# Measure Evaluation Criteria Refresher
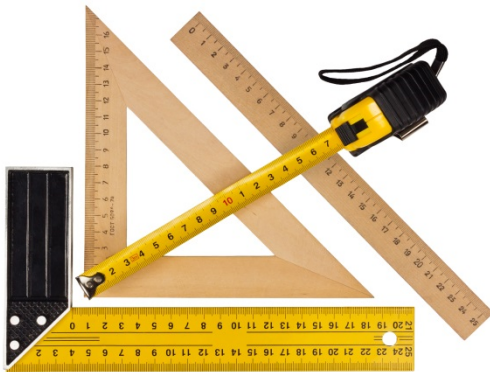
CSAC Informational Update

*April 24, 2019*

# Overview of Presentation

- Provide an overview of NQF's measure evaluation criteria
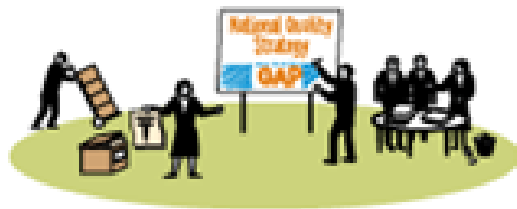
# Why Do We Measure?



The primary goal of healthcare performance measurement is to **improve the quality (and access and cost) of healthcare** received by patients (and ultimately, to **improve health**)
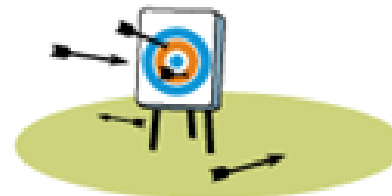
**Measurement is a quality improvement tool, not an end in and of itself**

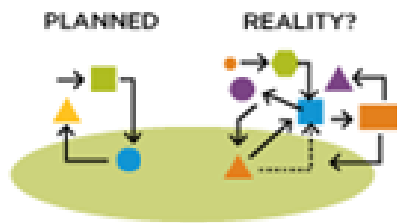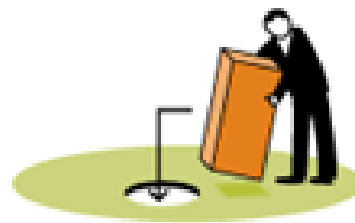# What Makes a Great Measure?
# NQF's Measure Evaluation Criteria



**1** IMPORTANCE TO MEASURE AND REPORT

**2** SCIENTIFIC ACCEPTABILITY OF MEASURE PROPERTIES

**3** FEASIBILITY

**4** USABILITY AND USE

**5** ASSESS RELATED AND COMPETING MEASURES

# NQF Measure Evaluation Criteria for Endorsement

**NQF endorses measures that are suitable for both quality improvement efforts as well as for accountability applications (public reporting, payment programs, accreditation, etc.)**

- Standardized evaluation criteria
- Criteria have evolved over time in response to stakeholder feedback
- The quality measurement enterprise is constantly growing and evolving—greater experience, lessons learned, expanding demands for measures—the criteria evolve to reflect the ongoing needs of stakeholders

# Major Endorsement Criteria

- **Importance to measure and report**: Goal is to measure those aspects with greatest potential of driving improvements; if not important, the other criteria are less meaningful (**must-pass**)

- **Reliability and Validity-scientific acceptability of measure properties**: Goal is to make valid conclusions about quality; if not reliable and valid, there is risk of improper interpretation (**must-pass**)

- **Feasibility**: Goal is to, ideally, cause as little burden as possible; if not feasible, consider alternative approaches

- **Usability and Use** (Use is **must-pass** for maintenance measures): Goal is to use for decisions related to accountability and improvement, to achieve high-quality, efficient healthcare for individuals or populations

- **Comparison to related or competing measures**

# Criterion #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 where there is variation in, or overall less-than-optimal, performance

*1a. Evidence:  the measure focus is evidence-based*

*1b.  Opportunity for Improvement:  demonstration of quality problems and opportunity for improvement (i.e., data demonstrating considerable variation, or overall less-than-optimal performance, in the quality of care across providers; and/or disparities in care across population groups*

*1c. Quality construct and rationale (composite measures only)*

# Subcriteron 1a:  Evidence

- ## Outcome measures
  - *Empirical data demonstrate a relationship between the outcome and at least one healthcare structure, process, intervention, or service.  If not available, wide variation in performance can be used as evidence, assuming the data are from a robust number of providers and results are not subject to systematic bias*

- ## Structure, process, intermediate outcome measures
  - *The quantity, quality, and consistency of the body of evidence underlying the measure should demonstrate that the measure focuses on those aspects of care known to influence desired patient outcomes*
    - » Empirical studies (expert opinion is not evidence)
    - » Systematic review and grading of evidence preferred

- ## For measures derived from "patient" report
  - *Evidence should demonstrate that the target population values the measured outcome, process, or structure and finds it meaningful*
  - *Current requirements for structure and process measures also apply to patient-reported structure/process measures*
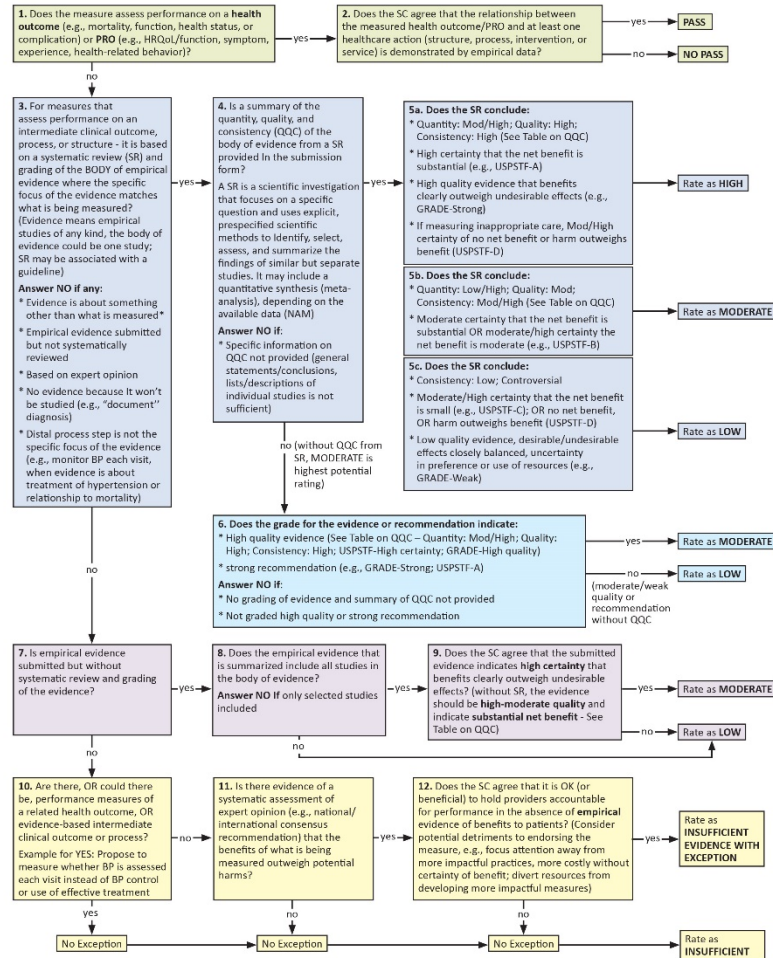
# Key Points for Evaluating Evidence

- The evaluation of the evidence subcriterion depends on the type of measure under consideration.

- Evidence should be presented about the relevant body of evidence—not selected individual studies.

- Ideally, measure developers will summarize a systematic review of the evidence that has been assembled, reviewed, and graded by others.

- Expert opinion is not considered empirical evidence, but evidence is not limited to randomized controlled trials.

- Measures with inconsistent or conflicting evidence should not pass the evidence subcriterion.

# Key Points for Evaluating Evidence

- When evaluating the quality of the evidence, consider the following:
  - *The study design itself (e.g., RCT, non-RCT) or flaws in the design or conduct of the study (e.g., lack of blinding; large losses to follow-up)*
  - *The directness/indirectness of the evidence to the measure as specified (e.g., regarding the population, intervention, comparators, and/or outcomes)*
  - *Imprecision in study results (i.e., wide confidence intervals due to few patients or events)*
- Under limited circumstances, an exception to the evidence subcriterion may be invoked and evaluated according to the evidence algorithm

# Algorithm for Rating Evidence



Algorithm 1. Guidance for Evaluating the Clinical Evidence

# Subcriteron 1b:  Opportunity for Improvement

- Underlying question:  is there a quality problem?
- We expect current data for the measure as specified; (relevant data from the literature also may be used, especially for initial endorsement)
- When evaluating whether there is opportunity for improvement, consider:
  - *The distribution of performance scores*
  - *The number and representativeness of the entities included in the measure performance data*
  - *The size of the population at risk, effectiveness of an intervention, likely occurrence of an outcome, and consequences of the quality problem*
  - *Data on disparities*
- Topped out measures: potential for Reserve Status

# Criterion #2: Reliability and Validity—Scientific Acceptability of Measure Properties

**Extent to which the measure, <u>as specified</u>, produces consistent (reliable) and credible (valid) results about the quality of healthcare delivery**

## 2a. Reliability (must-pass)

*2a1. Precise specifications*

*2a2. Reliability testing—data elements or measure score*

## 2b. Validity (must-pass)

*2b1. Validity testing—data elements or measure score*

*2b2. Justification of exclusions—relates to evidence*

*2b3. Risk adjustment—typically for outcome/cost/resource use*

*2b4. Identification of differences in performance*

*2b5. Comparability of data sources/methods*

*2b6. Missing data*

## 2c. Analysis support composite construction approach (must-pass)

# Questions These Subcriteria Address

**Reliability**

- Are the specifications clear so that everyone will calculate the measure in the same way?

- Is the variation between providers primarily due to real differences in performance? Or is it because there is a lot of "noise" in the measurement?

**Validity**

- Is the measure actually measuring what it is intended to measure (e.g., quality of care)?

- Do the results of the measurement allow for correct conclusions about quality of care?

# Key Points for Evaluating Scientific Acceptability

- Scientifically acceptable measures must be both reliable and valid

- Empirical demonstration of reliability and validity is expected, although for new measures, demonstration of face validity of the measure score as an indicator of quality also is allowed

- NQF is not prescriptive about how empirical measure testing is done; similarly, NQF does not set minimum thresholds for reliability or validity testing results

- Reliability and validity must be demonstrated for the measure as specified (including data source and level of analysis)

# Key Points for Evaluating Scientific Acceptability

- Depending on the measure type, NQF may allow testing at either the data element level (using patient-level data) or at the performance measure score level (using data that have been aggregated across providers).

- When evaluating measure testing results, the method of testing, the data used for testing (often from a sample), and the results of the testing must be considered

- All three subcriteria under Scientific Acceptability are "must-pass"; therefore, each must be met in order to be recommended for endorsement

# Some Additional Considerations

- We have rating algorithms for reliability and validity
- Most criteria/subcriteria involve a matter of degree rather than all-or-nothing determination—this requires both evidence and expert judgment
- Reliability
    - *The foundation for reliability is good specifications: definitions, codes, and instructions on how to calculate the measure*
    - *Evaluating testing for reliability can be tricky—several methods available, but no thresholds for results*
    - *Even so, assessing reliability is likely less subjective than assessing validity*
- Validity
    - *Assessing threats to validity even more important than testing*

# Some Additional Considerations: Accounting for Social Risk Factors

The Standing Committee will be asked to consider the following questions:

- Is there a conceptual relationship between the SDS factor and the measure focus?

- What are the patient-level sociodemographic variables that were available and analyzed during measure development?

- Does empirical analysis (as provided by the measure developer) show that the SDS factor has a significant and unique effect on the outcome in question?

- Does the reliability and validity testing match the final measure specifications?

# Criterion #3: Feasibility

Extent to which the required data are readily available, retrievable without undue burden, and can be implemented for performance measurement

*3a: Clinical data generated during care process*

*3b: Electronic sources*

*3c: Data collection strategy can be implemented*

- This is not a must-pass criterion
  - *HOWEVER, feasibility is critical for eCQMs*

# Criterion #4: Usability and Use

Extent to which potential audiences are using or could use performance results for both accountability and performance improvement to achieve the goal of high-quality, efficient healthcare for individuals or populations.

## Use (4a) Must-pass for maintenance measures

*4a1: Accountability and Transparency*: Performance results are used in at least one accountability application within three years after initial endorsement and are publicly reported within six years after initial endorsement

*4a2: Feedback by those being measured or others*: Those being measured have been given results and assistance in interpreting results; those being measured and others have been given opportunity for feedback; the feedback has been considered by developers

## Usability (4b)

*4b1: Improvement*: Progress toward achieving the goal of high-quality, efficient healthcare for individuals or populations is demonstrated.

*4b2: Benefits outweigh the harms*: The benefits of the performance measure in facilitating progress toward achieving high-quality, efficient healthcare for individuals or populations outweigh evidence of unintended negative consequences to individuals or populations (if such evidence exists)

# Criterion #5: Related or Competing Measures

If a measure meets the four criteria <u>and</u> there are endorsed/new **related** measures (same measure focus or same target population) or **competing** measures (both the same measure focus <u>and</u> same target population), the measures are compared to address harmonization and/or selection of the best measure.

- 5a. The measure specifications are harmonized with related measures **OR** the differences in specifications are justified.

- 5b. The measure is superior to competing measures (e.g., is a more valid or efficient way to measure) **OR** multiple measures are justified.