



TO: Consensus Standards Approval Committee (CSAC)  
Health IT Advisory Committee (HITAC)

FR: Helen Burstin, Karen Pace, and Christopher Millet

RE: Comments Received on Draft Report: *Review and Update of Guidance for Evaluating Evidence and Measure Testing*

DA: September 6, 2013

The CSAC and HITAC will review comments on the draft report: *Review and Update of Guidance for Evaluating Evidence and Measure Testing* during a joint conference call on September 10, 2013. The draft report was posted for a 21-day comment period from August 8 through August 30, 2013.

This memo includes background of the project, the major themes within the comments, and specific issues that require further CSAC/HITAC action.

#### AGENDA

- Review the themes in the order presented in this memo
- Public comment
- Helen Burstin will give a Kaizen report for the day

The following documents support this memo:

1. [Review and Update of Guidance for Evaluating Evidence and Measure Testing](#). This is the draft report that was posted for comment.
2. **Comment table**. This lists 128 comments received from 16 organizations and includes draft responses for your consideration (posted on CSAC SharePoint and distributed to HITAC).

#### CSAC/HITAC ACTION REQUIRED

- Review this briefing memo and be prepared to provide feedback and input on identified issues.
- Review the comments received and the proposed draft responses (see Excel file). We will not be able to review each individual comment on the call. Please notify us if you have suggested revisions to any of the responses or wish to discuss any specific comments on the call.
- Review the draft report as needed to further understand the comments or responses.

#### BACKGROUND

In 2010, the NQF convened two task forces to help provide guidance for evaluating the clinical evidence and measure testing for reliability and validity that is submitted in support of a performance measure being considered for endorsement. The approved recommendations were implemented in 2011. Testing of eMeasures also was addressed in the 2011 guidance and in some subsequent draft policy statements. Some challenges and inconsistencies in applying the guidance have been identified.



The purpose of this project was to review the implementation of the 2011 guidance on evaluating evidence and measure testing (including eMeasure testing requirements) and to propose modifications to address any major challenges. Modifications that would potentially increase consistency and clarity in the evaluation of performance measures for potential NQF endorsement also were considered.

The specific goals of the project included:

- promote consistency in evaluation across measures and projects;
- clarify common misunderstandings about the criteria and guidance;
- remain consistent with the criteria and principles from the 2011 guidance (i.e., do not change the “bar” for endorsement or the information requested for a measure submission); and
- address the current challenges with eMeasure testing.

The Consensus Standards Approval Committee (CSAC), Health IT Advisory Committee (HITAC), and subcommittees of both groups worked with NQF staff from March to August 2013 to consider the issues and propose potential solutions. The draft guidance was posted for public and member comment.

### **MAJOR THEMES AND ISSUES FOR DISCUSSION**

Of the 16 organizations that submitted comments, 14 were measure stewards or developers. The comments were very useful in identifying areas that require further clarification or action by the CSAC and HITAC.

Overall the comments were favorable to the proposed algorithms in terms of promoting greater consistency. Several commenters noted that steering committees need training and or support to implement the guidance and examples would be useful to measure developers. If approved, NQF will explain these algorithms as part of steering committee orientation activities. In addition, NQF staff will use this guidance to provide a first review of measure submissions to assist the steering committee. However, steering committees will continue to have responsibility for rating the criteria and making recommendations regarding endorsement. NQF attempts to include methodologists or statisticians on steering committees when possible and will continue to explore other mechanisms for committees to request statistical support. NQF also is providing some examples of [what good looks](#) like in terms of responding to questions on the measure submission form.

The major themes and issues for discussion and action are listed below. However, you also may identify additional areas for discussion from the comments and draft responses.

#### ***Theme 1: eMeasure endorsement - lack of support for required testing in 3 EHRs, each with 3 sites***

Commenters did not agree with setting this as a minimum, citing burden with finding sites and costs. ONC noted the difficulty it has encountered with trying to implement this requirement. One commenter suggested requiring 3 EHRs, but no minimum number of sites. If adopted, it would represent a higher bar than is currently required for other types of measures and would need to be phased in to allow developers time to adjust their procedures.

#### **ACTION ITEMS:**

- Should NQF specify a required number of EHR systems for testing? The HITAC also considered requiring 2 EHRs rather than 3. The draft 2012 policy did not specify a required number of EHR systems for testing and the feasibility assessment requirements referred to assessing feasibility in “multiple” EHRs (i.e., more than one).

### **Theme 2: eMeasure endorsement- need to confirm guidance**

Some commenters raised important issues or questions that would benefit from CSAC and HITAC specifically confirming the proposed guidance.

- Multiple versions of the same measure (eMeasure and some other data specification). Some thought that NQF should only endorse one version – the best way to measure performance. Others wanted to preserve older non-eMeasure versions.
- Some questioned if the XML cannot be directly applied to most EHRs, what are the expectations related to “testing a performance measure as specified” for eMeasures.
- Someone questioned how this guidance will apply to “retooled” measures.
- Some questioned whether identifying HQMF, QDM and VSAC value sets was too prescriptive given that standards may change and they may not accommodate all aspects of performance measures (also whether the VSAC is accepting new value sets).

#### **ACTION ITEMS:**

- Review and confirm the following. Changes from the draft report are redlined.

### **Endorsing eMeasures**

The following is a consolidation and clarification of the requirements for testing eMeasures that are submitted to NQF for endorsement (initial or endorsement maintenance). These requirements would apply to both new (de novo) eMeasures and previously endorsed measures (retooled).

- eMeasures must be specified in the accepted standard of HQMF format, and must use the Quality Data Model (QDM) and value sets vetted through the National Library of Medicine’s Value Set Authority Center (VSAC). Output from the Measure Authoring Tool (MAT) ensures that an eMeasure is in the HQMF format and uses the QDM (however, the MAT is not required to produce HQMF). Alternate forms of “e-specifications” other than HQMF are not considered eMeasures. *However, if HQMF or QDM does not support all aspects of a particular measure construct, those may be specified outside HQMF with an explanation and plans to request expansion of those standards. If particular value sets are not vetted by the VSAC, explain why they are used in the measure and describe plans to submit them to VSAC for approval. Please contact NQF staff to discuss format for measure specifications that the standards do not support.*
- A new requirement for a feasibility assessment will be implemented with projects beginning after July 1, 2013 (see the [eMeasure Feasibility Report](#)). The feasibility assessment addresses the data elements as well as the measure logic. (See Appendix C for feasibility criteria and example scorecard).
- All measures (including eMeasures) are subject to meeting the same evaluation criteria that are current at the time of initial submission or endorsement maintenance (regardless of meeting prior criteria and prior endorsement status). Algorithms 1, 2, and 3 apply to eMeasures.
  - Importance to Measure and Report (clinical evidence, performance gap, priority)
  - Scientific Acceptability of Measure Properties (reliability, validity)
  - Feasibility
  - Usability and Use (accountability/transparency, improvement)
  - Related and competing measures

- To be considered for NQF endorsement, All-all measures (including eMeasures) must be tested for reliability and validity using the data source that is specified. Therefore, eMeasures, whether new (de novo), previously respecified (retooled) but without eMeasure testing, or newly respecified, must be submitted with testing using the eMeasure specifications with the specified data source (e.g., EHRs, registry).
  - In the information provided on the data used for testing, indicate-describe how the eMeasure specifications were used to obtain the **electronic data used to compute the performance measure**. Often eMeasures cannot be directly applied to EHRs or databases from EHRs and additional programming is needed to identify the location of standardized data elements. However, in some instances, the eMeasure specifications might be used directly with EHRs.
- If testing of eMeasures occurs in a small number of sites, it may be best accomplished by focusing on patient-level data element validity (comparing data used in the measure to the authoritative source). However, as with other measures, testing at the level of the performance score is acceptable if data can be obtained from enough measured entities. The use of EHRs and the potential access to robust clinical data provides opportunities for other approaches to testing.
  - If the testing is focused on validating the accuracy of the electronic data, analyze agreement between the electronic data obtained using the eMeasure specifications and those obtained through abstraction of the entire electronic record (not just the fields used to obtain the electronic data), using statistical analysis such as sensitivity and specificity, positive predictive value, negative predictive value. The guidance on measure testing allows this type of validity testing to also satisfy reliability of patient-level data elements (see Algorithms 2 and 3).
  - Note that testing at the level of data elements requires that all critical data elements be tested (not just agreement of one final overall computation for all patients). A—at a minimum, the numerator, denominator, and exclusions (sometimes referred to as exceptions) must be assessed and reported separately.
  - Use of a simulated data set is no longer suggested for testing validity of data elements and is best suited for checking that the measure specifications and logic are working as intended.
  - NQF’s guidance has some flexibility; therefore, measure developers should consult with NQF staff if they think they have another reasonable approach to testing reliability and validity.
- For eMeasures, the sample for testing the patient-level data used in constructing the eMeasures should include a **minimum of two different EHR systems**. The general guidance on samples for testing all measures includes:
  - Testing may be conducted on a sample of the accountable entities (e.g., hospital, physician). The analytic unit specified for the particular measure (e.g., physician, hospital, home health agency) determines the sampling strategy for scientific acceptability testing.
  - The sample should represent the variety of entities whose performance will be measured. The Measure Testing Task Force recognized that the samples used for reliability and validity testing often have limited generalizability because measured entities volunteer to participate. Ideally, however, all types of entities whose performance will be measured should be included in reliability and validity testing.
  - The sample should include adequate numbers of units of measurement and adequate numbers of patients to answer the specific reliability or validity question with the chosen statistical method.

- When possible, units of measurement and patients within units should be randomly selected.
    - Is there any specific guidance regarding selecting the systems for testing (e.g., not customized, market share,?)?
- The following subcriteria under Scientific Acceptability of Measure Properties also apply to eMeasures.
  - Exclusion analysis (2b3). If exclusions (sometimes referred to as exceptions) are not based on the clinical evidence, analyses should identify the overall frequency of occurrence of the exclusions as well as variability across the measured entities to demonstrate the need to specify exclusions.
  - Risk adjustment (2b4). Outcome and resource use measures require testing of the risk adjustment approach.
  - Differences in performance (2b5). This criterion is about using the measure as specified to distinguish differences in performance across the entities that are being measured. The performance scores should be computed for all accountable entities for which you have eMeasure data are available (not just those on which reliability/validity testing was conducted) and then analyzed to identify differences in performance.
  - Comparability of performance scores if specified for multiple data sources (2b6) (e.g., EHRs, claims). If a performance measure is specified for more than one data source, it should be tested with each. The assumption is that measures specified for different data sources do not produce comparable results unless empirical analyses demonstrate comparability of computed scores and should be submitted as separate measures. Measures are endorsed only for the data specifications and levels of analyses for which they are tested. The measures specified for different data sources and the same levels of analysis will be evaluated as competing measures to determine whether one is superior to the other or whether there is justification for endorsing multiple versions of the same measure concepts.
  - Note: NQF needs to address how to number or relate this type of related measures.
  - Analysis of missing data (2b7). Approved recommendations from the 2012 projects on eMeasure feasibility assessment, composites, and patient-reported outcomes call for an assessment of missing data or nonresponses.

**Theme 3: Support for approval for trial use, but need further specificity**

The comments were generally favorable, but questions and suggestions indicated the need for further specificity regarding:

- Confirmation that approval would be based on a multistakeholder process (i.e., the CDP)
- Expiration of approval
- What is evaluated when submitted for trial use and then later for endorsement?
- Use of time-limited endorsement for eMeasures or trial use for all measures
- Flexible timing for submitting for endorsement

One commenter suggested that approval for trial use should not require specifications in HQMF format; others either agreed or did not comment.

One commenter suggested “approval for trial use” may imply it could be used in accountability applications and suggested renaming: “NQF Recommended for Implementation Testing.”

One commenter suggested that if adopted, NQF review all eMeasures to determine if they were tested as eMeasures and if not, reclassify as Recommended for Trial Implementation and Testing.

**ACTION ITEMS:**

- Discuss and approve the following criteria and processes and decide what to name this status. Changes from the draft report are redlined.
- Discuss whether Trial Implementation and Testing should be available to all measures and totally eliminate time-limited endorsement.
  - If so, should the current criteria for time-limited endorsement be retained for non-eMeasures? (An incumbent measure does not address the specific topic of interest in the proposed measure; a critical time line must be met (e.g., legislative mandate); t The measure is not complex (e.g., composite, require risk adjustment);
- Should NQF retro-actively apply this new designation to previously endorsed eMeasures as appropriate? Specifically, if they were not tested as eMeasures, then reclassify them as Recommended for Trial Implementation and Testing?

**Approve Trial Measures for Implementation and Testing (refer to as trial measures)**

*Approval of Trial Measures for Implementation and Testing* means the eMeasure has been judged to meet the criteria ~~indicating that indicate its readiness for implementation~~~~it is ready to be implemented~~ in real-world settings in order to generate the data required to assess reliability and validity. ~~Such measures~~ also could be used for internal performance improvement. However, ~~such measures would not have~~ ~~has not yet~~ been judged to meet all the criteria indicating it is suitable for use in accountability applications.

- Such measures will be considered Approved ~~for as~~ Trial Measures for Use Implementation and Testing, **NOT** Endorsed
- When sufficient data have been accumulated for adequate reliability and validity testing, the eMeasure can be submitted to NQF for potential endorsement (not all may progress to endorsement).

**Criteria for Approval of Trial Measures for Implementation and Testing**

- Must be eMeasures, meaning the measures are specified in the accepted standard of HQMF format, and must use the Quality Data Model (QDM). Output from the Measure Authoring Tool (MAT) ensures that an eMeasure is in the HQMF format and uses the QDM (however, the MAT is not required to produce HQMF). Alternate forms of “e-specifications” other than HQMF are not considered eMeasures. *However, if HQMF or QDM does not support all aspects of a particular measure construct, those may be specified outside HQMF. Please contact NQF staff to discuss format for measure specifications.*
- Must use value sets vetted through the National Library of Medicine’s Value Set Authority Center. This will help ensure appropriate use of codes and code systems and will help minimize value set harmonization issues in submitted eMeasures. If particular value sets are not vetted by the VSAC, explain why they are used in the measure and describe plans to submit them to VSAC for approval.
- Must meet all criteria under Importance to Measure and Report (clinical evidence, performance gap, priority).
- The feasibility assessment must be completed.

- Results from testing with a simulated (or test) data set demonstrate that the QDM and HQMF are used appropriately and that the measure logic performs as expected.
- There is a plan for use and discussion of how the measure will be useful for accountability and improvement.
- Related and competing measures are identified with a plan for harmonization or justification for developing a competing measure.

### Process

- Measures submitted for approval as Trial Measures for Implementation and Testing will be submitted and evaluated as other measures that are submitted for endorsement (i.e., using the multistakeholder CDP)
- Measures submitted as Trial Measures will be evaluated on Importance to Measure and Report, eMeasure specifications, Feasibility, Usability and Use, and Related and Competing Measures
- Trial Measure designation automatically expires 3 years after approved if not submitted for endorsement prior to that time.
  - The time to submit for endorsement is driven by success with testing. There is no expectation that every trial measure will be submitted for endorsement – some may fail during testing.
  - Should there be a process to extend status as trial measure? (for example, if a developer can show that they are actively engaged in the process of testing?)
- When submitted for endorsement, it will be evaluated through the multistakeholder process. Ideally, standing committees and/or more flexible schedules for submitting measures will prevent delays for the endorsement process.
- If submitted for endorsement prior to the 3-year expiration, consider the following options for evaluation and endorsement:
  - Option 1: Submit and evaluate only Scientific Acceptability of Measure Properties, including the final eMeasure specifications and all testing. Grant endorsement that will need to be maintained **3 years from the approval as trial measure**, at which time it will be submitted for endorsement maintenance and subject to evaluation on all criteria.
  - Option 2: Submit and evaluate all criteria. Grant endorsement which will need to be maintained 3 years from the endorsement date, at which time it will be submitted for endorsement maintenance and subject to evaluation on all criteria.
- If submitted for endorsement three or more years after the approval as a trial measure, the measure will be submitted and evaluated on all criteria just as any measure being submitted for initial endorsement or maintenance.

### ***Theme 4: Different views on evidence related to health outcome/patient-reported Outcomes (PROs)***

The comments reflect the differences in perspectives about evidence requirements for measures of health outcomes and PROs. Most commenters did not mention the issue of evidence for outcome measures; two commenters specifically agreed with the current approach (comments #33, 60); and three commenters suggested that measures of outcomes have the same evidence requirements or additional empirical analysis (comments # 16, 29, 44). Some commenters asked for greater clarity or specificity to evaluate a “plausible” rationale.

The CSAC recently reaffirmed the current criteria and guidance that requires a rationale that supports that the outcome is influenced by at least one healthcare structure, process, intervention, or service. The updated guidance was intended to reflect the current criteria, not change or “raise the bar.”

**ACTION ITEMS:**

- Discuss and finalize revision of Algorithm 1 (see suggestion below).
- Reaffirm that this new guidance will be implemented immediately, pending Board approval.
- Decide whether to revisit the evidence requirement for health outcomes and PROs and if so, discuss the timing and process for this reconsideration.

<p>1. Does the measure assess performance on a <b>health outcome</b> (e.g., mortality, function, <b>health status</b>, complication) or <b>PRO</b> (e.g., function, symptom, experience, <b>health-related behavior</b>)?</p>	<p>→ Yes →→→→→→→→→→</p>	<p>2. Does the SC agree that the relationship between the measured outcome/PRO and <b>at least one</b> healthcare action (structure, process, intervention, or service) is supported by: *the stated rationale; and *the diagram or description of the path between the outcome and healthcare processes or structures? Is there at least one healthcare process, intervention, service, or structure identified as influencing the outcome with a plausible rationale?</p> <p>→ No →→→→→→→→→→→→→→→</p>	<p>→ Yes          →</p>	<p><b>Pass</b>          <b>No pass</b></p>
---	-----------------------------	---	---	--

***Theme 5: Need for further clarification regarding guidance when information on quantity, quality, and consistency of the evidence is not available***

Most of the comments were favorable to the proposed algorithm. However, some comments indicated the need for further clarification. One commenter thought the suggested approach for when a summary of quantity, quality, and consistency is not provided accommodates poor measure submissions. One commenter asked whether the guidance means that a summary of quantity, quality, and consistency is not required.

**ACTION ITEMS:**

- Discuss and finalize a response:  
The current criteria and guidance require a summary of the quantity, quality, and consistency of the body of evidence and this updated guidance is not intended to change that requirement. Transparency about the evidence that does and does not exist is a core principle upon which the 2011 evidence task force recommendations were made and approved. Because of the number and variety of evidence review and grading systems, the need for a summary of the evidence still exists and is still required. The proposed guidance is intended to assist steering committees to implement the current guidance when developers identify that the information from the systematic review of the evidence is not available to provide the requested summary of quantity, quality, and consistency of the



evidence. Although the guidance cannot differentiate the reason for not providing the summary of quantity, quality, and consistency (e.g., not available, choose not to submit), it is important to note that without that information, moderate is the highest possible rating and that rating could apply only if the evidence grade/definition indicated that the evidence is high quality or results in high certainty. Without a summary of the quantity, quality, and consistency of the body of evidence, a grade indicating moderate quality would not pass the evidence subcriterion).

The language in the algorithm will be revised as follows.

<p>4. Is a summary of the quantity, quality, and consistency (QQC) of the body of evidence from a systematic review (SR) provided in the submission form?</p> <p>A SR is a scientific investigation that focuses on a specific question and uses explicit, pre-specified scientific methods to identify, select, assess, and summarize the findings of similar but separate studies. It may include a quantitative synthesis (meta-analysis), depending on the available data. (IOM)</p> <p><b>Answer NO if:</b>          *QQC not submitted (even if available)          *Details of SR not available (i.e., QQC, evidence tables)          *Specific information on QQC not available (i.e., general statements/conclusions, lists/descriptions of individual studies are not sufficient)</p>	→ Yes		
↓ No - without QQC from SR, moderate is highest potential rating			
<p>6. Does the grade for the evidence or recommendation indicate:</p> <p>*High quality evidence (e.g., Table 1 - Quant:Mod/Hi; Qual:Hi; Consist:Hi; USPSTF - High certainty; GRADE-High quality)          *Strong recommendation (e.g., GRADE -Strong; USPSTF-A)</p> <p><b>Answer NO if:</b>          *No grading of evidence and summary of QQC not provided          *Not graded high quality or strong recommendation</p>	→ Yes	Rate as Moderate	
→No (moderate/weak quality or recommendation without QQC) → →		→	Rate Low

**Theme 6: Difference of opinions about exceptions to the evidence**

Most commenters agreed that the guidance would help promote greater consistency. However, some thought the proposed guidance is still too subjective and asked for more specificity. Some commenters were concerned that the proposed guidance will not allow them to submit the measures they want for a specialty area where there is limited evidence.

**ACTION ITEMS:**

- Discuss and finalize an addition to the algorithm to address situations when empirical evidence exists but has not been systematically assessed (see below). In this case, a moderate rating would be the highest potential rating and would depend on the Steering

Committee assessment that the available evidence is of high or moderate quality and indicates a high certainty that benefits clearly outweigh undesirable effects.

- Discuss and finalize a response regarding exceptions:  
The question in Box 7 generally refers to measuring distal care processes (such as assessing blood pressure) when alternative processes or outcomes could be measured instead (e.g., BP control—an intermediate clinical outcome or effective treatment—a particular evidence-based process that is proximal to a desired outcome). The absence of performance measures for a particular condition, setting, or provider does not alter the criteria for endorsement. For example, although assessment is necessary, it alone is not sufficient to improve outcomes and often the assessment component can be incorporated into a process or outcome measure. Developers will need to explain why health outcomes cannot be measured or measures of intermediate clinical outcomes or interventions will not meet the evidence requirement. Note that in the case of high quality evidence for assessment or screening recommendations (such as those developed by the USPSTF), a performance measure of such a distal process could meet the evidence criterion.

<p>Is empirical evidence submitted but without systematic review and grading of the evidence?</p>	<p>→ <b>Yes</b></p>	<p>Does the empirical evidence that is summarized include all studies in the body of evidence?</p> <p><b>Answer NO if only selected studies included</b></p>	<p>→ <b>Yes</b></p>	<p>Does the SC agree that the submitted evidence indicates high certainty that benefits clearly outweigh undesirable effects? <i>(without SR, the evidence should be high-moderate quality and indicate substantial net benefit - See Table 1)</i></p>	<p>→ <b>Yes</b></p>	<p><b>Rate as Moderate</b></p>
<p>↓ No</p>		<p>No →</p>	<p>→</p>	<p>No →</p>	<p>→</p>	<p><b>Rate Low</b></p>
<p>7. Are there, OR could there be, performance measures of a related health outcome, OR evidence-based intermediate clinical outcome or process?</p> <p><i>Example for YES: Propose to measure whether BP is assessed each visit instead of BP control or use of effective treatment</i></p>	<p>→ <b>No</b></p>	<p>8. Is there evidence of a systematic assessment of expert opinion (e.g., national/international consensus recommendation) that the benefits of what is being measured outweigh potential harms?</p>	<p>→ <b>Yes</b></p>	<p>9. Does the SC agree that it is okay (or beneficial) to hold providers accountable for performance in the absence of empirical evidence of benefits to patients? <i>Are there any perceived detriments to endorsing the measure? (e.g., focus attention away from more impactful practices, more costly without certainty of benefit; divert resources from developing more impactful measures)</i></p>	<p>→ <b>Yes</b></p>	<p><b>Rate as Insufficient evidence with Exception</b></p>
<p>↓ Yes</p>		<p>↓ No</p>	<p>↓ No</p>			<p><b>Rate as Insufficient</b></p>
<p><b>No exception</b> → → → → →</p>	<p>→</p>	<p><b>No exception</b> → → → → →</p>	<p>→</p>	<p><b>No exception</b> →</p>	<p>→</p>	<p><b>Rate as Insufficient</b></p>

***Theme 7: Need for review of face validity***

Although empirical testing is preferable, NQF's criteria allow for use of face validity in lieu of empirical testing. Because face validity is the weakest form of validity, the updated guidance proposes that systematic assessment of face validity involve experts who were not involved in measure development. The assumption was that a developer would not submit a measure for potential endorsement if its experts did not think it had face validity. Commenters requested more clarity about defining "involvement in measure development" and two commenters specifically disagreed with requiring experts beyond those involved in measure development.

**ACTION ITEMS:**

- Discuss the role of face validity. Consider the following options:
  - No change to current guidance – developers choose who to use for systematic assessment of face validity.
  - Assume the measure has face validity with the developer and its experts at the time of initial endorsement. If the evidence subcriterion is met and the measure as specified is consistent with the evidence, does a systematic assessment of face validity offer any additional information, especially at the time of initial endorsement?
  - Should empirical validity testing be required? (for initial endorsement? For endorsement maintenance?)
    - For initial endorsement, this would represent a higher bar and would need to be phased in.
    - For endorsement maintenance, this is consistent with the approved recommendation from Measure Testing Task Force that testing be expanded on endorsement maintenance if a measure did not already achieve a high rating, so empirical testing would be required at that time. However, currently such expanded testing is often not conducted and NQF has not implemented this requirement to date.)

***Theme 8: Agreement that NQF should not specify minimum thresholds for sample sizes or reliability statistics***

Most commenters agreed that it is difficult or impossible to identify minimum thresholds that are applicable to all testing situations. Three commenters suggested considering power analysis to determine the appropriate sample size for the statistical test used. This will require further exploration.

**ACTION ITEMS:**

- Decide whether NQF should explore requiring a power analysis to justify sample size used for testing, and if so, discuss the timing and process for this exploration.

**ADDITIONAL DISCUSSION ON COMMENTS/RESPONSES**

Please identify any comments and draft responses that require discussion and resolution by the CSAC and/or HITAC. In the interest of using the call to discuss the substantive issues, we ask that you send us edits to the draft responses via email.