

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

January 2011

NATIONAL QUALITY FORUM

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

CONTENTS

OVERVIEW AND PURPOSE 1

BACKGROUND 3

 Evidence Issues Identified with Measures Submitted to NQF..... 3

 The Changing Environment 4

 Clinical Practice Guidelines 5

 Evidence Grading Systems..... 6

RECOMMENDATIONS 8

 Definitions 9

 Principles 9

 I. Recommendations for Selecting the Focus for Measure Development..... 10

 II. Recommendations on Sources of Evidence and Evidence Grading for the Present and the
 Future 11

 III. Recommendations for the Evidence Needed to Justify the Focus of a Quality Measure ... 13

 IV. Recommendations for Evaluating Criterion 1c—Quantity, Quality, Consistency of Body of
 Evidence 16

 V. Recommendations for Evaluating Importance to Measure and Report and the Other
 Subcriteria 21

 VI. Recommendations for Modifications to the NQF Evaluation Criteria 24

 VII. Recommendations for Modifications to the Measure Submission Form 27

 VIII. Recommendations for Evidence Required for Practices Considered for NQF
 Endorsement..... 30

APPENDIX A EVALUATION CRITERIA FOR MEASURES AND PRACTICES 34

APPENDIX B TASK FORCE MEMBERS 41

APPENDIX C U.S. PREVENTIVE SERVICES TASK FORCE SYSTEM FOR GRADING
EVIDENCE AND RECOMMENDATIONS 42

NATIONAL QUALITY FORUM

OVERVIEW AND PURPOSE

Steering committees have diverse backgrounds and expertise and could benefit from more guidance and support to consistently apply NQF measure evaluation criteria. Both evidence and expert judgment play a role in evaluating measures against criteria. However, judgment can best be applied when Steering Committees have a thorough understanding of the evidence that does or does not exist. Evidence comes in many different forms (e.g., peer-reviewed publications, practice guidelines from authoritative sources, expert assessments); it is often inconsistent and incomplete; and it can be difficult to interpret and reach conclusions about. In October 2009, the Board directed NQF to strengthen its processes for evaluating the synthesis and scoring of evidence and for presenting this information in ways that will be best understood and useful to Steering Committees. To comply with the Board's directive, NQF convened two task forces: one to evaluate the subcriteria under *Importance to Measure and Report*, particularly the evidence that supports the focus of measurement, which is the subject of this report, and the other to address the criterion of *Scientific Acceptability of Measure Properties*, which is the subject of another report.

NQF's [evaluation criteria](#) require a variety of evidence (see Table 1). Some of the most rigorous evidence is required to support the focus of measurement (subcriterion 1c), that is, the specific process, structure, or outcome that is being measured. Evidence refers to the information that is used to determine or demonstrate the truth of a hypothesis. The highest quality evidence available should be used to support the focus of quality performance measures. Evidence is not limited to quantitative studies, and the best type of evidence depends upon the question being studied (e.g., randomized controlled trials [RCTs] appropriate for studying drug efficacy are not well suited for complex system changes). A body of evidence includes all the evidence for a topic, which is systematically identified, based on pre-established criteria for relevance and quality.

NATIONAL QUALITY FORUM

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 2a.-2g. Reliability, validity, risk adjustment	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

NQF endorses measures that are intended for use in public reporting as well as for internal quality improvement activities, with the goal of improving the quality of American healthcare. The evidence that supports the focus of a quality measure is considered under the threshold, or “must-pass” criterion, *Importance to Measure and Report* because if the measure focus is not supported by evidence that it can facilitate gains in quality and health, then the use of limited resources for measuring and reporting on it would be questionable. For most healthcare quality measures, the evidence will be that of clinical effectiveness and the link to desired outcomes.

Task Force Charge

The Evidence Task Force was charged with the following tasks:

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

NATIONAL QUALITY FORUM

BACKGROUND

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

However, much of healthcare has not been subjected to research studies, much less with randomized controlled trials or comparative effectiveness studies. Lohr observed, “Perhaps no more than half, or even one-third, of services are supported by compelling evidence that benefits outweigh harms.”¹ For example, Tricoci et al. reviewed recommendations in American College of Cardiology/American Heart Association guidelines and found that only 314 of 2,711 recommendations were based on A-level evidence,² that is, evidence derived from multiple randomized trials with large numbers of patients. Many quality performance measures are based on clinical practice guidelines; however, not all guideline recommendations are appropriate for performance measure development, which depends on the strength of the evidence and the relationship to meaningful outcomes.³

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

Evidence Issues Identified with Measures Submitted to NQF

The NQF evaluation criteria (1c, footnotes 3 and 4) and submission questions may not provide enough direction to reviewers or measure developers. Consequently, measure submissions often lack sufficient information about the strength of the evidence or the strength of a guideline recommendation. Measures have been submitted with no evidence or low-quality evidence. Measures have been submitted with no systematic grading or incorrect grading of the evidence or

NATIONAL QUALITY FORUM

guideline recommendation. A grading system other than the recommended U.S. Preventive Services Task Force (USPSTF) system may be used without providing any explanation. In some cases, measure developers assigned a grade without using the associated methods to assess the body of evidence, or they believed that the USPSTF evidence grading system is only applicable to preventive services. Finally, some measures are focused on process steps that are far removed from the desired outcome, even though there is evidence for a particular intervention or intermediate outcome that is more directly linked to the desired outcome (e.g., measures to assess immunization status rather than measures to assess administration of vaccine).

NQF consensus projects were not intended to undertake systematic evidence reviews for the variety of measures that are submitted for consideration, nor would it be feasible for them to do so. Such detailed evidence reviews also have not been viewed by measure developers as an integral part of the measure development process. Measure developers generally rely on other sources of evidence reviews and grading, such as those found in clinical practice guidelines or published systematic reviews. However, the responsibility for basing quality performance measures on appropriate evidence ultimately lies with measure developers. NQF wishes to clearly signal, through this document and the measure submission form itself, that measure developers are responsible for reporting on the body of evidence that supports the focus of measures submitted to NQF for potential endorsement.

The Changing Environment

As guidelines and quality metrics are increasingly used not only for internal quality improvement but also for public reporting, the necessity for a strong evidence base has become more urgent and compelling. This need is further substantiated by the development of reimbursement programs that utilize such publicly reported metrics. Although public reporting and pay for performance have the potential to inform consumers, focus quality improvement activities, and reward high performance, unintended negative consequences might result if measures do not meet all the aspects of the importance criterion. Such consequences include confusion about the importance of particular elements of care to quality and the diversion of scarce resources to implement and measure aspects of care that have marginal or no impact on quality. To achieve the intended positive effects of quality measurement and to minimize the unintended negative

NATIONAL QUALITY FORUM

consequences, measures should be based on high-quality evidence that supports the focus of measurement. Quality measures also should conform to the measurement science principles, which are addressed under the criterion, *Scientific Acceptability of Measure Properties*. Recognizing the high stakes of performance measurement in an increasingly transparent environment, some measure developers have enhanced their evidence requirements.⁴

Clinical Practice Guidelines

Although they are not the only evidence base for performance measures, clinical practice guidelines are used by many measure developers to support the focus of measurement.^{3,4} There has been a proliferation of such guidelines, some overlapping or even contradictory. There also is substantial variability in the methodological rigor of review and grading of the evidence and recommendations. In 2000, Grilli and colleagues reported that of 431 specialty society guidelines reviewed, 82 percent did not apply explicit criteria to grade the scientific evidence used as a basis for recommendations, 87 percent did not report whether a systematic literature search was conducted, and 67 percent did not describe the professional involved.⁵ Some tools to assess clinical practice guidelines⁶⁻⁸ are available, and the Institute of Medicine (IOM) is studying [standards for the development of trustworthy guidelines](#).

At the January 11, 2010, IOM meeting on developing trustworthy guidelines, Vivian Coates [presented](#) the following information about the [National Guidelines Clearinghouse](#) (NGC):⁹

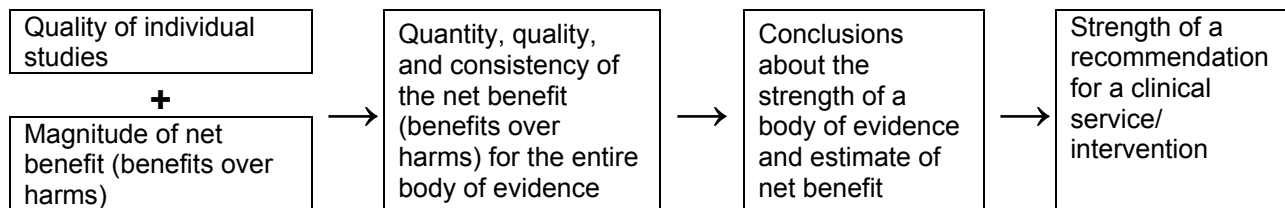
- Currently, NGC contains more than 2,500 guidelines from more than 200 measure developers.
- Most of the measure developers whose guidelines are represented in NGC (158 of 204; 77 percent) use some sort of rating scheme to grade the underlying evidence and/or strength of the recommendations. Of these:
 - Ten developers report using GRADE or modified GRADE.
 - Six report using the USPSTF approach, either as is, or modified.
 - The majority (142 developers) does not identify the origin of their rating schemes and appear to be using schemes unique to their organizations.

Evidence Grading Systems

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

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

Figure 1: Approach to Determining Quality of Evidence and Strength of Recommendation



Differences in terminology and grading scales may inhibit understanding of the strength of evidence. Differences can range from a rather minor but understandable difference in terminology (e.g., strength, quality, or level of evidence) to pronounced differences in the assignment of grades (e.g., a grade of A could indicate evidence based on consensus of opinion in one system and evidence based on meta-analyses of randomized controlled trials in another system). An international initiative to standardize grading evidence and recommendations, Grading of Recommendations Assessment, Development and Evaluation ([GRADE](#)),¹⁰⁻¹⁶ is now

NATIONAL QUALITY FORUM

supported by many [organizations](#) including the Cochrane Collaboration. The Agency for Healthcare Research and Quality (AHRQ) supports two evidence grading systems: the aforementioned one used by the USPSTF^{17, 18} and one used by the Evidence-based Practice Centers¹⁹ (consistent with GRADE). Table 2 provides examples of terminology used by four evidence grading systems. It is important to note that grading systems are tied to specific methods for reviewing and assessing the quality of evidence.

Systematic reviews and meta-analyses are used to assess a body of evidence. PRISMA (Preferred Reporting Items for Systematic Reviews and Meta-Analyses) focuses on the transparent and full reporting of such reviews.²⁰ The IOM has two consensus projects under way that relate to grading the quality of evidence for clinical interventions: [Standards for Developing Trustworthy Clinical Practice Guidelines](#) and [Standards for Systematic Reviews of Clinical Effectiveness Research](#); however, reports will not be ready until early 2011.

NATIONAL QUALITY FORUM

Table 2: Examples of Terminology in Selected Grading Scales

	<u>USPSTF</u>	<u>GRADE</u>	<u>AHRQ Evidence-based Practice Centers</u>	<u>American College of Cardiology Foundation/American Heart Association</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 populations, RCT, meta-analysis • B: limited population, single RCT or non-RCT • C: very limited population, consensus expert opinion, case studies Size of treatment effect <ul style="list-style-type: none"> • Class 1: 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 Moderate/Small • D - Recommend against: High or Moderate/Zero-Negative • 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

RECOMMENDATIONS

The Task Force identified some definitions and principles that guided its discussion and the recommendations that follow.

NATIONAL QUALITY FORUM

Definitions

Evidence refers to the information used to determine or demonstrate the truth of a hypothesis. The highest quality evidence available should be used to support the focus of quality performance measures. 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).

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

Principles

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

Measures that will be used for public reporting should meet a high standard of evidence for the focus of measurement. NQF endorses measures intended for public reporting, as well as internal quality improvement activities. Public reporting of measures often impact large numbers of providers and entail investment of significant resources in measurement and improvement. Therefore, measures that will be used for public reporting should meet a high standard of evidence for the focus of measurement. The focus of measurement should have evidence of a net benefit to patients that outweighs any potential harm to patients and also be clinically or practically meaningful to justify implementation. A lower standard of evidence may be deemed appropriate by those selecting measures for use in smaller scale, internal, quality improvement activities within a learning system that allows for rapid adjustments. Such measures, although potentially of value, are not considered by NQF because they are not appropriate for public reporting.

In the absence of strong evidence of certainty of net benefit for a structure or process being measured, expert judgment must conclude that *potential* benefits to patients clearly

NATIONAL QUALITY FORUM

outweigh *potential* harms to patients from the specific structure or process. Much of healthcare has not been subjected to research studies and thus does not have a strong evidence base. In the absence of strong evidence, structures or processes of care that are the focus of quality performance measures should be judged to provide benefits to patients that clearly outweigh any potential harms to patients.

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

Consistency with prior terminology, whenever possible, minimizes confusion. Terminology used in prior NQF documents should be changed only if it is incorrect or the change will lead to increased understanding. Whenever possible, narrative descriptions should be used instead of technical terminology.

I. Recommendations for Selecting the Focus for Measure Development

Based on its discussion and recommendations regarding evidence to support the measure focus, the Task Force made the following recommendations regarding selecting a focus for measure development.

- There is a hierarchical preference for outcome measures (when possible), followed by process measures, then structure measures. Outcome measures are preferred because improving health outcomes is a central goal of healthcare. However, both outcome and process measures have advantages and disadvantages²¹ and both have a place in quality assessment and the NQF portfolio.
- For process and structure measures, the focus of measurement should be on the aspect of care with the most direct evidence of a strong relationship to the desired outcome. For example, evidence about effective medication to control blood pressure is direct evidence

NATIONAL QUALITY FORUM

for the medication but only indirect evidence for the frequency of assessing blood pressure. Assessing blood pressure, although necessary, is not sufficient to achieving control. When there are multiple processes that affect a desired outcome, efforts should be made to include measures for all processes that have a strong relationship to the desired outcome.

- Specific drugs and devices included in quality performance measures should be Food and Drug Administration (FDA)-approved for the target condition.
- Structural measures are appropriate primarily when there are very well established structure-process-outcome relationships and when it is not feasible to directly measure the outcome or processes.
- For any topic area, measures based on the best evidence should be considered over measures based on lower quality evidence (e.g., expert opinion).

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

Following are the expectations for sources and grading of evidence used in support of measures submitted to NQF for potential endorsement. These recommendations are not intended to require measure developers to conduct primary reviews and grade the evidence. Rather the intent is to provide guidance on the expectations of evidence required to meet NQF's evaluation criteria.

- The preferred sources of evidence for quality performance measures are systematic reviews and grading of a body of evidence conducted by independent organizations such as [USPSTF](#), [AHRQ Evidence-based Practice Centers](#), and the [Cochrane Collaboration](#); or guidelines that meet national standards for trustworthy guidelines (as being developed by the IOM).
- Until such time when guidelines are certified as meeting a set standard, preferred guidelines are those developed with balanced representation beyond one specialty group and with full disclosure of biases and how they were addressed. Further, the evidence underlying a guideline recommendation must be accessible in order to provide the information necessary to meet the requirements set out in this report.

NATIONAL QUALITY FORUM

- An assigned evidence grade alone is not sufficient to evaluate whether the NQF subcriterion on evidence for the focus of measurement (1c) is met, either now or in the future. The specific information on the quantity, quality, and consistency of the body of evidence that was used to determine an overall grade should be summarized in the measure submission.
- Explicit, transparent information on the quantity, quality, and consistency of the body of evidence supporting a measure will facilitate identification of guideline recommendations that do not have acceptable evidence as the basis for performance measurement. Explicit information about the evidence also facilitates review by all stakeholders, although technical advisory panels and steering committees will continue to include experts that possess knowledge about the state of the science for a particular topic.

Current Expectations

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

- To promote transparency and standardization, NQF should require measure developers to provide specific information about the quantity, quality, and consistency of the body of evidence underlying a quality performance measure. Information should include who graded the evidence, the evidence grading system used, and the grade assigned. If the measure developer fails to provide this information, then NQF should not review the proposed measure.
- NQF prefers (but does not require) that submitted evidence be graded based on either the [USPSTF](#) or [GRADE](#) systems because such standardization facilitates broader understanding of the strength of the evidence.

Future Expectations

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

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

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

There has been widespread acceptance of Donabedian’s^{21, 22} structure-process-outcome model for assessing healthcare quality. These three approaches to measuring quality can be used with any healthcare topic, and the evidence required generally does not vary by topic. The required evidence to support the focus of a quality measure is for the links depicted by the red arrows in Figure 2. As shown under process, there are multiple steps prior to delivering an intervention; however, the evidence is most often about the relationship between the intervention and outcome, and, therefore, interventions are the preferred focus of process measures. Although

NATIONAL QUALITY FORUM

patient factors influence structures, processes, and outcomes, it is particularly important to consider those that influence outcomes for risk-adjustment of outcome measures.

Figure 2: Structure-Process-Outcome Model

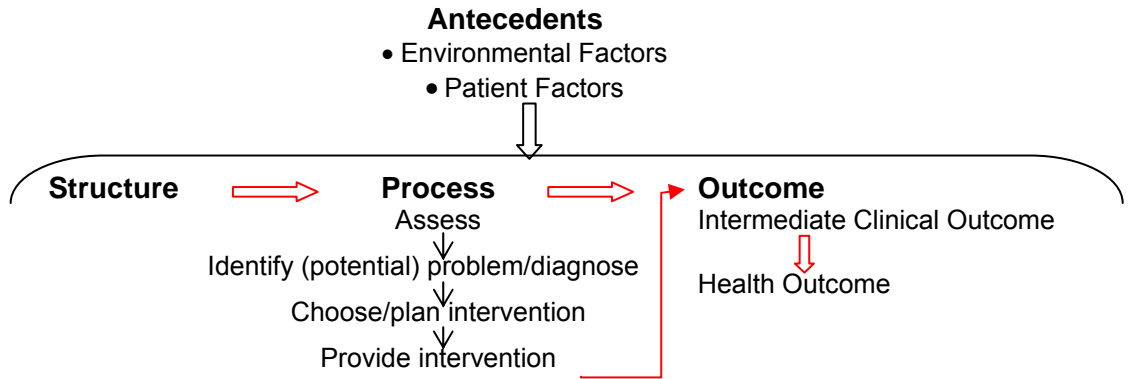


Table 3 outlines the evidence required to justify the structure, process, or outcome that is the focus of measurement. It also identifies special considerations related to certain quality topics. Subsequent tables lay out the approach for evaluating the evidence and using it to determine if the NQF subcriterion for evidence is met.

As noted by the Task Force and articulated by NQF’s Board, there is a preference for measures of health outcomes. Achieving or improving health outcomes is a central goal of healthcare treatments and services (e.g., health, function, survival, symptom control). Outcomes also are viewed as useful quality indicators because they integrate the influence of multiple care processes and disciplines involved in the care. Further, once they are measured and reported, many outcomes that were not thought to be modifiable tend to improve. This suggests that measurement stimulates identification and adoption of effective practices. Because multiple processes may influence a health outcome, several bodies of evidence could be relevant. For the reasons noted above, health outcomes do not necessarily require empirical evidence linking them to a known process or structure of care. Although such evidence is desirable, a rationale supporting the linkages between the measured health outcome and at least one healthcare structure, process, intervention, or service is sufficient.

NATIONAL QUALITY FORUM

Table 3: Evidence to Support the Focus of Measurement

Type of Measure	Evidence	Example of Measure Type and Evidence to Be Addressed
<p>Health Outcome An outcome of care is the health status 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 deterioration in health status).</p>	<p>A rationale supports the relationship of the health outcome to at least one healthcare structure, process, intervention, or service. See Table 5.</p>	<p>#0230 Acute myocardial Infarction 30-day mortality</p> <p>Survival is a goal of seeking and providing treatment for AMI.</p> <p>Rationale linking healthcare processes/ interventions (aspirin, reperfusion) to mortality/ survival</p> <p>#0171 Acute care hospitalization (risk-adjusted) [of home care patients]</p> <p>Improvement or stabilization of condition to remain at home is a goal of seeking and providing home care services.</p> <p>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</p> <p>Avoiding harm from treatment is a goal when seeking and providing healthcare.</p> <p>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]</p> <p>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 healthcare-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, then quantity, quality, and consistency of a body of evidence that the measured healthcare process does <i>not</i> lead to desired health outcomes in the target</p>	<p>#0551 ACE inhibitor/Angiotensin receptor blocker (ARB) use and persistence among members with coronary artery disease at high risk for coronary events</p> <p>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</p> <p>Evidence that antibiotics are not effective for acute bronchitis</p>

NATIONAL QUALITY FORUM

Type of Measure	Evidence	Example of Measure Type and Evidence to Be Addressed
	population. See Table 4.	
Structure Structure of care is a feature of a healthcare organization or clinician related to its capacity to provide high-quality healthcare.	Quantity, quality, and consistency of a body of evidence that the measured healthcare structure leads to desired health outcomes with benefits that outweigh harms (including evidence for the link to effective care processes and the link from the care processes to desired health outcomes). See Table 4.	#0190 Nurse staffing hours Evidence that higher nursing hours result in lower mortality or morbidity, or 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	<ul style="list-style-type: none"> • 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 <i>and</i> 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

IV. Recommendations for Evaluating Criterion 1c—Quantity, Quality, Consistency of Body of Evidence

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

- Evidence should be evaluated on the *quantity* of studies, *quality* of studies, and *consistency* in direction and magnitude of net benefit (clinically or practically meaningful)

NATIONAL QUALITY FORUM

benefits over harms to patients) of a *body of evidence* on a scale of High, Moderate, or Low.

- The dimensions of *quantity*, *quality*, and *consistency* of a body of evidence apply to measures based on guidelines as well as those for which guidelines may not exist (e.g., measures of care coordination or team functioning may not be based on guidelines, but often have bodies of evidence including nonclinical literature that should be systematically assessed).
- Measures without a clear description of the *quantity*, *quality*, and *consistency* of the supporting body of evidence or without any evidence should not pass subcriterion 1c and the threshold criterion of *Importance to Measure and Report*.
- Use of only selected studies rather than an entire body of evidence that meets pre-established criteria is not adequate to evaluate the evidence and should not pass subcriterion 1c and the threshold criterion of *Importance to Measure and Report*.
- Inconsistent and conflicting evidence should result in the measure not passing subcriterion 1c and the threshold criterion of *Importance to Measure and Report*.
- Outcome measures are considered an exception to the evidence requirement. A rationale should support the relationship of the outcome to processes of care and/or the importance of measuring the outcome.
- Expert opinion is not considered empirical evidence and will only be considered in exceptional circumstances when all of the following conditions are met.
 - No evidence is available.
 - Expert opinion is systematically assessed. That is, identified experts explicitly address the certainty or confidence that benefits to patients from the specific process or structure greatly outweigh potential harms, using a specified process that is transparent and open to peer review (e.g., modified Delphi, formal consensus process, [RAND Appropriateness Method](#)²³). The methods and results are reported for review.
 - There is a strong rationale for why the specific structure or process should be the focus of a quality performance measure.

NATIONAL QUALITY FORUM

Table 4 provides definitions and guidance on how to evaluate each of the dimensions of *quantity*, *quality*, and *consistency* for a body of quantitative evidence. Each dimension is rated on a scale of high, moderate, low, or inadequate to evaluate. A body of evidence could have different ratings for each dimension, e.g., high on quantity, low on quality, and moderate on consistency. Table 5 provides recommended decision rules for using the ratings for all three dimensions to make a decision on whether a measure should pass subcriterion 1c. Strong evidence usually requires multiple studies, each with sufficient numbers of patients to give precise estimates, but occasionally a large and representative study can provide adequate evidence. For example, one study (low quantity) that is a randomized controlled trial with a large representative sample of patients (high quality) and substantial estimates of net benefit would pass the subcriterion, whereas, a body of evidence with low consistency of estimates of net benefits should not pass the subcriterion regardless of the ratings for quantity and quality of studies.

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

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

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

Although RCTs remain the gold standard for evidence of efficacy of treatment, there are many areas where RCTs may not currently exist and are unlikely to be conducted. Furthermore, the strict eligibility and exclusion criteria for randomized trials may sometimes result in findings that are not fully generalizable in real-world applications. NQF recognizes the evidentiary value of well-conducted observational studies, particularly those that attempt to balance measured covariates (e.g., using propensity scores) and to account for other sources of bias as articulated in the [GRACE principles](#) (Good Research for Comparative Effectiveness).²⁴ This is particularly true when there are multiple observational studies that arrive at similar conclusions.

NATIONAL QUALITY FORUM

Qualitative studies often are used to gain understanding of people's attitudes, behaviors, and values and may be suited to evidence regarding patient experience with care. The descriptions of quality and consistency of the evidence in Table 4 do not apply to qualitative evidence. When qualitative studies are used, appropriate qualitative research criteria should be used to judge the strength of the evidence.²⁵

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

NATIONAL QUALITY FORUM

Table 4: Evaluation of Quantity, Quality, and Consistency of Body of Evidence for Structure, Process, and Intermediate Outcome Measures

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^a 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 ^b	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 ^b	<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, then the estimate of benefits greatly outweighs the estimate of potential harms to patients (one study cannot achieve high consistency rating)</p>
Low	0-1 studies ^b	<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 	<ul style="list-style-type: none"> • 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>If only one study, then estimate of benefits do not greatly outweigh harms to patients</p>
Insufficient to Evaluate (See Table 5 for exceptions.)	<ul style="list-style-type: none"> • No empirical evidence OR • Only selected studies from a larger body of evidence 	<ul style="list-style-type: none"> • 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

NATIONAL QUALITY FORUM

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

Study flaws that may bias estimates of effect include: lack of allocation concealment; lack of blinding; large losses to follow-up; failure to adhere to intention to treat analysis; stopping early for benefit; and failure to report important outcomes.

Imprecision with wide confidence intervals around estimates of effects can occur in studies involving few patients and few events.

Indirectness of evidence includes: indirect comparisons (e.g., two drugs compared to placebos rather than head-to-head); and differences between the population, intervention, comparator interventions, and outcome of interest and those included in the relevant studies.¹⁵

^bThe suggested number of studies for rating levels of quantity is considered a general guideline.

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 Results of Body of Evidence	Pass Subcriterion 1c
Moderate-High	Moderate-High	Moderate-High	Yes
Low	Moderate-High	Moderate (if only one study, high consistency not possible)	Yes, but only if it is judged 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 it is judged that potential benefits to patients clearly outweigh potential harms; otherwise, No
Low-Moderate-High	Low-Moderate-High	Low	No
Low	Low	Low	No
Exception to Empirical Body of Evidence for Health Outcome For a health outcome measure: A rationale supports the relationship of the health outcome to at least one healthcare structure, process, intervention, or service			Yes, if it is judged that the rationale supports the relationship of the health outcome to at least one healthcare structure, process, intervention, or service
Potential Exception to Empirical Body of Evidence for Other Types of Measures If there is no empirical evidence, expert opinion is systematically assessed with agreement that the benefits to patients greatly outweigh potential harms.			Yes, but only if it is judged that potential benefits to patients clearly outweigh potential harms; otherwise, No

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

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

NATIONAL QUALITY FORUM

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

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

The Task Force noted that, at the time of endorsement maintenance review, if measure performance data indicate overall high performance with little variation, then justification would be required for continued endorsement of the measure. The Consensus Standards Approval Committee (CSAC) added that the default action should be to remove endorsement unless there is a strong justification to continue endorsement. If a measure fails opportunity for improvement (1b), then it does not pass the threshold criterion, *Importance to Measure and Report*, and is therefore not suitable for endorsement. The CSAC noted that opportunity for improvement also could be considered during the review of measures with time-limited endorsement if there were enough data to make such a judgment.

NATIONAL QUALITY FORUM

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

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

- At the time of initial endorsement, evidence for opportunity for improvement generally will be based on research studies, or on epidemiologic or resource use data. However, at the time of review for endorsement maintenance, the primary interest is on the endorsed measure as specified, and the *evidence for opportunity for improvement should be based on data for the specific endorsed measure*.
- When assessing measure performance data for opportunity for improvement, the following factors should be considered:
 - number and representativeness of the entities included in the measure performance data; and
 - size of the population at risk, effectiveness of an intervention, likely occurrence of an outcome, and consequences of the quality problem.
- At the time of review for endorsement maintenance, an overall high level of performance with little variation in the endorsed measure scores should result in removal of endorsement. If other evidence (e.g., epidemiologic or research) is consistent with the measure performance data, then it confirms the lack of opportunity for improvement. If other evidence is not consistent with the measure performance data, then it is suggestive of potential problems with the measure as specified.
- In exceptional situations, a strong justification for continuing endorsement could be considered (e.g., *evidence* that overall performance will likely deteriorate if not monitored and of the magnitude of potential harm if outcomes deteriorate while not being monitored).

NATIONAL QUALITY FORUM

Table 6: Evidence for Evaluating Importance to Measure and Report

Pass Criterion, Importance to Measure and Report?			
All three 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)	<ul style="list-style-type: none"> Addresses a <i>specific national health goal/priority</i> identified by the Secretary of DHHS or the NPP <p>OR</p> <ul style="list-style-type: none"> 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 	<p>#0140 Ventilator-associated pneumonia for ICU and high-risk nursery (HRN) patients</p> <p>NPP goal: Focus relentlessly on continually reducing and seeking to eliminate all healthcare-associated infections (HAIs)</p> <p>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)</p>	<p>Yes— Demonstrated at least one of the aspects of high impact</p> <p>No—Did not demonstrate at least one of the aspects of high impact</p>
Opportunity for improvement (1b)	<p>Initial Endorsement Epidemiologic or resource use data or health services research demonstrating considerable variation or overall less than optimal performance for the focus of measurement across providers and/or population groups (disparities in care)</p> <p>Review for Endorsement Maintenance Data for the measure as specified and endorsed demonstrating considerable variation or overall less than optimal performance</p>	<p>#0432 Influenza vaccination of nursing home/skilled nursing facility residents</p> <p>NPP goal: All Americans will receive the most effective preventive services recommended by the U.S. Preventive Services Task Force.</p> <p>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)</p>	<p>Yes— Demonstrated either variation or overall less than optimal performance</p> <p>No—Did not demonstrate either variation or overall less than optimal performance</p>
Evidence for the focus of measurement (1c)	See Table 3	See Table 3	See Table 4 and Table 5

VI. Recommendations for Modifications to the NQF Evaluation Criteria

Table 7 presents modifications to the criteria to reflect the Task Force’s recommendations, including the recommendation that all three subcriteria must be met to pass the threshold criterion of *Importance to Measure and Report*. The Task Force identified that consequences of

NATIONAL QUALITY FORUM

measurement are not the same as the consequences of implementing the measured structure or process, that is, the benefits or harms to the patient related to the specific topic of measurement. Therefore, subcriterion 4d on susceptibility to inaccuracies, errors, or unintended consequences of measurement was not addressed in these recommendations related to evidence and *Importance to Measure and Report*.

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., data¹ demonstrating 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 evidence³ that supports the specific measure focus as follows: <ul style="list-style-type: none"> ○ <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. ○ <u>Process</u> – evidence that the measured clinical or administrative process leads to improved 	<p>1. Importance to measure and report: Extent to which the specific measure focus is evidence-based and important to making significant gains in healthcare 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., large numbers of patients affected 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 for improvement, i.e., data¹ 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 a health outcome or is evidence-based, demonstrated as follows:</p> <ul style="list-style-type: none"> • <u>health outcome</u>:² a rationale supports the relationship of the health outcome to processes or structures of care; OR • evidence-based as demonstrated by a systematic assessment and grading of the quantity, quality, and consistency of the body of evidence.³ <ul style="list-style-type: none"> ○ <u>Intermediate clinical outcome</u>: evidence that the measured intermediate clinical outcome leads to a desired health outcome. ○ <u>Process</u>:⁴ evidence that the measured healthcare

NATIONAL QUALITY FORUM

Current Measure Evaluation Criteria	Modified Measure Evaluation Criteria
<p>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).</p> <ul style="list-style-type: none"> ○ <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. ○ <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. ○ <u>Access</u> – evidence that an association exists between access to a health service and the outcomes of, or experience with, care. ○ <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</p> <p>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.</p> <p>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>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</p>	<p>process leads to desired health outcomes.</p> <ul style="list-style-type: none"> ○ <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). ○ <u>Efficiency</u>:⁵ evidence for the quality component as noted above. <p>OR</p> <ul style="list-style-type: none"> ● <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. <p>Footnotes</p> <p>¹Examples 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.</p> <p>²Generally, 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.</p> <p>³The preferred systems for grading the evidence are the USPSTF grading definitions and methods, or GRADE.</p> <p>⁴Clinical 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, then the step with the strongest evidence for the link to the desired outcome should be selected as the focus of measurement.</p> <p>⁵Measures of efficiency combine the concepts of resource use <i>and</i> quality (NQF’s Measurement Framework: Evaluating Efficiency Across Episodes of Care; AQA Principles of Efficiency Measures).</p>

NATIONAL QUALITY FORUM

Current Measure Evaluation Criteria	Modified Measure Evaluation Criteria
<p>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>	

VII. Recommendations for Modifications to the Measure Submission Form

The information requested on NQF’s measure submission form is consistent with that identified in a 2009 collaborative effort undertaken with AHRQ, the Centers for Medicare & Medicaid Services (CMS), The Joint Commission, National Committee for Quality Assurance (NCQA), and the Physician Consortium for Performance Improvement (PCPI) convened by the American Medical Association to identify common data fields. The Task Force suggested modifications to the information requested on the NQF [measure submission form](#) to implement the above recommendations. The intent is full transparency about the supporting evidence for the submitted measure. This will facilitate understanding of the adequacy of the evidence presented (selected evidence versus a body of evidence) and the measure developer’s representation of the quality of the evidence. Currently, evidence graded using the USPSTF or GRADE systems may not be available; however, an accurate description of the evidence and any grading system used should still be expected. The items in Table 8 pertain to the recommendations related to evidence (subcriterion 1c under *Importance to Measure and Report*).

NATIONAL QUALITY FORUM

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) <i>(for NQF staff use)</i> Specific NPP 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:</p> <p>1a.3. Summary of Evidence of High Impact</p> <p>1a.4. Citations for Evidence of High Impact</p> <p>Opportunity for Improvement (Measure evaluation criterion 1b) 1b.1. Briefly explain the benefits (improvements in quality) envisioned by use of this measure</p> <p>1b.2. Summary of Data Demonstrating Performance Gap <i>(Variation or overall poor performance across providers)</i></p> <p>1b.3. Citations for Data on Performance Gap</p> <p>1b.4. Summary of Data on Disparities by Population Group</p> <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</p>	<p>High Impact (Measure evaluation criterion 1a) <i>(for NQF staff use)</i> Specific priority goal: 1a.1. Demonstrated High-Impact Aspect of Healthcare Large numbers of patients affected Leading cause of morbidity/mortality Severity of illness Patient/societal consequences of poor quality Frequently performed procedure High resource use Other:</p> <p>1a.3. Summary of Evidence of High Impact <i>(Provide epidemiologic or resource use data.)</i></p> <p>1a.4. Citations for Evidence of High Impact</p> <p>Opportunity for Improvement (Measure evaluation subcriterion 1b) 1b.1. Briefly explain the benefits (improvements in quality) envisioned by use of this measure</p> <p>1b.2. Summary of Data Demonstrating Performance Gap <i>(variation or overall less than optimal performance across providers)</i></p> <p>1b.3. Citations for Data on Performance Gap</p> <p>1b.4. Summary of Data on Disparities by Population Group</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 and 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</p>

NATIONAL QUALITY FORUM

Current Measure Submission (4.1) Items	Modified Measure Submission Items
<p>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>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>1c.2. Source of Evidence Clinical practice guideline Systematic review of body of evidence (other than within guideline development) Selected individual studies (rather than entire body of evidence) Other (1c.3).</p> <p>1c.4. Summary of Body of Evidence Quantity of Studies in Body of Evidence <i>(total number of studies, not articles)</i></p> <p>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>Directness to focus of measurement and 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></p> <p>Consistency of Results Across Studies <i>(Summarize the consistency of the magnitude and direction of the effect.):</i></p> <p>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: 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</i></p>

NATIONAL QUALITY FORUM

Current Measure Submission (4.1) Items	Modified Measure Submission Items
<p>1c.7. Summary of Controversy/Contradictory Evidence</p> <p>1c.8. Citations for Evidence (<i>Other than guidelines</i>)</p> <p>1c.9. Quote 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><i>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? (<i>Rate each as High, Moderate, or Low.</i>) Quantity: Quality: Consistency:</p>

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

NQF also endorses practices such as [safe practices](#) and care coordination practices. The criteria for endorsing practices include evidence of effectiveness.²⁶ The Task Force recommends that the

NATIONAL QUALITY FORUM

same evidence requirements as indicated for process measures (Tables 3) be applied to practices considered for NQF endorsement (Table 9). Therefore, the rating categories for the quantity, quality and consistency of a body of evidence for quality measures as presented in Table 4 and the conclusion about the adequacy of the evidence presented in Table 5 also apply to rating the evidence for practices that are considered for NQF endorsement.

Table 9: Evidence to Support a Practice

Evidence to Support a Practice	Example of Practice and 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	<p>Safe Practice 16 Safe Adoption of Computerized Prescriber Order Entry</p> <p>Evidence that computerized order entry systems are associated with lower medication errors and adverse events</p>

Table 10 presents modifications to the evidence criterion to reflect the Task Force’s recommendations. The other criteria used to evaluate practices for potential endorsement (specificity, benefit, generalizability, and readiness) were not addressed in this project.

Table 10: Current and Modified Practice Evaluation Evidence Criterion

Current Practice Evaluation Criteria	Modified Practice Evaluation Criteria
<p>Evidence of Effectiveness. There must be clear evidence that the practice would be effective in reducing patient safety events. Such evidence may take various forms, including the following:</p> <ul style="list-style-type: none"> • Research studies showing a direct connection between improved clinical outcomes (e.g., reduced mortality or morbidity) and the practice; • experiential data (including broad expert agreement, widespread opinion, or professional consensus) showing the practice is "obviously beneficial" or self-evident (i.e., the practice absolutely constrains a potential problem or forces an improvement to occur, reduces reliance on memory, standardizes equipment or process steps, or promotes teamwork); or • Research findings or experiential data from non-healthcare industries that should be substantially transferable to healthcare (e.g., repeat-back of verbal orders or standardizing abbreviations). 	<p>Evidence of Effectiveness. A practice is evidence-based as demonstrated by a systematic assessment of the quantity, quality, and consistency of the body of evidence and standardized grading of the body of evidence. The preferred systems for grading the evidence are the USPSTF grading definitions and methods, or GRADE. Evidence from non-healthcare industries that should be substantially transferable to healthcare (e.g., safety practices of repeat-back of verbal orders or standardizing abbreviations) also may be considered.</p>

NOTES

1. Lohr KN, Rating the strength of scientific evidence: relevance for quality improvement programs, *Int J Qual Health Care*, 2004;16(1):9-18.
2. Tricoci P, Allen JM, Kramer JM, et al., Scientific evidence underlying the ACC/AHA clinical practice guidelines, *JAMA*, 2009;301(8):831-841.
3. Spertus JA, Eagle KA, Krumholz HM, et al., American College of Cardiology and American Heart Association methodology for the selection and creation of performance measures for quantifying the quality of cardiovascular care, *Circulation*, 2005;111(13):1703-1712.
4. Physician Consortium for Performance Improvement, *Physician Consortium for Performance Improvement® (PCPI) Position Statement - The Evidence Base Required for Measures Development*, Chicago, IL: American Medical Association, 2009. Available at www.ama-assn.org/ama/pub/physician-resources/clinical-practice-improvement/clinical-quality/physician-consortium-performance-improvement/position-papers.shtml. Last accessed March 2010.
5. Grilli R, Magrini N, Penna A, et al., Practice guidelines developed by specialty societies: the need for a critical appraisal, *Lancet*, 2000;355(9198):103-106.
6. Shiffman RN, Shekelle P, Overhage JM, et al., Standardized reporting of clinical practice guidelines: a proposal from the Conference on Guideline Standardization, *Ann Intern Med*, 2003;139(6):493-498.
7. The AGREE Collaboration, *Appraisal of Guidelines for Research and Evaluation AGREE Instrument*, 2001. Available at www.agreecollaboration.org/instrument. Last accessed February 2010.
8. The AGREE Collaboration, Development and validation of an international appraisal instrument for assessing the quality of clinical practice guidelines: the AGREE project, *Qual Saf Health Care*, 2003;12(1):18-23.
9. Coates VH, "National Guideline Clearinghouse" Presented at Workshop on Standards for Clinical Practice Guidelines, Institute of Medicine, Washington, DC; January 11, 2010. Available at <http://iom.edu/Activities/Quality/ClinicPracGuide/2010-JAN-11.aspx>. Last accessed February 2010.
10. Atkins D, Eccles M, Flottorp S, et al., Systems for grading the quality of evidence and the strength of recommendations I: critical appraisal of existing approaches The GRADE Working Group, *BMC Health Serv Res*, 2004;4(1):38.
11. Atkins D, Best D, Briss PA, et al., Grading quality of evidence and strength of recommendations, *BMJ*, 2004;328(7454):1490-1494.
12. Guyatt GH, Oxman AD, Vist GE, et al., GRADE: an emerging consensus on rating quality of evidence and strength of recommendations, *BMJ*, 2008;336(7650):924-926.
13. Guyatt GH, Oxman AD, Kunz R, et al., Incorporating considerations of resources use into grading recommendations, *BMJ*, 2008;336(7654):1170-1173.
14. Guyatt GH, Oxman AD, Kunz R, et al., Going from evidence to recommendations, *BMJ*, 2008;336(7652):1049-1051.
15. Guyatt GH, Oxman AD, Kunz R, et al., What is "quality of evidence" and why is it important to clinicians?, *BMJ*, 2008;336(7651):995-998.
16. Guyatt GH, Oxman AD, Vist GE, et al., GRADE: an emerging consensus on rating quality of evidence and strength of recommendations, *BMJ*, 2008;336(7650):924-926.

NATIONAL QUALITY FORUM

17. Harris RP, Helfand M, Woolf SH, et al., Current methods of the US Preventive Services Task Force: a review of the process, *Am J Prev Med*, 2001;20(3 Suppl):21-35.
18. Sawaya GF, Guirguis-Blake J, LeFevre M, et al., Update on the methods of the U.S. Preventive Services Task Force: estimating certainty and magnitude of net benefit, *Ann Intern Med*, 2007;147(12):871-875.
19. Owens DK, Lohr KN, Atkins D, et al., AHRQ series paper 5: grading the strength of a body of evidence when comparing medical interventions--agency for healthcare research and quality and the effective health-care program, *J Clin Epidemiol*, 2010;63(5):513-523.
20. Liberati A, Altman DG, Tetzlaff J, et al., The PRISMA statement for reporting systematic reviews and meta-analyses of studies that evaluate health care interventions: explanation and elaboration, *Ann Intern Med*, 2009;151(4):W65-W94.
21. Donabedian A, The role of outcomes in quality assessment and assurance, *Qual Rev Bull*, 1992;18(11):356-360.
22. Donabedian A, *An Introduction to Quality Assurance in Health Care*, New York, NY: Oxford University Press, 2003.
23. Fitch K, Bernstein SJ, Aguilar MS, et al., *The RAND/UCLA Appropriateness Method User's Manual*, Santa Monica, CA: RAND Health, 2000. Available at www.rand.org/pubs/monograph_reports/MR1269.html. Last accessed November 2010.
24. Dreyer NA, Schneeweiss S, McNeil BJ, et al., GRACE principles: recognizing high-quality observational studies of comparative effectiveness, *Am J Manag Care*, 2010;16(6):467-471.
25. Cohen DJ, Crabtree BF, Evaluative criteria for qualitative research in health care: controversies and recommendations, *Ann Fam Med*, 2008;6(4):331-339.
26. The National Quality Forum (NQF), *Safe Practices for Better Healthcare*, Washington, DC: NQF, 2009. Available at www.qualityforum.org/News_And_Resources/Press_Kits/Safe_Practices_for_Better_Health_care.aspx. Last accessed February 2010.

APPENDIX A EVALUATION CRITERIA FOR MEASURES AND PRACTICES

Evaluation Criteria for Measures (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).

NATIONAL QUALITY FORUM

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 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:
 - Intermediate outcome – evidence that the measured intermediate outcome (e.g., blood pressure, HbA1c) leads to improved health/avoidance of harm or cost/benefit.
 - 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).
 - 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.
 - 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.
 - Access – evidence that an association exists between access to a health service and

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.

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.

NATIONAL QUALITY FORUM

the outcomes of, or experience with, care.

- 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.

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:

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), 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.

NATIONAL QUALITY FORUM

- 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}

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,

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 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.

NATIONAL QUALITY FORUM

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 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.

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.

NATIONAL QUALITY FORUM

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).

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.

NATIONAL QUALITY FORUM

Prior Evaluation Criteria for Practices²⁶

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

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

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

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

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

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

NATIONAL QUALITY FORUM

APPENDIX B TASK FORCE MEMBERS

David Shahian, MD (Chair)

Chair, Society of Thoracic Surgeons Workforce on National Databases
Center for Quality and Safety and Department of Surgery, Massachusetts General Hospital;
Professor of Surgery, Harvard Medical School

Kristine Martin Anderson, MBA (CSAC member)

Senior Vice President
Booz Allen Hamilton
Rockville, MD

David Atkins, MD, MPH

Director of Quality Enhancement Research Initiative (QUERI)
Department of Veterans Affairs, Health Services Research & Development Service
Washington, D.C.

Arthur Levin, MPH (CSAC member)

Director
Center for Medical Consumers
New York, NY

Mary Naylor, PhD, RN (Board Member)

Marian S. Ware Professor in Gerontology
University of Pennsylvania School of Nursing
Philadelphia, PA

Greg Pawlson, MD, MPH

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

Eric Schneider, MD, MSc, FACP

Senior Scientist and Director, RAND Boston; Associate Professor, Division of General Medicine and Primary Care, Brigham and Women's Hospital; Associate Professor, Department of Health Policy and Management, Harvard School of Public Health
Boston, MA

NATIONAL QUALITY FORUM

**APPENDIX C
U.S. PREVENTIVE SERVICES TASK FORCE SYSTEM FOR GRADING EVIDENCE
AND RECOMMENDATIONS**

The following tables are found on the web page [Update on Methods](#) for estimating certainty and magnitude of net benefit by the U.S. Preventive Services Task Force (USPSTF).

Table 1: U.S. Preventive Services Task Force Recommendation Grid*

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

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

Table 2: Questions Considered by the U.S. Preventive Services Task Force for Evaluating Evidence Related 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)? |
|---|

NATIONAL QUALITY FORUM

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: 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: 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.

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

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

NATIONAL QUALITY FORUM

The following table is from the USPSTF web page [Grade Definitions](#).

Table 5. What the Grades Mean and Suggestions for Practice

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.