

MEASURE WORKSHEET

This document summarizes the evaluation of the measure as it progresses through NQF's Consensus Development Process (CDP). The information submitted by measure developers/stewards is included after the Brief Measure Information, Preliminary Analysis, and Pre-meeting Public and Member Comments sections.

To navigate the links in the worksheet: Click to go to the link. ALT + LEFT ARROW to return

Purple text represents the responses from measure developers.

Red text denotes developer information that has changed since the last measure evaluation review.

Brief Measure Information

NQF #: 3294

Measure Title: STS Lobectomy for Lung Cancer Composite Score

Measure Steward: The Society of Thoracic Surgeons

Brief Description of Measure: The STS Lobectomy Composite Score comprises two domains:

1. Operative Mortality (death during the same hospitalization as surgery or within 30 days of the procedure)
2. Presence of at least one of these major complications: pneumonia, acute respiratory distress syndrome, bronchopleural fistula, pulmonary embolus, initial ventilator support greater than 48 hours, reintubation/respiratory failure, tracheostomy, myocardial infarction, or unexpected return to the operating room.

The composite score is created by a weighted combination of the above two domains resulting in a single composite score. In addition to receiving a numeric score, participants are assigned to rating categories designated by the following:

- 1 star: lower-than expected performance
- 2 stars: as-expected-performance
- 3 start: higher-than-expected-performance

Developer Rationale: n/a

Numerator Statement: The STS Lobectomy Composite Score comprises two domains:

1. Operative Mortality (death during the same hospitalization as surgery or within 30 days of the procedure)
2. Presence of at least one of these major complications: pneumonia, acute respiratory distress syndrome, bronchopleural fistula, pulmonary embolus, initial ventilator support greater than 48 hours, reintubation/respiratory failure, tracheostomy, myocardial infarction, or unexpected return to the operating room.

The composite score is created by a weighted combination of the above two domains resulting in a single composite score. Operative mortality and major complications were weighted inversely by their respective standard deviations across participants. This procedure is equivalent to first rescaling mortality and complications by their respective standard deviations and then assigning equal weighting to the rescaled mortality rate and rescaled complication rate. This is the same methodology used for other STS composite measures.

In addition to receiving a numeric score, participants are assigned to rating categories designated by the following:

- 1 star: lower-than expected performance
- 2 stars: as-expected-performance

3 start: higher-than-expected-performance

Patient Population: The STS GTSD was queried for all patients treated with lobectomy for lung cancer between January 1, 2014, and December 31, 2016. We excluded patients with non-elective status, occult or stage 0 tumors, American Society of Anesthesiologists class VI, and with missing data for age, sex, or discharge mortality status.

Time Window: 01/01/2014 - 12/31/2016

Model variables: Variables in the model: age, sex, year of operation, body mass index, hypertension, steroid therapy, congestive heart failure, coronary artery disease, peripheral vascular disease, reoperation, preoperative chemotherapy within 6 months, cerebrovascular disease, diabetes mellitus, renal failure, dialysis, past smoker, current smoker, forced expiratory volume in 1 second percent of predicted, Zubrod score (linear plus quadratic), American Society of Anesthesiologists class (linear plus quadratic), and pathologic stage.

Denominator Statement: Number of patients greater than or equal to 18 years of age undergoing elective lobectomy for lung cancer

Denominator Exclusions: Patients were excluded with non-elective status, occult or stage 0 tumors, American Society of Anesthesiologists class VI, and with missing data for age, sex, or discharge mortality status.

Measure Type: Composite

Data Source: Other, Registry Data

Level of Analysis: Facility

IF Endorsement Maintenance – Original Endorsement Date: Most Recent Endorsement Date:

Staff Preliminary Analysis: New Measure

Criteria 1: Importance to Measure and Report

1a. [Evidence](#)

1a. Evidence. The evidence requirements for a health outcome measure include providing empirical data that demonstrate a relationship between the outcome and at least one healthcare structure, process, intervention, or service; if these data not available, data demonstrating wide variation in performance, assuming the data are from a robust number of providers and results are not subject to systematic bias. For measures derived from patient report, evidence also should demonstrate that the target population values the measured outcome, process, or structure and finds it meaningful.

Evidence Summary

- This new measure assesses the operative mortality and the presence of at least one of 9 major [complications](#) of lobectomy, the most frequently performed lung resection procedure. The developer reports that data in the STS General Thoracic Surgery database (GTSD) show a reduction in perioperative morbidity and equivalent long term survival when minimally invasive approaches for lobectomy are used.
- The developer provided the performance data below, 0.95 to 0.98, for approximately 200-300 participants and 24,000+ operations from 2013 to 2016.
- *Empirical data* demonstrating a relationship between the outcome to at least one healthcare process is now required. NQF guidance states that a 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.

Question for the Committee:

- *Is there at least one thing that the provider can do to achieve a change in the measure results?*

Guidance from the Evidence Algorithm: Measure assesses performance on a health outcome (Box 1) → The relationship between the outcome and the intervention demonstrated by performance data (Box 2) → Pass

Preliminary rating for evidence: ☒ Pass ☐ No Pass

1b. [Gap in Care/Opportunity for Improvement](#) and 1b. [Disparities](#)

Maintenance measures – increased emphasis on gap and variation

1b. Performance Gap. The performance gap requirements include demonstrating quality problems and opportunity for improvement.

- Data were collected in two overlapping 3 year time periods: January 1, 2014 – December 31, 2016 and January 1, 2013 – December 31, 2015.

	January 1, 2013 – December 31, 2015		January 1, 2014 – December 31, 2016	
No. participants	242	185	233	286
No. of operations	23,574	22,572	24,912	24,318
Mean	0.972	0.972	0.973	0.974
Standard Deviation	0.007	0.008	0.006	0.007
IQR	0.008	0.009	0.007	0.009
Minimum	0.945	0.945	0.953	0.953
Maximum	0.988	0.988	0.987	0.987

Disparities

- The developer provides descriptive data of the sampled population, but disparities data for these groups are not provided.

Questions for the Committee:

- Does the Committee think there is enough variation among providers to justify a national performance measure?

Preliminary rating for opportunity for improvement: ☐ High ☒ Moderate ☐ Low ☐ Insufficient

1c. Composite – [Quality Construct and Rationale](#)

Maintenance measures – same emphasis on quality construct and rationale as for new measures.

1c. Composite Quality Construct and Rationale. The quality construct and rationale should be explicitly articulated and logical; a description of how the aggregation and weighting of the components is consistent with the quality construct and rationale also should be explicitly articulated and logical.

Quality construct

- This measure is based on a combination of an operative mortality outcome and the risk adjusted occurrence of any of nine major complications. Operative mortality is described as death during the same hospitalization as surgery or within 30 days of the procedure. Complications include:
 - Pneumonia
 - Acute respiratory distress syndrome
 - Branchopleural fistula
 - Pulmonary embolus
 - Initial ventilator support greater than 48 hours
 - Reintubation/respiratory failure
 - Tracheostomy
 - Myocardial infarction

- Unexpected return to the operating room
- Participants are scored for each domain (mortality and complication), and an overall composite score which is created by a weighted combination of the two domains. Participants are also assigned a rating designated by one to three stars:
 - 1 star: lower-than expected performance
 - 2 stars: as-expected performance
 - 3 stars: higher than expected performance
- The developer reports that since mortality rates for thoracic surgery have declined, it can be difficult to differentiate performance based on mortality alone since it fails to take into account that not all operative survivors received equal quality care. Therefore, a composite score from a weighted combination of mortality and operative complications provides a more comprehensive measure of overall surgical quality.
- Operative mortality is weighted approximately four times that of a major complication in the composite. The developer reports this weighting is consistent with STS adult cardiac measures.

Questions for the Committee:

- *Are the quality construct and a rationale for the composite explicitly stated and logical?*
- *Is the method for aggregation and weighting of the components explicitly stated and logical?*

Preliminary rating for composite quality construct and rationale: ☒ High ☐ Moderate ☐ Low ☐ Insufficient

Committee Pre-evaluation Comments: Criteria 1: Importance to Measure and Report (including 1a, 1b, 1c)

1a. Evidence

**STS General Thoracic Surgery database (GTSD) 200-300 patients 24,000 patients (?) PASS

**good evidence

**Adequate evidence

**An important measure for public accountability, as already illustrated by improvement in outcomes over the course of the registry. Weighting death 4x morbidities is somewhat arbitrary, but reasonable and consistent with other such measures.

1b. Performance Gap

**Improvement, Two year data, MODERATE per NQF reviewer "

**PG present

**Increasing morbidity associated with lobectomy clearly justifies this composite measure

**Minimal gap (91% average performance), so limited opportunity for quality improvement. But as noted above, important for public accountability.

1c. Composite Quality Construct

**Operative mortality is weighted approximately four times that of a major complication in the composite, consistent with the STS adult cardiac surgery quality measures. The STS General Thoracic Surgery Database working group believes this is an improvement from its previous lung cancer resection model in which mortality and major morbidity were weighted equally. Logical

**High quality composite construct

**The construct makes good sense and adds value to the individual components

**See above. Well constructed statistically, but balance between different outcomes will always be arbitrary.

Criteria 2: Scientific Acceptability of Measure Properties

2a. Reliability: [Specifications](#) and [Testing](#)

2b. Validity: [Testing](#); [Exclusions](#); [Risk-Adjustment](#); [Meaningful Differences](#); [Comparability](#) [Missing Data](#)

2c. For composite measures: [empirical analysis support composite approach](#)

Reliability

2a1. Specifications requires the measure, as specified, to produce consistent (reliable) and credible (valid) results about the quality of care when implemented. For maintenance measures – no change in emphasis – specifications should be evaluated the same as with new measures.

2a2. Reliability testing demonstrates if the measure data elements are repeatable, producing the same results a high proportion of the time when assessed in the same population in the same time period and/or that the measure score is precise enough to distinguish differences in performance across providers. For maintenance measures – less emphasis if no new testing data provided.

Validity

2b2. Validity testing should demonstrate the measure data elements are correct and/or the measure score correctly reflects the quality of care provided, adequately identifying differences in quality. For maintenance measures – less emphasis if no new testing data provided.

2b2-2b6. Potential threats to validity should be assessed/addressed.

Composite measures only:

2d. Empirical analysis to support composite construction. Empirical analysis should demonstrate that the component measures add value to the composite and that the aggregation and weighting rules are consistent with the quality construct.

Complex measure evaluated by Scientific Methods Panel? ☒ Yes ☐ No

Evaluators: Jennifer Perloff, Ron Walters, Joe Kunisch, David Cella, Karen Maddox

Evaluation of Reliability and Validity (and composite construction, if applicable):

[Evaluation A](#)

[Evaluation B](#)

[Evaluation C](#)

[Evaluation D](#)

[Evaluation E](#)

Questions for the Committee regarding reliability:

- Do you have any concerns that the measure can be consistently implemented (i.e., are measure specifications adequate)?
- The Scientific Methods Panel is satisfied with the reliability testing for the measure. Does the Committee think there is a need to discuss and/or vote on reliability?

Questions for the Committee regarding validity:

- Do you have any concerns regarding the validity of the measure (e.g., exclusions, risk-adjustment approach, etc.)?
- The Scientific Methods Panel is satisfied with the validity analyses for the measure. Does the Committee think there is a need to discuss and/or vote on validity?

Questions for the Committee regarding composite construction:

- Do you have any concerns regarding the composite construction approach (e.g., do the component measures fit the quality construct and add value to the overall composite? Are the aggregation and weighting rules consistent with the quality construct and rationale while achieving the related objective of simplicity to the extent possible?)?
- The Scientific Methods Panel is satisfied with the composite construction. Does the Committee think there is a need to discuss the composite construction approach?

Preliminary rating for reliability: ☐ High ☒ Moderate ☐ Low ☐ Insufficient

Preliminary rating for validity: ☐ High ☒ Moderate ☐ Low ☐ Insufficient

Note: While score-level validity testing is desired, data element testing is accepted because this is a new measure. For future maintenance evaluations, score-level testing will be required.

Preliminary rating for composite construction: ☐ High ☒ Moderate ☐ Low ☐ Insufficient

Evaluation A: Scientific Acceptability

Extent to which the measure, as specified, produces consistent (reliable) and credible (valid) results about the quality of care when implemented. **Measures must be judged to meet the subcriteria for both reliability and validity to pass this criterion.**

Instructions:

1. Please complete this form for each measure you are evaluating.
2. Please pay close attention to the skip logic directions.
3. If you are unable to check a box, please highlight or shade the box for your response.
4. You must answer the “overall rating” item for both Reliability and Validity. Also, be sure to answer the composite measure question at the end of the form if your measure is a composite.
5. We have provided TIPS to help you answer the questions.
6. We’ve designed this form to try to minimize the amount of writing that you have to do. That said, **it is critical that you explain your thinking/rationale if you check boxes where we ask for an explanation** (because this is a Word document, you can just add your explanation below the checkbox). Feel free to add additional explanation, even if an explanation is not requested (but please type this underneath the appropriate checkbox).
7. This form is based on Algorithms 2 and 3 in the Measure Evaluation Criteria and Guidance document (see pages 18-24). These algorithms provide guidance to help you rate the Reliability and Validity subcriteria. **We ask that you refer to this document when you are evaluating your measures.**
8. Please contact Methods Panel staff if you have questions (methodspanel@qualityforum.org).

Measure Number: NQF#3294

Measure Title: STS Lobectomy for Lung Cancer Composite Score

RELIABILITY

9. Are submitted specifications precise, unambiguous, and complete so that they can be consistently implemented?

NOTE: NQF staff will conduct a separate, more technical, check of eMeasure (eCQM) specifications, value sets, logic, and feasibility, so no need to consider these in your evaluation.

TIPS: Consider the following: Are all the data elements clearly defined? Are all appropriate codes included? Is the logic or calculation algorithm clear? Is it likely this measure can be consistently implemented?

☒ Yes (go to Question #2)

☐ No (please explain below, and go to Question #2) NOTE that even though **non-precise**

specifications should result in an overall LOW rating for reliability, we still want you to look at the testing results.

The measure clearly defines the intended outcomes measures, mortality and complications of care. These have been utilized by this database for many years and are consistent with prior definitions in previous

implementations. The measure is intended at the facility level AND the individual group/practice level, however, data is provided for 233 participants (facilities) and 24,912 patient records. The data elements are clearly defined and annually audited for data completeness and accuracy.

Data presented, however, seems to be at the facility level (233/24,912) and I could not find comparable testing at the group/practice level from the submission. The term participant in both the submission and the publication referenced appears to apply to the facility level only.

10. Was empirical reliability testing (at the data element or measure score level) conducted using statistical tests with the measure as specified?

TIPS: Check the 2nd "NO" box below if: only descriptive statistics provided; only describes process for data management/cleaning/computer programming; testing does not match measure specifications (i.e. data, eMeasure, level of analysis, patients)

☒ Yes (go to Question #4)

☐ No, there is reliability testing information, but *not* using statistical tests and/or not for the measure as specified OR there is no reliability testing (please explain below then go to Question #3)

Reliability at the data element level (in 2a2.4) of 44.6% for all and higher for increasing number of cases (up to 68.0%) and the score level (one star to three stars) with the weighted composite is tested and indicated in 2d1.2.

11. Was **empirical VALIDITY testing** of patient-level data conducted?

☐ Yes (use your rating from data element validity testing – Question #16- under Validity Section)

☐ No (please explain below and rate Question #11: OVERALL RELIABILITY as INSUFFICIENT and proceed to the [VALIDITY SECTION](#))

Data element validity testing is stated in 1.7 as being via an annual audit of data completeness and accuracy for randomly selected surgical records at randomly selected participant sites, described in 2b1.2. A data element quality report is generated and provided to the participant for action, if required. Agreement was 97.78% in 2016.

12. Was reliability testing conducted with computed performance measure scores for each measured entity?

TIPS: Answer no if: only one overall score for all patients in sample used for testing patient-level data

☒ Yes (go to Question #5)

☐ No (go to Question #8)

Yes at the facility level. No at the group/practice level.

13. Was the method described and appropriate for assessing the proportion of variability due to real differences among measured entities? *NOTE: If multiple methods used, at least one must be appropriate.*

TIPS: Examples of appropriate methods include signal-to-noise analysis (e.g. Adams/RAND tutorial); random split-half correlation; other accepted method with description of how it assesses reliability of the performance score.

☒ Yes (go to Question #6)

☐ No (please explain below then go to Question #8)

Section 2a2.2 provides the methodology and the results.

14. **RATING (score level)** - What is the level of certainty or confidence that the performance measure scores are reliable?

TIPS: Consider the following: Is the test sample adequate to generalize for widespread implementation?

Do the results demonstrate sufficient reliability so that differences in performance can be identified?

☒ High (go to Question #8)

☐ Moderate (go to Question #8)

☐ Low (please explain below then go to Question #7)

At the facility level, score reliability does separate out those with high mortality, 1.2% and complications, 16.2% (one star) from those with low mortality, 0.4% and complications, 3.2%, (three star). An expert panel provided an assessment of validity. The methodology is described in 2b4.1.

15. Was other reliability testing reported?

☐ Yes (go to Question #8)

☐ No (rate Question #11: OVERALL RELIABILITY as LOW and proceed to the [VALIDITY SECTION](#))

16. Was reliability testing conducted with patient-level data elements that are used to construct the performance measure?

TIPS: Prior reliability studies of the same data elements may be submitted; if comparing abstraction to “authoritative source/gold standard” see Validity Section Question #15)

☐ Yes (go to Question #9)

☒ No (if there is score-level testing, rate Question #11: OVERALL RELIABILITY based on score-level rating from Question #6; otherwise, rate Question #11: OVERALL RELIABILITY as INSUFFICIENT. Then proceed to the [VALIDITY SECTION](#))

The inter-rater or intra-rater reliability testing is not specifically given in this measure submission but I suspect is known from prior experience with this dataset. The referenced data in the submission is to the audits for data accuracy. Thus, though I suspect the answer to this question is YES, I cannot state this from the data given.

17. Was the method described and appropriate for assessing the reliability of ALL critical data elements?

TIPS: For example: inter-abtractor agreement (ICC, Kappa); other accepted method with description of how it assesses reliability of the data elements

Answer no if: only assessed percent agreement; did not assess separately for all data elements (at least numerator, denominator, exclusions)

☐ Yes (go to Question #10)

☐ No (if no, please explain below and rate Question #10 as INSUFFICIENT)

18. **RATING (data element)** – Based on the reliability statistic and scope of testing (number and representativeness of patients and entities), what is the level of certainty or confidence that the data used in the measure are reliable?

TIPS: Consider the following: Is the test sample adequate to generalize for widespread implementation? Can data elements be collected consistently?

☐ Moderate (if score-level testing was NOT conducted, rate Question #11: OVERALL RELIABILITY as MODERATE)

☐ Low (if score-level testing was NOT conducted, rate Question #11: OVERALL RELIABILITY as LOW)

☐ Insufficient (go to Question #11)

11. OVERALL RELIABILITY RATING

OVERALL RATING OF RELIABILITY taking into account precision of specifications and all testing results:

☐ High (NOTE: Can be HIGH only if score-level testing has been conducted)

☐ Moderate (NOTE: Moderate is the highest eligible rating if score-level testing has not been conducted)

☐ Low (please explain below) [NOTE: Should rate LOW if you believe specifications are NOT precise, unambiguous, and complete]

☒ Insufficient (please explain below) [NOTE: For most measure types, testing at both the score level and the

data element level is not required]

Information is not provided about the inter-rater reliability at the data element level and is substituted by the results of random audits of data elements.

If, by the term “participant”, both facility and group/practice level is the intention and has been performed (see Question 1), and if there is prior evidence of inter-rater and intra-rater reliability testing historically, not based on random audits, then I would be willing to consider changing the overall rating to high. I could not infer this from the submission and these points should be clarified and discussed.

VALIDITY

ASSESSMENT OF THREATS TO VALIDITY

1. Were all potential threats to validity that are relevant to the measure empirically assessed?

TIPS: Threats to validity include: exclusions; need for risk adjustment; Able to identify statistically significant and meaningful differences; multiple sets of specifications; missing data/nonresponse.

☒ Yes (go to Question #2)

☐ No (please explain below and go to Question #2) [NOTE that even if **non-assessment of applicable threats should result in an overall INSUFFICIENT rating for validity**, we still want you to look at the testing results]

There was acknowledgement of the potential impact of missing data elements. A conscious decision was made to either impute the data value from other elements present, to the median, or to the value indicating absence of the risk factor for some of the data elements, and to exclude others. The range of missing value was between 1% and 3.5%. The conclusion was that this did not lead to bias in the measure.

2. Analysis of potential threats to validity: Any concerns with measure exclusions?

TIPS: Consider the following: Are the exclusions consistent with the evidence? Are any patients or patient groups inappropriately excluded from the measure? Are the exclusions/exceptions of sufficient frequency and variation across providers to be needed (and outweigh the data collection burden)? If patient preference (e.g., informed decisionmaking) is a basis for exclusion, does it impact performance and if yes, is the measure specified so that the information about patient preference and the effect on the measure is transparent?

☒ Yes (please explain below then go to Question #3)

☐ No (go to Question #3)

☐ Not applicable (i.e., there are no exclusions specified for the measure; go to Question #3)

The conscious decision regarding social risk factors is discussed in the submission and below in Question 3.

3. Analysis of potential threats to validity: Risk-adjustment (applies to all outcome, cost, and resource use measures; may also apply to other types of measure)

☐ Not applicable (e.g., structure or process measure that is not risk-adjusted; go to Question #4)

a. Is a conceptual rationale for social risk factors included? ☒ Yes ☐ No

b. Are social risk factors included in risk model? ☐ Yes ☒ No

c. Any concerns regarding the risk-adjustment approach?

*TIPS: Consider the following: If a justification for **not risk adjusting** is provided, is there any evidence that contradicts the developer's rationale and analysis? If the developer asserts there is **no conceptual basis** for adjusting this measure for social risk factors, do you agree with the rationale? **If risk adjusted:** Are the candidate and final variables included in the risk adjustment model adequately described for the measure to be implemented? Are the candidate and final variables included in the risk adjustment model adequately described for the measure to be implemented? Are all of the risk adjustment variables present at the start of care (if not, do you agree with the rationale)? If social risk factors are not included in the risk-adjustment approach, do you agree with the developer's decision? Is an appropriate risk-adjustment strategy included in the measure (e.g.,*

adequate model discrimination and calibration)? Are all statistical model specifications included, including a “clinical model only” if social risk factors are included in the final model?

☒ Yes (please explain below then go to Question #4)

☐ No (go to Question #4)

Social risk factors are not collected in the database and therefore not included in the risk adjustment. It is possible that the data elements collected override other social risk factors, or account for them, but it would be nice to see some statement to that effect. Payer status as a proxy is a part of the database but analysis has not been performed as to its additive value to the model.

4. Analysis of potential threats to validity: Any concerns regarding ability to identify meaningful differences in performance or overall poor performance?

☐ Yes (please explain below then go to Question #5)

☒ No (go to Question #5)

Despite the above considerations, the large sample size and the historical usage of this database does lead to confidence in the assumptions. And, despite the above, there were statistically meaningful differences in performance demonstrated in the measure score.

5. Analysis of potential threats to validity: Any concerns regarding comparability of results if multiple data sources or methods are specified?

☐ Yes (please explain below then go to Question #6)

☐ No (go to Question #6)

☒ Not applicable (go to Question #6)

Sole data source is the abstracted STS.

6. Analysis of potential threats to validity: Any concerns regarding missing data?

☐ Yes (please explain below then go to Question #7)

☒ No (go to Question #7)

See Question 1 above and note Section 2b6 describes the analysis and subsequent attribution of missing data elements and efforts to minimize their impact.

ASSESSMENT OF MEASURE TESTING

7. Was empirical validity testing conducted using the measure as specified and appropriate statistical test?

Answer no if: face validity; only refer to clinical evidence; only descriptive statistics; only describe process for data management/cleaning/computer programming; testing does not match measure specifications (i.e. data, eMeasure, level, setting, patients).

☒ Yes (go to Question #10) [NOTE: If appropriate empirical testing has been conducted, then evaluation of face validity is not necessary. Go to Question #8 **only if** there is insufficient information provided to evaluate data element and score-level testing.]

☐ No (please explain below then go to Question #8)

Section 2d1.x describes the methodology used to assess the weighting and the effect on the metric of star ratings. It is empirical in that it is applied to hospitals (participants?) with more than 30 lobectomies and results in valid separation between one star and three start ratings for both operative mortality and complication rates. Morbidity is noted to explain more of the variation in the score. Sections 2d2.1 and 2d2.2 describe the derivation of the weight distribution. This fit expert panel expectations.

8. Was face validity systematically assessed by recognized experts to determine agreement on whether the computed performance measure score from the measure as specified can be used to distinguish good and poor quality?

TIPS: Answer no if: focused on data element accuracy/availability/feasibility/other topics; the degree of consensus and any areas of disagreement not provided/discussed.

☒ Yes (go to Question #9)

☐ No (please explain below and rate Question #17: OVERALL VALIDITY as INSUFFICIENT)

Although not required as the assessment of Question 7 was "YES", the measure score was assessed by a panel of experts and the methodology was felt to accurately portray the relative contribution of mortality and morbidities to the overall score. It did result in statistically significant differences between those with one-star and three-star ratings.

9. **RATING (face validity)** - Do the face validity testing results indicate substantial agreement that the performance measure score from the measure as specified can be used to distinguish quality AND potential threats to validity are not a problem, OR are adequately addressed so results are not biased?

☒ Yes (if a NEW measure, rate Question #17: OVERALL VALIDITY as MODERATE)

☐ Yes (if a MAINTENANCE measure, do you agree with the justification for not conducting empirical testing? If no, rate Question #17: OVERALL VALIDITY as

INSUFFICIENT; otherwise, rate Question #17: OVERALL VALIDITY as MODERATE)

☐ No (please explain below and rate Question #17: OVERALL VALIDITY AS LOW)

10. Was validity testing conducted with computed performance measure scores for each measured entity?

TIPS: Answer no if: one overall score for all patients in sample used for testing patient-level data.

☒ Yes (go to Question #11)

☐ No (please explain below and go to Question #13)

11. Was the method described and appropriate for assessing conceptually and theoretically sound hypothesized relationships?

TIPS: For example: correlation of the performance measure score on this measure and other performance measures; differences in performance scores between groups known to differ on quality; other accepted method with description of how it assesses validity of the performance score

☒ Yes (go to Question #12)

☐ No (please explain below, rate Question #12 as INSUFFICIENT and then go to Question #14)

Section 2d1.2 demonstrates that the score results do separate the groups in to high, medium and low performers (star ratings), and that the score does reflect both components of mortality and major complications. The score does demonstrate that the components included in the composite are consistent with the described quality construct and add value to the overall composite.

The question is, when the components of the composite score are THE two most important quality outcomes pertinent to the patient, in this case, mortality and complications, can those themselves be used as indicators of quality from validity testing perspective, which is shown by the table presented. I cannot think of more important quality indicators against which these two could be tested and, therefore, have to say yes to the methodological appropriateness.

12. **RATING (measure score)** - Based on the measure score results (significance, strength) and scope of testing (number of measured entities and representativeness) and analysis of potential threats, what is the level of certainty or confidence that the performance measure scores are a valid indicator of quality?

☐ High (go to Question #14)

☒ Moderate (go to Question #14)

☐ Low (please explain below then go to Question #13)

☐ Insufficient

It would have been nice to see some data about the effect of different weightings on the validity of the composite score. The predominant test was face validity and the score (and star ratings) derived from the model.

13. Was other validity testing reported?

☒ Yes (go to Question #14)

☐ No (please explain below and rate Question #17: OVERALL VALIDITY as LOW)

Face validity with a panel of experts was used to assess the validity of the model, who said that an 82.7/17.3 ratio was intuitively supported.

14. Was validity testing conducted with patient-level data elements?

TIPS: Prior validity studies of the same data elements may be submitted

☒ Yes (go to Question #15)

☐ No (please explain below and rate Question #17: OVERALL VALIDITY as INSUFFICIENT if no score-level testing was conducted, otherwise, rate Question #17: OVERALL VALIDITY based on score-level rating from Question #12)

Apparently, though not stated, probably due to resource requirements, the data field validity was capped at 500 maximum denominator. It would be nice to state that.

15. Was the method described and appropriate for assessing the accuracy of ALL critical data elements? *NOTE that data element validation from the literature is acceptable.*

TIPS: For example: Data validity/accuracy as compared to authoritative source- sensitivity, specificity, PPV, NPV; other accepted method with description of how it assesses validity of the data elements.

Answer No if: only assessed percent agreement; did not assess separately for all data elements (at least numerator, denominator, exclusions)

☐ Yes (go to Question #16)

☒ No (please explain below and rate Question #16 as INSUFFICIENT)

It is noted that due to the absence of access to all of the data results, a kappa statistic could not be provided. Generally, percent agreement is not sufficient while easily understood. Another option would have been sensitivity/specificity calculations.

16. **RATING (data element)** - Based on the data element testing results (significance, strength) and scope of testing (number and representativeness of patients and entities) and analysis of potential threats, what is the level of certainty or confidence that the data used in the measure are valid?

☐ Moderate (if score-level testing was NOT conducted, rate Question #17: OVERALL VALIDITY as MODERATE)

☐ Low (please explain below) (if score-level testing was NOT conducted, rate Question #17: OVERALL VALIDITY as LOW)

☒ Insufficient (go to Question #17)

17. OVERALL VALIDITY RATING

OVERALL RATING OF VALIDITY taking into account the results and scope of all testing and analysis of potential threats.

☐ High (NOTE: Can be HIGH only if score-level testing has been conducted)

☒ Moderate (NOTE: Moderate is the highest eligible rating if score-level testing has NOT been conducted)

☐ Low (please explain below) [NOTE: Should rate LOW if you believe that there are threats to validity and/or threats to validity were not assessed]

☐ Insufficient (if insufficient, please explain below) [NOTE: For most measure types, testing at both the score level and the data element level is not required] [NOTE: If rating is INSUFFICIENT for all empirical testing, then go back to Question #8 and evaluate any face validity that was conducted, then reconsider this overall rating.]

The statistic for data element validity is not the best available. The authors did mention their lack of ability to calculate a kappa statistic. There is a conscious lack of the use of social risk factors. See Question 11 above for the discussion regarding score level validity.

FOR COMPOSITE MEASURES ONLY: Empirical analyses to support composite construction

What is the level of certainty or confidence that the empirical analysis demonstrates that the component measures add value to the composite and that the aggregation and weighting rules are consistent with the quality construct?

TIPS: Consider the following: Do the component measures fit the quality construct? Are the objectives of parsimony and simplicity achieved while supporting the quality construct?

☐ High

☒ Moderate

☐ Low (please explain below)

☐ Insufficient (please explain below)

See Questions 11, 12 and 13 above. Clinical rationale is good. Precisely how the 83/17 ratio was derived and why it is the most applicable one is not clear from the data submitted.

Evaluation B: Scientific Acceptability

Extent to which the measure, as specified, produces consistent (reliable) and credible (valid) results about the quality of care when implemented. **Measures must be judged to meet the subcriteria for both reliability and validity to pass this criterion.**

Instructions:

19. Please complete this form for each measure you are evaluating.
20. Please pay close attention to the skip logic directions.
21. If you are unable to check a box, please highlight or shade the box for your response.
22. You must answer the "overall rating" item for both Reliability and Validity. Also, be sure to answer the composite measure question at the end of the form if your measure is a composite.
23. We have provided TIPS to help you answer the questions.
24. We've designed this form to try to minimize the amount of writing that you have to do. That said, **it is critical that you explain your thinking/rationale if you check boxes where we ask for an explanation** (because this is a Word document, you can just add your explanation below the checkbox). Feel free to add additional explanation, even if an explanation is not requested (but please type this underneath the appropriate checkbox).
25. This form is based on Algorithms 2 and 3 in the Measure Evaluation Criteria and Guidance document (see pages 18-24). These algorithms provide guidance to help you rate the Reliability and Validity subcriteria. **We ask that you refer to this document when you are evaluating your measures.**
26. Please contact Methods Panel staff if you have questions (methodspanel@qualityforum.org).

Measure Number: 3294

Measure Title: STS Lobectomy for Lung Cancer Composite Score

RELIABILITY

27. Are submitted specifications precise, unambiguous, and complete so that they can be consistently implemented?

NOTE: NQF staff will conduct a separate, more technical, check of eMeasure (eCQM) specifications, value sets, logic, and feasibility, so no need to consider these in your evaluation.

TIPS: Consider the following: Are all the data elements clearly defined? Are all appropriate codes included? Is the logic or calculation algorithm clear? Is it likely this measure can be consistently implemented?

☒ Yes (go to Question #2)

☐ No (please explain below, and go to Question #2) *NOTE that even though **non-precise specifications should result in an overall LOW rating for reliability**, we still want you to look at the testing results.*

28. Was empirical reliability testing (at the data element or measure score level) conducted using statistical tests with the measure as specified?

TIPS: Check the 2nd "NO" box below if: only descriptive statistics provided; only describes process for data management/cleaning/computer programming; testing does not match measure specifications (i.e. data, eMeasure, level of analysis, patients)

☒ Yes (go to Question #4)

☐ No, there is reliability testing information, but *not* using statistical tests and/or not for the measure as specified OR there is no reliability testing (please explain below then go to Question #3)

29. Was **empirical** **VALIDITY** testing of patient-level data conducted?

☐ Yes (use your rating from data element validity testing – Question #16- under Validity Section)

☐ No (please explain below and rate Question #11: OVERALL RELIABILITY as INSUFFICIENT and proceed to the [VALIDITY SECTION](#))

30. Was reliability testing conducted with computed performance measure scores for each measured entity?

TIPS: Answer no if: only one overall score for all patients in sample used for testing patient-level data

☒ Yes (go to Question #5)

☐ No (go to Question #8)

31. Was the method described and appropriate for assessing the proportion of variability due to real differences among measured entities? *NOTE: If multiple methods used, at least one must be appropriate.*

TIPS: Examples of appropriate methods include signal-to-noise analysis (e.g. Adams/RAND tutorial); random split-half correlation; other accepted method with description of how it assesses reliability of the performance score.

☒ Yes (go to Question #6)

☐ No (please explain below then go to Question #8)

32. **RATING (score level)** - What is the level of certainty or confidence that the performance measure scores are reliable?

TIPS: Consider the following: Is the test sample adequate to generalize for widespread implementation?

Do the results demonstrate sufficient reliability so that differences in performance can be identified?

☐ High (go to Question #8)

☒ Moderate (go to Question #8)

☐ Low (please explain below then go to Question #7)

33. Was other reliability testing reported?

☐ Yes (go to Question #8)

☐ No (rate Question #11: OVERALL RELIABILITY as LOW and proceed to the [VALIDITY SECTION](#))

34. Was reliability testing conducted with patient-level data elements that are used to construct the performance measure?

TIPS: Prior reliability studies of the same data elements may be submitted; if comparing abstraction to “authoritative source/gold standard” see Validity Section Question #15)

☐ Yes (go to Question #9)

☒ No (if there is score-level testing, rate Question #11: OVERALL RELIABILITY based on score-level rating from Question #6; otherwise, rate Question #11: OVERALL RELIABILITY as INSUFFICIENT. Then proceed to the [VALIDITY SECTION](#))

35. Was the method described and appropriate for assessing the reliability of ALL critical data elements?

TIPS: For example: inter-abtractor agreement (ICC, Kappa); other accepted method with description of how it assesses reliability of the data elements

Answer no if: only assessed percent agreement; did not assess separately for all data elements (at least numerator, denominator, exclusions)

☐ Yes (go to Question #10)

☐ No (if no, please explain below and rate Question #10 as INSUFFICIENT)

36. **RATING (data element)** – Based on the reliability statistic and scope of testing (number and representativeness of patients and entities), what is the level of certainty or confidence that the data used in the measure are reliable?

TIPS: Consider the following: Is the test sample adequate to generalize for widespread implementation? Can data elements be collected consistently?

☐ Moderate (if score-level testing was NOT conducted, rate Question #11: OVERALL RELIABILITY as MODERATE)

☐ Low (if score-level testing was NOT conducted, rate Question #11: OVERALL RELIABILITY as LOW)

☐ Insufficient (go to Question #11)

11. OVERALL RELIABILITY RATING

OVERALL RATING OF RELIABILITY taking into account precision of specifications and all testing results:

☐ High (NOTE: Can be HIGH only if score-level testing has been conducted)

☒ Moderate (NOTE: Moderate is the highest eligible rating if score-level testing has not been conducted)

☐ Low (please explain below) [NOTE: Should rate LOW if you believe specifications are NOT precise, unambiguous, and complete]

☐ Insufficient (please explain below) [NOTE: For most measure types, testing at both the score level and the data element level is not required]

VALIDITY

ASSESSMENT OF THREATS TO VALIDITY

17. Were all potential threats to validity that are relevant to the measure empirically assessed?

TIPS: Threats to validity include: exclusions; need for risk adjustment; Able to identify statistically significant and meaningful differences; multiple sets of specifications; missing data/nonresponse.

☒ Yes (go to Question #2)

☐ No (please explain below and go to Question #2) [NOTE that even if **non-assessment of applicable**

threats should result in an overall INSUFFICIENT rating for validity, we still want you to look at the testing results]

18. Analysis of potential threats to validity: Any concerns with measure exclusions?

TIPS: Consider the following: Are the exclusions consistent with the evidence? Are any patients or patient groups inappropriately excluded from the measure? Are the exclusions/exceptions of sufficient frequency and variation across providers to be needed (and outweigh the data collection burden)? If patient preference (e.g., informed decisionmaking) is a basis for exclusion, does it impact performance and if yes, is the measure specified so that the information about patient preference and the effect on the measure is transparent?

☐ Yes (please explain below then go to Question #3)

☒ No (go to Question #3)

☐ Not applicable (i.e., there are no exclusions specified for the measure; go to Question #3)

19. Analysis of potential threats to validity: Risk-adjustment (applies to all outcome, cost, and resource use measures; may also apply to other types of measure)

☐ Not applicable (e.g., structure or process measure that is not risk-adjusted; go to Question #4)

a. Is a conceptual rationale for social risk factors included? ☐ Yes ☒ No

b. Are social risk factors included in risk model? ☐ Yes ☒ No

c. Any concerns regarding the risk-adjustment approach?

*TIPS: Consider the following: If a justification for **not risk adjusting** is provided, is there any evidence that contradicts the developer's rationale and analysis? If the developer asserts there is **no conceptual basis** for adjusting this measure for social risk factors, do you agree with the rationale? **If risk adjusted:** Are the candidate and final variables included in the risk adjustment model adequately described for the measure to be implemented? Are the candidate and final variables included in the risk adjustment model adequately described for the measure to be implemented? Are all of the risk adjustment variables present at the start of care (if not, do you agree with the rationale)? If social risk factors are not included in the risk-adjustment approach, do you agree with the developer's decision? Is an appropriate risk-adjustment strategy included in the measure (e.g., adequate model discrimination and calibration)? Are all statistical model specifications included, including a "clinical model only" if social risk factors are included in the final model?*

☒ Yes (please explain below then go to Question #4)

☐ No (go to Question #4)

Only concern is with the use of a random effects model for a procedure in which there may be a significant volume effect. Because such models can artificially shrink low-volume providers to the mean, they can alter the ordering of performance significantly, and mask any poor performance associated with low volume status.

20. Analysis of potential threats to validity: Any concerns regarding ability to identify meaningful differences in performance or overall poor performance?

☒ Yes (please explain below then go to Question #5)

☐ No (go to Question #5)

As above

21. Analysis of potential threats to validity: Any concerns regarding comparability of results if multiple data sources or methods are specified?

☐ Yes (please explain below then go to Question #6)

☐ No (go to Question #6)

☒ Not applicable (go to Question #6)

22. Analysis of potential threats to validity: Any concerns regarding missing data?

☐ Yes (please explain below then go to Question #7)

☒ No (go to Question #7)

ASSESSMENT OF MEASURE TESTING

23. Was empirical validity testing conducted using the measure as specified and appropriate statistical test?

Answer no if: face validity; only refer to clinical evidence; only descriptive statistics; only describe process for data management/cleaning/computer programming; testing does not match measure specifications (i.e. data, eMeasure, level, setting, patients).

☒ Yes (go to Question #10) [NOTE: If appropriate empirical testing has been conducted, then evaluation of face validity is not necessary. Go to Question #8 **only if** there is insufficient information provided to evaluate data element and score-level testing.]

☐ No (please explain below then go to Question #8)

24. Was face validity systematically assessed by recognized experts to determine agreement on whether the computed performance measure score from the measure as specified can be used to distinguish good and poor quality?

TIPS: Answer no if: focused on data element accuracy/availability/feasibility/other topics; the degree of consensus and any areas of disagreement not provided/discussed.

☐ Yes (go to Question #9)

☐ No (please explain below and rate Question #17: OVERALL VALIDITY as INSUFFICIENT)

25. **RATING (face validity)** - Do the face validity testing results indicate substantial agreement that the performance measure score from the measure as specified can be used to distinguish quality AND potential threats to validity are not a problem, OR are adequately addressed so results are not biased?

☐ Yes (if a NEW measure, rate Question #17: OVERALL VALIDITY as MODERATE)

☐ Yes (if a MAINTENANCE measure, do you agree with the justification for not conducting empirical testing? If no, rate Question #17: OVERALL VALIDITY as

INSUFFICIENT; otherwise, rate Question #17: OVERALL VALIDITY as MODERATE)

☐ No (please explain below and rate Question #17: OVERALL VALIDITY AS LOW)

26. Was validity testing conducted with computed performance measure scores for each measured entity?

TIPS: Answer no if: one overall score for all patients in sample used for testing patient-level data.

☐ Yes (go to Question #11)

☒ No (please explain below and go to Question #13)

Confidence interval testing was shown, but there is no validity testing of the measure score that meets the NQF recommendations for such ("Examples of validity testing of the measure score include, but are not limited to: testing hypotheses that the measures scores indicate quality of care, e.g., measure scores are different for groups known to have differences in quality assessed by another valid quality measure or method; correlation of measure scores with another valid indicator of quality for the specific topic; or relationship to conceptually related measures (e.g., scores on process measures to scores on outcome measures.)")

27. Was the method described and appropriate for assessing conceptually and theoretically sound hypothesized relationships?

TIPS: For example: correlation of the performance measure score on this measure and other performance measures; differences in performance scores between groups known to differ on quality; other accepted method with description of how it assesses validity of the performance score

☐ Yes (go to Question #12)

☐ No (please explain below, rate Question #12 as INSUFFICIENT and then go to Question #14)

28. **RATING (measure score)** - Based on the measure score results (significance, strength) and scope of testing (number of measured entities and representativeness) and analysis of potential threats, what is the level of certainty or confidence that the performance measure scores are a valid indicator of quality?

☐ High (go to Question #14)

☐ Moderate (go to Question #14)

☐ Low (please explain below then go to Question #13)

☐ Insufficient

29. Was other validity testing reported?

☒ Yes (go to Question #14)

☐ No (please explain below and rate Question #17: OVERALL VALIDITY as LOW)

30. Was validity testing conducted with patient-level data elements?

TIPS: Prior validity studies of the same data elements may be submitted

☒ Yes (go to Question #15)

☐ No (please explain below and rate Question #17: OVERALL VALIDITY as INSUFFICIENT if no score-level testing was conducted, otherwise, rate Question #17: OVERALL VALIDITY based on score-level rating from Question #12)

31. Was the method described and appropriate for assessing the accuracy of ALL critical data elements? *NOTE that data element validation from the literature is acceptable.*

TIPS: For example: Data validity/accuracy as compared to authoritative source- sensitivity, specificity, PPV, NPV; other accepted method with description of how it assesses validity of the data elements.

Answer No if: only assessed percent agreement; did not assess separately for all data elements (at least numerator, denominator, exclusions)

☐ Yes (go to Question #16)

☒ No (please explain below and rate Question #16 as INSUFFICIENT)

Only assessed percent agreement –this is OK this time given the high agreement, but will need to use other listed methods above in future submissions.

32. **RATING (data element)** - Based on the data element testing results (significance, strength) and scope of testing (number and representativeness of patients and entities) and analysis of potential threats, what is the level of certainty or confidence that the data used in the measure are valid?

☒ Moderate (if score-level testing was NOT conducted, rate Question #17: OVERALL VALIDITY as MODERATE)

☐ Low (please explain below) (if score-level testing was NOT conducted, rate Question #17: OVERALL VALIDITY as LOW)

☐ Insufficient (go to Question #17)

Please see note above – ordinarily testing only percent agreement would be unacceptable, but will rate as moderate if measure developers submit more appropriate testing with future submissions.

17. OVERALL VALIDITY RATING

OVERALL RATING OF VALIDITY taking into account the results and scope of all testing and analysis of potential threats.

☐ High (NOTE: Can be HIGH only if score-level testing has been conducted)

- ☒ Moderate (NOTE: Moderate is the highest eligible rating if score-level testing has NOT been conducted)
- ☐ Low (please explain below) [NOTE: Should rate LOW if you believe that there are threats to validity and/or threats to validity were not assessed]
- ☐ Insufficient (if insufficient, please explain below) [NOTE: For most measure types, testing at both the score level and the data element level is not required] [NOTE: If rating is INSUFFICIENT for all empirical testing, then go back to Question #8 and evaluate any face validity that was conducted, then reconsider this overall rating.]

FOR COMPOSITE MEASURES ONLY: Empirical analyses to support composite construction

What is the level of certainty or confidence that the empirical analysis demonstrates that the component measures add value to the composite and that the aggregation and weighting rules are consistent with the quality construct?

TIPS: Consider the following: Do the component measures fit the quality construct? Are the objectives of parsimony and simplicity achieved while supporting the quality construct?

- ☐ High
- ☒ Moderate
- ☐ Low (please explain below)
- ☐ Insufficient (please explain below)

Evaluation C: Scientific Acceptability

Extent to which the measure, as specified, produces consistent (reliable) and credible (valid) results about the quality of care when implemented. **Measures must be judged to meet the subcriteria for both reliability and validity to pass this criterion.**

Instructions:

37. Please complete this form for each measure you are evaluating.
38. Please pay close attention to the skip logic directions.
39. If you are unable to check a box, please highlight or shade the box for your response.
40. You must answer the "overall rating" item for both Reliability and Validity. Also, be sure to answer the composite measure question at the end of the form if your measure is a composite.
41. We have provided TIPS to help you answer the questions.
42. We've designed this form to try to minimize the amount of writing that you have to do. That said, **it is critical that you explain your thinking/rationale if you check boxes where we ask for an explanation** (because this is a Word document, you can just add your explanation below the checkbox). Feel free to add additional explanation, even if an explanation is not requested (but please type this underneath the appropriate checkbox).
43. This form is based on Algorithms 2 and 3 in the Measure Evaluation Criteria and Guidance document (see pages 18-24). These algorithms provide guidance to help you rate the Reliability and Validity subcriteria. **We ask that you refer to this document when you are evaluating your measures.**
44. Please contact Methods Panel staff if you have questions (methodspanel@qualityforum.org).

Measure Number: 3294

Measure Title: Risk-Adjusted Morbidity and Mortality for Lung Resection for Lung Cancer

RELIABILITY

45. Are submitted specifications precise, unambiguous, and complete so that they can be consistently implemented?
NOTE: NQF staff will conduct a separate, more technical, check of eMeasure (eCQM) specifications, value sets, logic, and feasibility, so no need to consider these in your evaluation.

TIPS: Consider the following: Are all the data elements clearly defined? Are all appropriate codes included? Is the logic or calculation algorithm clear? Is it likely this measure can be consistently implemented?

☒ Yes (go to Question #2)

Reliability of data elements was supported by external audit of the General Thoracic Surgery Database (GTSD) demonstrating high agreement rates and validation of data accuracy.

☐ No (please explain below, and go to Question #2) NOTE that even though **non-precise**

specifications should result in an overall LOW rating for reliability, we still want you to look at the testing results.

46. Was empirical reliability testing (at the data element or measure score level) conducted using statistical tests with the measure as specified?

TIPS: Check the 2nd "NO" box below if: only descriptive statistics provided; only describes process for data management/cleaning/computer programming; testing does not match measure specifications (i.e. data, eMeasure, level of analysis, patients)

☒ Yes (go to Question #4)

☐ No, there is reliability testing information, but *not* using statistical tests and/or not for the measure as specified OR there is no reliability testing (please explain below then go to Question #3)

47. Was **empirical VALIDITY testing** of patient-level data conducted?

☐ Yes (use your rating from data element validity testing – Question #16- under Validity Section)

☐ No (please explain below and rate Question #11: OVERALL RELIABILITY as INSUFFICIENT and proceed to the [VALIDITY SECTION](#))

48. Was reliability testing conducted with computed performance measure scores for each measured entity?

TIPS: Answer no if: only one overall score for all patients in sample used for testing patient-level data

☐ Yes (go to Question #5)

☒ No (go to Question #8)

49. Was the method described and appropriate for assessing the proportion of variability due to real differences among measured entities? NOTE: If multiple methods used, at least one must be appropriate.

TIPS: Examples of appropriate methods include signal-to-noise analysis (e.g. Adams/RAND tutorial); random split-half correlation; other accepted method with description of how it assesses reliability of the performance score.

☐ Yes (go to Question #6)

☐ No (please explain below then go to Question #8)

50. **RATING (score level)** - What is the level of certainty or confidence that the performance measure scores are reliable?

TIPS: Consider the following: Is the test sample adequate to generalize for widespread implementation?

Do the results demonstrate sufficient reliability so that differences in performance can be identified?

☐ High (go to Question #8)

☐ Moderate (go to Question #8)

☐ Low (please explain below then go to Question #7)

51. Was other reliability testing reported?

☐ Yes (go to Question #8)

☐ No (rate Question #11: OVERALL RELIABILITY as LOW and proceed to the [VALIDITY SECTION](#))

52. Was reliability testing conducted with patient-level data elements that are used to construct the performance measure?

TIPS: Prior reliability studies of the same data elements may be submitted; if comparing abstraction to “authoritative source/gold standard” see Validity Section Question #15)

☒ Yes (go to Question #9)

☐ No (if there is score-level testing, rate Question #11: OVERALL RELIABILITY based on score-level rating from Question #6; otherwise, rate Question #11: OVERALL RELIABILITY as INSUFFICIENT. Then proceed to the [VALIDITY SECTION](#))

53. Was the method described and appropriate for assessing the reliability of ALL critical data elements?

TIPS: For example: inter-abtractor agreement (ICC, Kappa); other accepted method with description of how it assesses reliability of the data elements

Answer no if: only assessed percent agreement; did not assess separately for all data elements (at least numerator, denominator, exclusions)

☐ Yes (go to Question #10)

☒ No (if no, please explain below and rate Question #10 as INSUFFICIENT)

Only agreement rates were provided in the analysis.

54. **RATING (data element)** – Based on the reliability statistic and scope of testing (number and representativeness of patients and entities), what is the level of certainty or confidence that the data used in the measure are reliable?

TIPS: Consider the following: Is the test sample adequate to generalize for widespread implementation? Can data elements be collected consistently?

☐ Moderate (if score-level testing was NOT conducted, rate Question #11: OVERALL RELIABILITY as MODERATE)

☐ Low (if score-level testing was NOT conducted, rate Question #11: OVERALL RELIABILITY as LOW)

☒ Insufficient (go to Question #11)

11. OVERALL RELIABILITY RATING

OVERALL RATING OF RELIABILITY taking into account precision of specifications and all testing results:

☐ High (NOTE: Can be HIGH only if score-level testing has been conducted)

☐ Moderate (NOTE: Moderate is the highest eligible rating if score-level testing has not been conducted)

☒ Low (please explain below) [NOTE: Should rate LOW if you believe specifications are NOT precise, unambiguous, and complete]

The submitters demonstrated a robust analysis of inter-abtractor agreement across the hospitals examined. Analysis would be much stronger if they obtained the case level data to compute a Kappa statistic to test interrater reliability.

☐ Insufficient (please explain below) [NOTE: For most measure types, testing at both the score level and the data element level is not required]

VALIDITY

ASSESSMENT OF THREATS TO VALIDITY

33. Were all potential threats to validity that are relevant to the measure empirically assessed?

TIPS: Threats to validity include: exclusions; need for risk adjustment; Able to identify statistically significant and meaningful differences; multiple sets of specifications; missing data/nonresponse.

☒ Yes (go to Question #2)

☐ No (please explain below and go to Question #2) [NOTE that even if *non-assessment of applicable threats should result in an overall INSUFFICIENT rating for validity*, we still want you to look at the testing results]

34. Analysis of potential threats to validity: Any concerns with measure exclusions?

TIPS: Consider the following: Are the exclusions consistent with the evidence? Are any patients or patient groups inappropriately excluded from the measure? Are the exclusions/exceptions of sufficient frequency and variation across providers to be needed (and outweigh the data collection burden)? If patient preference (e.g., informed decisionmaking) is a basis for exclusion, does it impact performance and if yes, is the measure specified so that the information about patient preference and the effect on the measure is transparent?

☐ Yes (please explain below then go to Question #3)

☒ No (go to Question #3)

☐ Not applicable (i.e., there are no exclusions specified for the measure; go to Question #3)

35. Analysis of potential threats to validity: Risk-adjustment (applies to all outcome, cost, and resource use measures; may also apply to other types of measure)

☐ Not applicable (e.g., structure or process measure that is not risk-adjusted; go to Question #4)

a. Is a conceptual rationale for social risk factors included? ☐ Yes ☒ No

b. Are social risk factors included in risk model? ☐ Yes ☒ No

c. Any concerns regarding the risk-adjustment approach?

*TIPS: Consider the following: If a justification for **not risk adjusting** is provided, is there any evidence that contradicts the developer's rationale and analysis? If the developer asserts there is **no conceptual basis** for adjusting this measure for social risk factors, do you agree with the rationale? **If risk adjusted:** Are the candidate and final variables included in the risk adjustment model adequately described for the measure to be implemented? Are the candidate and final variables included in the risk adjustment model adequately described for the measure to be implemented? Are all of the risk adjustment variables present at the start of care (if not, do you agree with the rationale)? If social risk factors are not included in the risk-adjustment approach, do you agree with the developer's decision? Is an appropriate risk-adjustment strategy included in the measure (e.g., adequate model discrimination and calibration)? Are all statistical model specifications included, including a "clinical model only" if social risk factors are included in the final model?*

☐ Yes (please explain below then go to Question #4)

☒ No (go to Question #4)

Risk adjustment for the clinical indicators is strongly supported. I agree with the submitters that social risk data is not available in the GTSD but would encourage the Society of Thoracic Surgeons to consider adding social risk factors to their data collection tools. Currently the GTSD does collect Primary and Secondary Payor information which could be used for Dual Eligibility stratification and possibly used as a risk adjustment.

The multivariable logistic models demonstrated statistical significance in all patient level data except Diabetes and Hypertension in all 3 models. I would question the value of leaving these in the models.

36. Analysis of potential threats to validity: Any concerns regarding ability to identify meaningful differences in performance or overall poor performance?

☐ Yes (please explain below then go to Question #5)

☒ No (go to Question #5)

The submitters validated a difference in performance using the Bayesian modeling to compare the Standardized Incidence Ratio between 231 hospitals.

37. Analysis of potential threats to validity: Any concerns regarding comparability of results if multiple data sources or methods are specified?

☐ Yes (please explain below then go to Question #6)

☐ No (go to Question #6)

☒ Not applicable (go to Question #6)

38. Analysis of potential threats to validity: Any concerns regarding missing data?

☐ Yes (please explain below then go to Question #7)

☒ No (go to Question #7)

Investigators adequately address missing data in their analysis.

ASSESSMENT OF MEASURE TESTING

39. Was empirical validity testing conducted using the measure as specified and appropriate statistical test?

Answer no if: face validity; only refer to clinical evidence; only descriptive statistics; only describe process for data management/cleaning/computer programming; testing does not match measure specifications (i.e. data, eMeasure, level, setting, patients).

☒ Yes (go to Question #10) [NOTE: If appropriate empirical testing has been conducted, then evaluation of face validity is not necessary. Go to Question #8 **only if** there is insufficient information provided to evaluate data element and score-level testing.]

☐ No (please explain below then go to Question #8)

40. Was face validity systematically assessed by recognized experts to determine agreement on whether the computed performance measure score from the measure as specified can be used to distinguish good and poor quality?

TIPS: Answer no if: focused on data element accuracy/availability/feasibility/other topics; the degree of consensus and any areas of disagreement not provided/discussed.

☐ Yes (go to Question #9)

☐ No (please explain below and rate Question #17: OVERALL VALIDITY as INSUFFICIENT)

41. **RATING (face validity)** - Do the face validity testing results indicate substantial agreement that the performance measure score from the measure as specified can be used to distinguish quality AND potential threats to validity are not a problem, OR are adequately addressed so results are not biased?

☐ Yes (if a NEW measure, rate Question #17: OVERALL VALIDITY as MODERATE)

☐ Yes (if a MAINTENANCE measure, do you agree with the justification for not conducting empirical testing? If no, rate Question #17: OVERALL VALIDITY as

INSUFFICIENT; otherwise, rate Question #17: OVERALL VALIDITY as MODERATE)

☐ No (please explain below and rate Question #17: OVERALL VALIDITY AS LOW)

42. Was validity testing conducted with computed performance measure scores for each measured entity?

TIPS: Answer no if: one overall score for all patients in sample used for testing patient-level data.

☐ Yes (go to Question #11)

☒ No (please explain below and go to Question #13)

No evidence that validity of performance score was tested. If the submitters have performed performance score testing for their previous risk-adjusted models, I would recommend updating the performance score testing with the proposed risk adjusted models.

43. Was the method described and appropriate for assessing conceptually and theoretically sound hypothesized relationships?
- TIPS: For example: correlation of the performance measure score on this measure and other performance measures; differences in performance scores between groups known to differ on quality; other accepted method with description of how it assesses validity of the performance score*
- ☐ Yes (go to Question #12)
- ☐ No (please explain below, rate Question #12 as INSUFFICIENT and then go to Question #14)
44. **RATING (measure score)** - Based on the measure score results (significance, strength) and scope of testing (number of measured entities and representativeness) and analysis of potential threats, what is the level of certainty or confidence that the performance measure scores are a valid indicator of quality?
- ☐ High (go to Question #14)
- ☐ Moderate (go to Question #14)
- ☐ Low (please explain below then go to Question #13)
- ☐ Insufficient
45. Was other validity testing reported?
- ☒ Yes (go to Question #14)
- ☐ No (please explain below and rate Question #17: OVERALL VALIDITY as LOW)
46. Was validity testing conducted with patient-level data elements?
- TIPS: Prior validity studies of the same data elements may be submitted*
- ☒ Yes (go to Question #15)
- ☐ No (please explain below and rate Question #17: OVERALL VALIDITY as INSUFFICIENT if no score-level testing was conducted, otherwise, rate Question #17: OVERALL VALIDITY based on score-level rating from Question #12)
47. Was the method described and appropriate for assessing the accuracy of ALL critical data elements? *NOTE that data element validation from the literature is acceptable.*
- TIPS: For example: Data validity/accuracy as compared to authoritative source- sensitivity, specificity, PPV, NPV; other accepted method with description of how it assesses validity of the data elements.*
- Answer No if: only assessed percent agreement; did not assess separately for all data elements (at least numerator, denominator, exclusions)*
- ☐ Yes (go to Question #16)
- ☒ No (please explain below and rate Question #16 as INSUFFICIENT)
- Only agreement rates were assessed.
48. **RATING (data element)** - Based on the data element testing results (significance, strength) and scope of testing (number and representativeness of patients and entities) and analysis of potential threats, what is the level of certainty or confidence that the data used in the measure are valid?
- ☐ Moderate (if score-level testing was NOT conducted, rate Question #17: OVERALL VALIDITY as MODERATE)
- ☒ Low (please explain below) (if score-level testing was NOT conducted, rate Question #17: OVERALL VALIDITY as LOW)
- It would be a much stronger analysis if the developer obtained the case level results to provide a kappa statistic.
- ☐ Insufficient (go to Question #17)

17. OVERALL VALIDITY RATING

OVERALL RATING OF VALIDITY taking into account the results and scope of all testing and analysis of potential threats.

☐ High (NOTE: Can be HIGH only if score-level testing has been conducted)

☐ Moderate (NOTE: Moderate is the highest eligible rating if score-level testing has NOT been conducted)

☒ Low (please explain below) [NOTE: Should rate LOW if you believe that there are threats to validity and/or threats to validity were not assessed]

No testing for threats to validity evident in the information provided by the submitters

☐ Insufficient (if insufficient, please explain below) [NOTE: For most measure types, testing at both the score level and the data element level is not required] [NOTE: If rating is INSUFFICIENT for all empirical testing, then go back to Question #8 and evaluate any face validity that was conducted, then reconsider this overall rating.]

FOR COMPOSITE MEASURES ONLY: Empirical analyses to support composite construction

What is the level of certainty or confidence that the empirical analysis demonstrates that the component measures add value to the composite and that the aggregation and weighting rules are consistent with the quality construct?

TIPS: Consider the following: Do the component measures fit the quality construct? Are the objectives of parsimony and simplicity achieved while supporting the quality construct?

☐ High

☒ Moderate

The statistical analysis supports the use of the Mortality or Major Morbidity Composite Model for risk adjustment and performance measurement. Although, the referenced article did show only fair performance of the composite model using the C-statistic results. I would recommend the submitters include the referenced article in their submission materials.

Fernandez FG, Kosinski AS, Burfeind W, Park B, DeCamp MM, Seder C, Marshall B, Magee MJ, Wright CD, Kozower BD. The Society of Thoracic Surgeons Lung Cancer Resection Risk Model: Higher Quality Data and Superior Outcomes. Ann Thorac Surg. 2016 Aug;102(2):370-7.

☐ Low (please explain below)

☐ Insufficient (please explain below)

Evaluation D: Scientific Acceptability

Extent to which the measure, as specified, produces consistent (reliable) and credible (valid) results about the quality of care when implemented. **Measures must be judged to meet the subcriteria for both reliability and validity to pass this criterion.**

Instructions:

55. Please complete this form for each measure you are evaluating.
56. Please pay close attention to the skip logic directions.
57. If you are unable to check a box, please highlight or shade the box for your response.
58. You must answer the "overall rating" item for both Reliability and Validity. Also, be sure to answer the composite measure question at the end of the form if your measure is a composite.
59. We have provided TIPS to help you answer the questions.
60. We've designed this form to try to minimize the amount of writing that you have to do. That said, **it is critical that you explain your thinking/rationale if you check boxes where we ask for an explanation** (because this is a Word

document, you can just add your explanation below the checkbox). Feel free to add additional explanation, even if an explanation is not requested (but please type this underneath the appropriate checkbox).

61. This form is based on Algorithms 2 and 3 in the Measure Evaluation Criteria and Guidance document (see pages 18-24). These algorithms provide guidance to help you rate the Reliability and Validity subcriteria. **We ask that you refer to this document when you are evaluating your measures.**
62. Please contact Methods Panel staff if you have questions (methodspanel@qualityforum.org).

Measure Number: 3294

Measure Title: STS Lobectomy for Lung Cancer Composite Score

RELIABILITY

63. Are submitted specifications precise, unambiguous, and complete so that they can be consistently implemented?

NOTE: NQF staff will conduct a separate, more technical, check of eMeasure (eCQM) specifications, value sets, logic, and feasibility, so no need to consider these in your evaluation.

TIPS: Consider the following: Are all the data elements clearly defined? Are all appropriate codes included? Is the logic or calculation algorithm clear? Is it likely this measure can be consistently implemented?

☒ Yes (go to Question #2)

☐ No (please explain below, and go to Question #2) NOTE that even though **non-precise**

specifications should result in an overall LOW rating for reliability, we still want you to look at the testing results.

64. Was empirical reliability testing (at the data element or measure score level) conducted using statistical tests with the measure as specified?

TIPS: Check the 2nd "NO" box below if: only descriptive statistics provided; only describes process for data management/cleaning/computer programming; testing does not match measure specifications (i.e. data, eMeasure, level of analysis, patients)

☒ Yes (go to Question #4)

☐ No, there is reliability testing information, but *not* using statistical tests and/or not for the measure as specified OR there is no reliability testing (please explain below then go to Question #3)

65. Was **empirical VALIDITY testing** of patient-level data conducted?

☐ Yes (use your rating from data element validity testing – Question #16- under Validity Section)

☐ No (please explain below and rate Question #11: OVERALL RELIABILITY as INSUFFICIENT and proceed to the [VALIDITY SECTION](#))

66. Was reliability testing conducted with computed performance measure scores for each measured entity?

TIPS: Answer no if: only one overall score for all patients in sample used for testing patient-level data

☒ Yes (go to Question #5)

☐ No (go to Question #8)

67. Was the method described and appropriate for assessing the proportion of variability due to real differences among measured entities? *NOTE: If multiple methods used, at least one must be appropriate.*

TIPS: Examples of appropriate methods include signal-to-noise analysis (e.g. Adams/RAND tutorial); random split-half correlation; other accepted method with description of how it assesses reliability of the performance score.

☒ Yes (go to Question #6)

☐ No (please explain below then go to Question #8)

68. **RATING (score level)** - What is the level of certainty or confidence that the performance measure scores are reliable?

TIPS: Consider the following: Is the test sample adequate to generalize for widespread implementation?

Do the results demonstrate sufficient reliability so that differences in performance can be identified?

☐ High (go to Question #8)

☒ Moderate (go to Question #8)

☐ Low (please explain below then go to Question #7)

69. Was other reliability testing reported?

☐ Yes (go to Question #8)

☐ No (rate Question #11: OVERALL RELIABILITY as LOW and proceed to the [VALIDITY SECTION](#))

70. Was reliability testing conducted with patient-level data elements that are used to construct the performance measure?

TIPS: Prior reliability studies of the same data elements may be submitted; if comparing abstraction to "authoritative source/gold standard" see Validity Section Question #15)

☐ Yes (go to Question #9)

☒ No (if there is score-level testing, rate Question #11: OVERALL RELIABILITY based on score-level rating from Question #6; otherwise, rate Question #11: OVERALL RELIABILITY as INSUFFICIENT. Then proceed to the [VALIDITY SECTION](#))

71. Was the method described and appropriate for assessing the reliability of ALL critical data elements?

TIPS: For example: inter-abtractor agreement (ICC, Kappa); other accepted method with description of how it assesses reliability of the data elements

Answer no if: only assessed percent agreement; did not assess separately for all data elements (at least numerator, denominator, exclusions)

☐ Yes (go to Question #10)

☐ No (if no, please explain below and rate Question #10 as INSUFFICIENT)

My one concern with the reliability of the data elements is changes in the registry reporting platform over time. Opening up new reporting options may reduce reliability of data over time.

72. **RATING (data element)** – Based on the reliability statistic and scope of testing (number and representativeness of patients and entities), what is the level of certainty or confidence that the data used in the measure are reliable?

TIPS: Consider the following: Is the test sample adequate to generalize for widespread implementation? Can data elements be collected consistently?

☐ Moderate (if score-level testing was NOT conducted, rate Question #11: OVERALL RELIABILITY as MODERATE)

☐ Low (if score-level testing was NOT conducted, rate Question #11: OVERALL RELIABILITY as LOW)

☐ Insufficient (go to Question #11)

11. OVERALL RELIABILITY RATING

OVERALL RATING OF RELIABILITY taking into account precision of specifications and all testing results:

☐ High (NOTE: Can be HIGH only if score-level testing has been conducted)

☒ Moderate (NOTE: Moderate is the highest eligible rating if score-level testing has not been conducted)

☐ Low (please explain below) [NOTE: Should rate LOW if you believe specifications are NOT precise,

unambiguous, and complete]

☐ Insufficient (please explain below) [NOTE: For most measure types, testing at both the score level and the data element level is not required]

The measure is tested at the hospital level. The measure summary form indicates that it can be used for hospitals or group practices, but I do not see any evidence of reliability testing with group practice data.

VALIDITY

ASSESSMENT OF THREATS TO VALIDITY

49. Were all potential threats to validity that are relevant to the measure empirically assessed?

TIPS: Threats to validity include: exclusions; need for risk adjustment; Able to identify statistically significant and meaningful differences; multiple sets of specifications; missing data/nonresponse.

☒ Yes (go to Question #2)

☐ No (please explain below and go to Question #2) [NOTE that even if **non-assessment of applicable threats should result in an overall INSUFFICIENT rating for validity**, we still want you to look at the testing results]

50. Analysis of potential threats to validity: Any concerns with measure exclusions?

TIPS: Consider the following: Are the exclusions consistent with the evidence? Are any patients or patient groups inappropriately excluded from the measure? Are the exclusions/exceptions of sufficient frequency and variation across providers to be needed (and outweigh the data collection burden)? If patient preference (e.g., informed decisionmaking) is a basis for exclusion, does it impact performance and if yes, is the measure specified so that the information about patient preference and the effect on the measure is transparent?

☒ Yes (please explain below then go to Question #3)

☐ No (go to Question #3)

☐ Not applicable (i.e., there are no exclusions specified for the measure; go to Question #3)

I was slightly concerned about dropping cases with missing discharge mortality status because I cannot tell if this introduces selection bias or offers an opportunity for gaming. I'm assuming this is a relatively rare event, although I didn't see the number of cases dropped in either the Composite Measure Testing worksheet or the journal article.

51. Analysis of potential threats to validity: Risk-adjustment (applies to all outcome, cost, and resource use measures; may also apply to other types of measure)

☐ Not applicable (e.g., structure or process measure that is not risk-adjusted; go to Question #4)

a. Is a conceptual rationale for social risk factors included? ☒ Yes ☐ No

b. Are social risk factors included in risk model? ☐ Yes ☒ No

c. Any concerns regarding the risk-adjustment approach?

*TIPS: Consider the following: If a justification for **not risk adjusting** is provided, is there any evidence that contradicts the developer's rationale and analysis? If the developer asserts there is **no conceptual basis** for adjusting this measure for social risk factors, do you agree with the rationale? **If risk adjusted:** Are the candidate and final variables included in the risk adjustment model adequately described for the measure to be implemented? Are the candidate and final variables included in the risk adjustment model adequately described for the measure to be implemented? Are all of the risk adjustment variables present at the start of care (if not, do you agree with the rationale)? If social risk factors are not included in the risk-adjustment approach, do you agree with the developer's decision? Is an appropriate risk-adjustment strategy included in the measure (e.g., adequate model discrimination and calibration)? Are all statistical model specifications included, including a "clinical model only" if social risk factors are included in the final model?*

☐ Yes (please explain below then go to Question #4)

☒ No (go to Question #4)

The one thing to note with a Bayesian risk adjustment model is the tendency for scores to fall in the middle of the distribution. We see this here with 91.4 percent of cases ending up with 2 stars. One strength of the risk model is that covariates were selected on an a-priori or theoretical basis and retained in the model regardless of impact rather than through a data driven process. The model c-statistics are modest, but not unexpected for clinical data.

52. Analysis of potential threats to validity: Any concerns regarding ability to identify meaningful differences in performance or overall poor performance?

☐ Yes (please explain below then go to Question #5)

☒ No (go to Question #5)

53. Analysis of potential threats to validity: Any concerns regarding comparability of results if multiple data sources or methods are specified?

☐ Yes (please explain below then go to Question #6)

☐ No (go to Question #6)

☒ Not applicable (go to Question #6)

54. Analysis of potential threats to validity: Any concerns regarding missing data?

☒ Yes (please explain below then go to Question #7)

☐ No (go to Question #7)

As I mentioned above, I have concerns about potential selection bias for sites with missing mortality information. It would be helpful to know the number of excluded cases – I assume it is small and random.

ASSESSMENT OF MEASURE TESTING

55. Was empirical validity testing conducted using the measure as specified and appropriate statistical test?

Answer no if: face validity; only refer to clinical evidence; only descriptive statistics; only describe process for data management/cleaning/computer programming; testing does not match measure specifications (i.e. data, eMeasure, level, setting, patients).

☒ Yes (go to Question #10) [NOTE: If appropriate empirical testing has been conducted, then evaluation of face validity is not necessary. Go to Question #8 **only if** there is insufficient information provided to evaluate data element and score-level testing.]

☐ No (please explain below then go to Question #8)

56. Was face validity systematically assessed by recognized experts to determine agreement on whether the computed performance measure score from the measure as specified can be used to distinguish good and poor quality?

TIPS: Answer no if: focused on data element accuracy/availability/feasibility/other topics; the degree of consensus and any areas of disagreement not provided/discussed.

☐ Yes (go to Question #9)

☐ No (please explain below and rate Question #17: OVERALL VALIDITY as INSUFFICIENT)

57. **RATING (face validity)** - Do the face validity testing results indicate substantial agreement that the performance measure score from the measure as specified can be used to distinguish quality AND potential threats to validity are not a problem, OR are adequately addressed so results are not biased?

☐ Yes (if a NEW measure, rate Question #17: OVERALL VALIDITY as MODERATE)

☐ Yes (if a MAINTENANCE measure, do you agree with the justification for not

conducting empirical testing? If no, rate Question #17: OVERALL VALIDITY as INSUFFICIENT; otherwise, rate Question #17: OVERALL VALIDITY as MODERATE)

☐ No (please explain below and rate Question #17: OVERALL VALIDITY AS LOW)

58. Was validity testing conducted with computed performance measure scores for each measured entity?

TIPS: Answer no if: one overall score for all patients in sample used for testing patient-level data.

☒ Yes (go to Question #11)

☐ No (please explain below and go to Question #13)

59. Was the method described and appropriate for assessing conceptually and theoretically sound hypothesized relationships?

TIPS: For example: correlation of the performance measure score on this measure and other performance measures; differences in performance scores between groups known to differ on quality; other accepted method with description of how it assesses validity of the performance score

☒ Yes (go to Question #12)

☐ No (please explain below, rate Question #12 as INSUFFICIENT and then go to Question #14)

This was done by looking at the relationship between observed rates of the two outcomes (mortality and major complications) and the overall star rating for the hospital. As the authors point out, there is a clear linear relationship between observed components. Worth noting, the 95% confidence intervals for the mortality measure almost overlap for the 1-star and 2-star groups. If hospitals with lower volume were included in the analysis these two groups may not be distinct.

Grouping measure scores by star rating helps confirm that the composite is not driven by a single measure and that both measures move together. However, as the authors point out, the major morbidity measure drives the variance. This is not surprising since it is made up of 9 medical complications and is itself a composite of sorts. It would be helpful to see a confirmatory factor analysis or structural measurement model to better understand how all 10 items relate to each other. Finally, it would be beneficial to have an external measure of adverse events after lobectomy or a broader category of lung surgeries to group hospitals (i.e., an independent measure that is not part of the composite).

60. **RATING (measure score)** - Based on the measure score results (significance, strength) and scope of testing (number of measured entities and representativeness) and analysis of potential threats, what is the level of certainty or confidence that the performance measure scores are a valid indicator of quality?

☐ High (go to Question #14)

☒ Moderate (go to Question #14)

☐ Low (please explain below then go to Question #13)

☐ Insufficient

61. Was other validity testing reported?

☐ Yes (go to Question #14)

☐ No (please explain below and rate Question #17: OVERALL VALIDITY as LOW)

62. Was validity testing conducted with patient-level data elements?

TIPS: Prior validity studies of the same data elements may be submitted

☒ Yes (go to Question #15)

☐ No (please explain below and rate Question #17: OVERALL VALIDITY as INSUFFICIENT if no score-level testing was conducted, otherwise, rate Question #17: OVERALL VALIDITY based on score-level rating from Question #12)

The authors provide information on overall validity testing for the General Thoracic Surgery Database in 2016, 2011 and 2010. In the narrative they refer to auditing 10% of sites for completeness, but only 15 lobectomy cases for accuracy. This leads to confusion with the table shown on pages 9-10 that shows a total of 500 cases for many data elements. It is not clear what this table is reporting at the data element level.

63. Was the method described and appropriate for assessing the accuracy of ALL critical data elements? *NOTE that data element validation from the literature is acceptable.*

TIPS: For example: Data validity/accuracy as compared to authoritative source- sensitivity, specificity, PPV, NPV; other accepted method with description of how it assesses validity of the data elements.

Answer No if: only assessed percent agreement; did not assess separately for all data elements (at least numerator, denominator, exclusions)

☒ Yes (go to Question #16)

☐ No (please explain below and rate Question #16 as INSUFFICIENT)

The method is appropriate, but as noted above, it is difficult to know if the agreement rates shown in the table are correct given the miss-match between the numbers in the table and text. It is also not clear why there is no one who has both the auditor's rating and the site level data to calculate a kappa statistic. This seems like a key component in assessing and maintain database integrity over time.

64. **RATING (data element)** - Based on the data element testing results (significance, strength) and scope of testing (number and representativeness of patients and entities) and analysis of potential threats, what is the level of certainty or confidence that the data used in the measure are valid?

☐ Moderate (if score-level testing was NOT conducted, rate Question #17: OVERALL VALIDITY as MODERATE)

☒ Low (please explain below) (if score-level testing was NOT conducted, rate Question #17: OVERALL VALIDITY as LOW)

☐ Insufficient (go to Question #17)

This rating is based on the concerns with the miss-match between the agreement rates in the text and tables as well as the lack of a kappa statistic. In addition, it would be helpful to know if the current data reporting options and auditing requirements for 2016 will carry forward to 2017 and beyond. Changes in these methods could adversely affect the validity of future data in the STS database.

17. OVERALL VALIDITY RATING

OVERALL RATING OF VALIDITY taking into account the results and scope of all testing and analysis of potential threats.

☐ High (NOTE: Can be HIGH only if score-level testing has been conducted)

☒ Moderate (NOTE: Moderate is the highest eligible rating if score-level testing has NOT been conducted)

☐ Low (please explain below) [NOTE: Should rate LOW if you believe that there are threats to validity and/or threats to validity were not assessed]

☐ Insufficient (if insufficient, please explain below) [NOTE: For most measure types, testing at both the score level and the data element level is not required] [NOTE: If rating is INSUFFICIENT for all empirical testing, then go back to Question #8 and evaluate any face validity that was conducted, then reconsider this overall rating.]

FOR COMPOSITE MEASURES ONLY: Empirical analyses to support composite construction

What is the level of certainty or confidence that the empirical analysis demonstrates that the component measures add value to the composite and that the aggregation and weighting rules are consistent with the quality construct?

TIPS: Consider the following: Do the component measures fit the quality construct? Are the objectives of parsimony and simplicity achieved while supporting the quality construct?

- ☐ High
- ☒ Moderate
- ☐ Low (please explain below)
- ☐ Insufficient (please explain below)

Overall this is a well thought out measure. It builds on the STS registry, which captures the vast majority of cases among participating members and is subjected to an independent auditing process. As the authors point out, not all lobectomies are performed by cardio-thoracic surgeons. From a 'public benefit' perspective, it would be helpful to include all relevant surgeries in the measure, not just the ones performed by a specific type of surgeon. Obviously this is not possible with the risk adjustment model used for the measure, but would be something to consider for the future.

It is important that the measure includes a minimum number of cases (N=30) since the reliability is modest for low case volumes. The composite score is a logical combination of a number of closely related outcomes. The standardization and weighting are strengths of the overall measure. The reliability testing was appropriate and shows modest reliability with relatively low sample sizes. The distribution of participant's composite scores for lobectomy in Figure 1 of Kozower et al. (2016) shows graphically that the measure is able to differentiate performance above and below the mean. Worth noting, composite scores are already relatively high, offering relatively limited room for improvement. Also the Bayesian risk adjustment results push many hospitals to the middle of the distribution, resulting in clear differentiation between high and low performers for a relatively small percent of the overall sample. Finally, this is a composite measure made up of two different measures, each of which captures adverse events after surgery. The measure could be improved by allowing all 10 adverse events to be standardized and weighted individually.

Kozower BD, O'Brein SM, Kosinski AS, et al. (2016). The Society of Thoracic Surgeons Composite Score for Rating Program Performance for Lobectomy for Lung Cancer. *Annals of Thoracic Surgery*, 101: 1379-1387.

Evaluation E: Scientific Acceptability

Extent to which the measure, as specified, produces consistent (reliable) and credible (valid) results about the quality of care when implemented. **Measures must be judged to meet the subcriteria for both reliability and validity to pass this criterion.**

Instructions:

73. Please complete this form for each measure you are evaluating.
74. Please pay close attention to the skip logic directions.
75. If you are unable to check a box, please highlight or shade the box for your response.
76. You must answer the "overall rating" item for both Reliability and Validity. Also, be sure to answer the composite measure question at the end of the form if your measure is a composite.
77. We have provided TIPS to help you answer the questions.
78. We've designed this form to try to minimize the amount of writing that you have to do. That said, ***it is critical that you explain your thinking/rationale if you check boxes where we ask for an explanation*** (because this is a Word document, you can just add your explanation below the checkbox). Feel free to add additional explanation, even if an explanation is not requested (but please type this underneath the appropriate checkbox).
79. This form is based on Algorithms 2 and 3 in the Measure Evaluation Criteria and Guidance document (see pages 18-24). These algorithms provide guidance to help you rate the Reliability and Validity subcriteria. ***We ask that you refer to this document when you are evaluating your measures.***
80. Please contact Methods Panel staff if you have questions (methodspanel@qualityforum.org).

Measure Number: 3294

Measure Title: STS Lobectomy for Lung Cancer Composite Score

RELIABILITY

81. Are submitted specifications precise, unambiguous, and complete so that they can be consistently implemented?

NOTE: NQF staff will conduct a separate, more technical, check of eMeasure (eCQM) specifications, value sets, logic, and feasibility, so no need to consider these in your evaluation.

TIPS: Consider the following: Are all the data elements clearly defined? Are all appropriate codes included? Is the logic or calculation algorithm clear? Is it likely this measure can be consistently implemented?

☒ Yes (go to Question #2)

☐ No (please explain below, and go to Question #2) *NOTE that even though **non-precise specifications should result in an overall LOW rating for reliability**, we still want you to look at the testing results.*

82. Was empirical reliability testing (at the data element or measure score level) conducted using statistical tests with the measure as specified?

TIPS: Check the 2nd "NO" box below if: only descriptive statistics provided; only describes process for data management/cleaning/computer programming; testing does not match measure specifications (i.e. data, eMeasure, level of analysis, patients)

☒ Yes (go to Question #4)

☐ No, there is reliability testing information, but *not* using statistical tests and/or not for the measure as specified OR there is no reliability testing (please explain below then go to Question #3)

83. Was **empirical** **VALIDITY** testing of patient-level data conducted?

☐ Yes (use your rating from data element validity testing – Question #16- under Validity Section)

☐ No (please explain below and rate Question #11: OVERALL RELIABILITY as INSUFFICIENT and proceed to the [VALIDITY SECTION](#))

84. Was reliability testing conducted with computed performance measure scores for each measured entity?

TIPS: Answer no if: only one overall score for all patients in sample used for testing patient-level data

☒ Yes (go to Question #5)

☐ No (go to Question #8)

85. Was the method described and appropriate for assessing the proportion of variability due to real differences among measured entities? *NOTE: If multiple methods used, at least one must be appropriate.*

TIPS: Examples of appropriate methods include signal-to-noise analysis (e.g. Adams/RAND tutorial); random split-half correlation; other accepted method with description of how it assesses reliability of the performance score.

☒ Yes (go to Question #6)

☐ No (please explain below then go to Question #8)

86. **RATING (score level)** - What is the level of certainty or confidence that the performance measure scores are reliable?

TIPS: Consider the following: Is the test sample adequate to generalize for widespread implementation?

Do the results demonstrate sufficient reliability so that differences in performance can be identified?

☐ High (go to Question #8)

☒ Moderate (go to Question #8)

☐ Low (please explain below then go to Question #7)

87. Was other reliability testing reported?

☐ Yes (go to Question #8)

☐ No (rate Question #11: OVERALL RELIABILITY as LOW and proceed to the [VALIDITY SECTION](#))

88. Was reliability testing conducted with patient-level data elements that are used to construct the performance measure?

TIPS: Prior reliability studies of the same data elements may be submitted; if comparing abstraction to “authoritative source/gold standard” see Validity Section Question #15)

☐ Yes (go to Question #9)

☒ No (if there is score-level testing, rate Question #11: OVERALL RELIABILITY based on score-level rating from Question #6; otherwise, rate Question #11: OVERALL RELIABILITY as INSUFFICIENT. Then proceed to the [VALIDITY SECTION](#))

Data not provided in the submission, but may be available in STS database.

89. Was the method described and appropriate for assessing the reliability of ALL critical data elements?

TIPS: For example: inter-abtractor agreement (ICC, Kappa); other accepted method with description of how it assesses reliability of the data elements

Answer no if: only assessed percent agreement; did not assess separately for all data elements (at least numerator, denominator, exclusions)

☐ Yes (go to Question #10)

☐ No (if no, please explain below and rate Question #10 as INSUFFICIENT)

90. **RATING (data element)** – Based on the reliability statistic and scope of testing (number and representativeness of patients and entities), what is the level of certainty or confidence that the data used in the measure are reliable?

TIPS: Consider the following: Is the test sample adequate to generalize for widespread implementation? Can data elements be collected consistently?

☐ Moderate (if score-level testing was NOT conducted, rate Question #11: OVERALL RELIABILITY as MODERATE)

☐ Low (if score-level testing was NOT conducted, rate Question #11: OVERALL RELIABILITY as LOW)

☐ Insufficient (go to Question #11)

11. OVERALL RELIABILITY RATING

OVERALL RATING OF RELIABILITY taking into account precision of specifications and all testing results:

☐ High (NOTE: Can be HIGH only if score-level testing has been conducted)

☒ Moderate (NOTE: Moderate is the highest eligible rating if score-level testing has not been conducted)

☐ Low (please explain below) [NOTE: Should rate LOW if you believe specifications are NOT precise, unambiguous, and complete]

☐ Insufficient (please explain below) [NOTE: For most measure types, testing at both the score level and the data element level is not required]

It appears reliability is moderately good, and improves, as expected, with increasing number of cases.

VALIDITY

ASSESSMENT OF THREATS TO VALIDITY

65. Were all potential threats to validity that are relevant to the measure empirically assessed?

TIPS: Threats to validity include: exclusions; need for risk adjustment; Able to identify statistically significant and meaningful differences; multiple sets of specifications; missing data/nonresponse.

☒ Yes (go to Question #2)

☐ No (please explain below and go to Question #2) [NOTE that even if *non-assessment of applicable threats should result in an overall INSUFFICIENT rating for validity*, we still want you to look at the testing results]

66. Analysis of potential threats to validity: Any concerns with measure exclusions?

TIPS: Consider the following: Are the exclusions consistent with the evidence? Are any patients or patient groups inappropriately excluded from the measure? Are the exclusions/exceptions of sufficient frequency and variation across providers to be needed (and outweigh the data collection burden)? If patient preference (e.g., informed decisionmaking) is a basis for exclusion, does it impact performance and if yes, is the measure specified so that the information about patient preference and the effect on the measure is transparent?

☐ Yes (please explain below then go to Question #3)

☒ No (go to Question #3)

☐ Not applicable (i.e., there are no exclusions specified for the measure; go to Question #3)

No major concerns, but social disparities explicitly ignored, with explanation.

67. Analysis of potential threats to validity: Risk-adjustment (applies to all outcome, cost, and resource use measures; may also apply to other types of measure)

☐ Not applicable (e.g., structure or process measure that is not risk-adjusted; go to Question #4)

a. Is a conceptual rationale for social risk factors included? ☐ Yes ☒ No

b. Are social risk factors included in risk model? ☐ Yes ☒ No

c. Any concerns regarding the risk-adjustment approach?

*TIPS: Consider the following: If a justification for **not risk adjusting** is provided, is there any evidence that contradicts the developer's rationale and analysis? If the developer asserts there is **no conceptual basis** for adjusting this measure for social risk factors, do you agree with the rationale? **If risk adjusted:** Are the candidate and final variables included in the risk adjustment model adequately described for the measure to be implemented? Are the candidate and final variables included in the risk adjustment model adequately described for the measure to be implemented? Are all of the risk adjustment variables present at the start of care (if not, do you agree with the rationale)? If social risk factors are not included in the risk-adjustment approach, do you agree with the developer's decision? Is an appropriate risk-adjustment strategy included in the measure (e.g., adequate model discrimination and calibration)? Are all statistical model specifications included, including a "clinical model only" if social risk factors are included in the final model?*

☐ Yes (please explain below then go to Question #4)

☒ No (go to Question #4)

68. Analysis of potential threats to validity: Any concerns regarding ability to identify meaningful differences in performance or overall poor performance?

☐ Yes (please explain below then go to Question #5)

☒ No (go to Question #5)

69. Analysis of potential threats to validity: Any concerns regarding comparability of results if multiple data sources or methods are specified?

☐ Yes (please explain below then go to Question #6)

☐ No (go to Question #6)

☒ Not applicable (go to Question #6)

70. Analysis of potential threats to validity: Any concerns regarding missing data?

☐ Yes (please explain below then go to Question #7)

☒ No (go to Question #7)

ASSESSMENT OF MEASURE TESTING

71. Was empirical validity testing conducted using the measure as specified and appropriate statistical test?

Answer no if: face validity; only refer to clinical evidence; only descriptive statistics; only describe process for data management/cleaning/computer programming; testing does not match measure specifications (i.e. data, eMeasure, level, setting, patients).

☒ Yes (go to Question #10) [NOTE: If appropriate empirical testing has been conducted, then evaluation of face validity is not necessary. Go to Question #8 **only if** there is insufficient information provided to evaluate data element and score-level testing.]

☐ No (please explain below then go to Question #8)

72. Was face validity systematically assessed by recognized experts to determine agreement on whether the computed performance measure score from the measure as specified can be used to distinguish good and poor quality?

TIPS: Answer no if: focused on data element accuracy/availability/feasibility/other topics; the degree of consensus and any areas of disagreement not provided/discussed.

☐ Yes (go to Question #9)

☐ No (please explain below and rate Question #17: OVERALL VALIDITY as INSUFFICIENT)

73. **RATING (face validity)** - Do the face validity testing results indicate substantial agreement that the performance measure score from the measure as specified can be used to distinguish quality AND potential threats to validity are not a problem, OR are adequately addressed so results are not biased?

☐ Yes (if a NEW measure, rate Question #17: OVERALL VALIDITY as MODERATE)

☐ Yes (if a MAINTENANCE measure, do you agree with the justification for not conducting empirical testing? If no, rate Question #17: OVERALL VALIDITY as

INSUFFICIENT; otherwise, rate Question #17: OVERALL VALIDITY as MODERATE)

☐ No (please explain below and rate Question #17: OVERALL VALIDITY AS LOW)

74. Was validity testing conducted with computed performance measure scores for each measured entity?

TIPS: Answer no if: one overall score for all patients in sample used for testing patient-level data.

☐ Yes (go to Question #11)

☒ No (please explain below and go to Question #13)

NQF recommends testing hypotheses that the measure scores indicate quality of care, e.g., measure scores differ by groups known to have differences in quality assessed by another valid quality measure or method; or by correlation of measure scores with another valid indicator of quality for a specific topic; or relationship to conceptually similar measures. This submission reported (mostly) separated confidence intervals but no 'anchor' against which to judge validity.

75. Was the method described and appropriate for assessing conceptually and theoretically sound hypothesized relationships?

TIPS: For example: correlation of the performance measure score on this measure and other performance measures; differences in performance scores between groups known to differ on quality; other accepted method with description of how it assesses validity of the performance score

☐ Yes (go to Question #12)

☐ No (please explain below, rate Question #12 as INSUFFICIENT and then go to Question #14)

76. **RATING (measure score)** - Based on the measure score results (significance, strength) and scope of testing (number of measured entities and representativeness) and analysis of potential threats, what is the level of certainty or confidence that the performance measure scores are a valid indicator of quality?

- ☐ High (go to Question #14)
- ☐ Moderate (go to Question #14)
- ☐ Low (please explain below then go to Question #13)
- ☐ Insufficient

77. Was other validity testing reported?

- ☒ Yes (go to Question #14)
- ☐ No (please explain below and rate Question #17: OVERALL VALIDITY as LOW)

78. Was validity testing conducted with patient-level data elements?

TIPS: Prior validity studies of the same data elements may be submitted

- ☒ Yes (go to Question #15)
- ☐ No (please explain below and rate Question #17: OVERALL VALIDITY as INSUFFICIENT if no score-level testing was conducted, otherwise, rate Question #17: OVERALL VALIDITY based on score-level rating from Question #12)

79. Was the method described and appropriate for assessing the accuracy of ALL critical data elements? *NOTE that data element validation from the literature is acceptable.*

TIPS: For example: Data validity/accuracy as compared to authoritative source- sensitivity, specificity, PPV, NPV; other accepted method with description of how it assesses validity of the data elements.

Answer No if: only assessed percent agreement; did not assess separately for all data elements (at least numerator, denominator, exclusions)

- ☐ Yes (go to Question #16)
- ☒ No (please explain below and rate Question #16 as INSUFFICIENT)

Kappa statistic for case-level data would help clarify confusion in the submission regarding agreement rates. This can probably be clarified in a follow-up submission

80. **RATING (data element)** - Based on the data element testing results (significance, strength) and scope of testing (number and representativeness of patients and entities) and analysis of potential threats, what is the level of certainty or confidence that the data used in the measure are valid?

- ☒ Moderate (if score-level testing was NOT conducted, rate Question #17: OVERALL VALIDITY as MODERATE)
- ☐ Low (please explain below) (if score-level testing was NOT conducted, rate Question #17: OVERALL VALIDITY as LOW)
- ☐ Insufficient (go to Question #17)

17. OVERALL VALIDITY RATING

OVERALL RATING OF VALIDITY taking into account the results and scope of all testing and analysis of potential threats.

- ☐ High (NOTE: Can be HIGH only if score-level testing has been conducted)
- ☒ Moderate (NOTE: Moderate is the highest eligible rating if score-level testing has NOT been conducted)
- ☐ Low (please explain below) [NOTE: Should rate LOW if you believe that there are threats to validity and/or threats to validity were not assessed]
- ☐ Insufficient (if insufficient, please explain below) [NOTE: For most measure types, testing at both the

score level and the data element level is not required] [NOTE: If rating is INSUFFICIENT for all empirical testing, then go back to Question #8 and evaluate any face validity that was conducted, then reconsider this overall rating.]

I suspect that further detail from available information will render this as an acceptable, reliable and valid measure
FOR COMPOSITE MEASURES ONLY: Empirical analyses to support composite construction

What is the level of certainty or confidence that the empirical analysis demonstrates that the component measures add value to the composite and that the aggregation and weighting rules are consistent with the quality construct?

TIPS: Consider the following: Do the component measures fit the quality construct? Are the objectives of parsimony and simplicity achieved while supporting the quality construct?

☐ High

☒ Moderate

☐ Low (please explain below)

☐ Insufficient (please explain below)

Some concern that there are wide and almost overlapping confidence intervals for the mortality outcome between 1-star and 2-star hospitals. Low volume hospitals, with lower reliability, would likely overlap. With so many hospitals classified in the middle group, this may not be a highly-differentiating outcome measure at the end of the day....but it seems conceptually and structurally sound

Committee Pre-evaluation Comments: Criteria 2: Scientific Acceptability of Measure Properties (including all 2a, 2b, and 2c)

2a1. Reliability specifications

**STS database. No concerns about implementation.

**Reliable

**Well-defined

**No issues.

2a2. Reliability testing

**STS database. No

**No

**No concerns

**No issues

2b1. Validity Testing

**The most recent audits of the General Thoracic Surgery Database have demonstrated a high degree of data validity.

Overall data accuracy rates have increased substantially since audits of the GTSD were first conducted in 2010;

agreement ranges have also narrowed, indicating greater consistency in data accuracy among audited sites.

The rates of missing data were low and are getting lower. We therefore concluded that systematic missing data did not lead to bias in our measure no threat to validity

**No concerns

**Data abstracted from clinical records - minimal concerns re: data validity

**No issues

2b2.-3. Other threats to validity

**Risk adjustment rigorous

**Adequate – the usual problem with random effects models of squishing outcomes towards the mean, especially for low-volume groups

2c. Composite Analysis

**Fits quality construct and rationale

STS's combined mortality and morbidity model for pulmonary resection for lung cancer is important and appropriate for public reporting for the following reasons:

- 1.) within the broad category of lung cancer resections, lobectomy is the single most common major procedure that a thoracic surgeon performs;
- 2.) These procedures are therefore useful and appropriate to use as a benchmark for performance by general thoracic surgery programs. By providing surgeons and teams with risk-adjusted results, they can identify how they are performing compared with other programs in the STS General Thoracic Database,

Most recently, STS surgeon members have expressed interest in real-time, online data updates, which has led to the development of dashboard-type reporting on STS.org. The general thoracic dashboard is scheduled for launch in 2018.

**Composite measure credibly reflects pt. experience

**Yes and yes

**No issues.

Criterion 3. [Feasibility](#)

3. Feasibility is the extent to which the specifications including measure logic, require data that are readily available or could be captured without undue burden and can be implemented for performance measurement.

- The developer reports that data are generated or collected by and used by healthcare personnel during the provision of care; coded by someone other than the person obtaining original information; and abstracted from a record by someone other than the person obtaining the original information.
- All data are in defined fields in a combination of electronic sources
- Data are collected continuously by the participating sites and harvested by DCRI twice a year; reports are then sent back to participating sites about three months after harvest. Participating sites generally have data managers on staff.
- The developer reports that STS GTSD participant surgeons pay an annual participant fee of \$550 or \$700 depending on whether the participant is an STS member.

Questions for the Committee:

- *Are the required data elements routinely generated and used during care delivery?*
- *Are the required data elements available in electronic form, e.g., EHR or other electronic sources?*
- *Is the data collection strategy ready to be put into operational use?*

Preliminary rating for feasibility: ☐ High ☒ Moderate ☐ Low ☐ Insufficient

Committee Pre-evaluation Comments: Criteria 3: Feasibility

3. Feasibility

**STS database feasibility good

**feasible

**Feasible through the GTSD

**Requires participation in the registry -- impossible to replicate/participate otherwise.

Criterion 4: [Usability and Use](#)

4a. Use (4a1. Accountability and Transparency; 4a2. Feedback on measure)

4a. Use evaluate the extent to which audiences (e.g., consumers, purchasers, providers, policymakers) use or could use performance results for both accountability and performance improvement activities.

4a.1. 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 (or the data on performance results are available). If not in use at the time of initial endorsement, then a credible plan for implementation within the specified timeframes is provided.

Current uses of the measure

Publicly reported? ☒ Yes ☐ No

Current use in an accountability program? ☒ Yes ☐ No ☐ UNCLEAR

OR

Planned use in an accountability program? ☐ Yes ☐ No

Accountability program details

- The measure results are shared with participants in the STS General Thoracic Surgery Database (GTSD) for quality improvement purposes. In addition, the developer reports active promotion of STS measures through the STS Public Reporting Task force. The task force develops public report cards that are consumer centric.

4a.2. Feedback on the measure by those being measured or others. Three criteria demonstrate feedback: 1) those being measured have been given performance results or data, as well as assistance with interpreting the measure results and data; 2) those being measured and other users have been given an opportunity to provide feedback on the measure performance or implementation; 3) this feedback has been considered when changes are incorporated into the measure

Feedback on the measure by those being measured or others

- The developer states that STS surgeon members have expressed interest in real-time, online data updates which led to the development of a general thoracic dashboard. The dashboard is scheduled for launch in 2018.
- The developer states that given the recent launch of public reporting that they have not received sufficient feedback from non-participants to be able to assess the impact of the public reporting initiative.

Additional Feedback:

- The developer reports that surgeons on the STS General Thoracic Surgery Task Force meet periodically to discuss participant reports and discuss enhancements to the GTS database. Additions and clarifications to the data collection form and the content/format of participant reports are discussed and implemented as appropriate.

Questions for the Committee:

- How can the performance results be used to further the goal of high-quality, efficient healthcare?*
- How has the measure been vetted in real-world settings by those being measured or others?*

Preliminary rating for Use: ☒ Pass ☐ No Pass

4b. Usability (4a1. Improvement; 4a2. Benefits of measure)

4b. Usability evaluate the extent to which audiences (e.g., consumers, purchasers, providers, policymakers) use or could use performance results for both accountability and performance improvement activities.

4b.1 Improvement. Progress toward achieving the goal of high-quality, efficient healthcare for individuals or populations is demonstrated.

Improvement results

- The developer reports that operative mortality in the STS General Thoracic Surgery Database (GTSD) decreased from 2.2% (from 2002-2008) to 1.4% (from 2012-2014). Further, when data from the GTSD were compared with

the Nationwide Inpatient Sample database from 2002 to 2008, patients in the GTSD had lower unadjusted mortality rates, median length of stay, and lower pulmonary complication rates for lobectomy.

4b2. Benefits vs. harms. 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).

Unexpected findings (positive or negative) during implementation

- The developer reports they are unaware of any unexpected findings associated with the implementation of this measure.

Potential harms

- The developer reports that the rate of major morbidity has increased from 8.6% to 9.1% from 2002 to 2008 which is potentially explained by more complete coding of complications by data abstractors and inclusion of unexpected return to the operating room for any reason.

Questions for the Committee:

- *How can the performance results be used to further the goal of high-quality, efficient healthcare?*
- *Do the benefits of the measure outweigh any potential unintended consequences?*

Preliminary rating for Usability and use: ☐ High ☒ Moderate ☐ Low ☐ Insufficient

Committee Pre-evaluation Comments: Criteria 4: Usability and Use

4a1. Use

**Most recently, STS surgeon members have expressed interest in real-time, online data updates, which has led to the development of dashboard-type reporting on STS.org. The general thoracic dashboard is scheduled for launch in 2018. Star ratings for surgeons and hospitals will be developed

**usable

**Measure is not being used in an accountability program but is being publicly reported

**Already publicly reported

4b1. Usability

**Believe the benefits outweigh unintended consequences. Recommend Approval

**No concerns. Separately I am worried about additive value of this measure compared to measure 1790

**Overall benefits outweigh harms

**Yes

Criterion 5: [Related and Competing Measures](#)

Related or competing measures

- 1790 Risk-Adjusted Morbidity and Mortality for Lung Resection for Lung Cancer
- The developer notes that NQF 1790 is related conceptually to 3294 and that the numerators for both measures include the same list of postoperative complications, but the outcomes for the Lobectomy Composite measure are grouped into two domains (operative mortality and major complications) and the measure is structured to provide general thoracic surgeons with a "star rating."
- Measure #1790 includes a broader range of lung resection procedures than the Lobectomy Composite, and therefore includes a larger number of cases and potentially provides performance data to more general thoracic surgeons.

Harmonization

- The developer reports that NQF 1790 and 3294 are harmonized to the extent possible.

Public and Member Comments

Comments and Member Support/Non-Support Submitted as of: January 23, 2018

- No NQF members have submitted support/non-support choices as of this date. No comments have been submitted as of this date.

1. Evidence and Performance Gap – Importance to Measure and Report

Extent to which the specific measure focus is evidence-based, important to making significant gains in healthcare quality, and improving health outcomes for a specific high-priority (high-impact) aspect of healthcare where there is variation in or overall less-than-optimal performance. **Measures must be judged to meet all sub criteria to pass this criterion and be evaluated against the remaining criteria.**

1a. Evidence to Support the Measure Focus – See attached Evidence Submission Form

NQF_evidence_attachment_STS-3294-111517.docx

1a.1 For Maintenance of Endorsement: Is there new evidence about the measure since the last update/submission?

Do not remove any existing information. If there have been any changes to evidence, the Committee will consider the new evidence. Please use the most current version of the evidence attachment (v7.1). Please use red font to indicate updated evidence.

No

1a Evidence (subcriterion 1a)

Measure Number (if previously endorsed):

Measure Title: STS Lobectomy for Lung Cancer Composite Score

IF the measure is a component in a composite performance measure, provide the title of the Composite Measure here:

Date of Submission: 11/15/2017

Instructions

- Complete 1a.1 and 1a.2 for all measures. If instrument-based measure, complete 1a.3.
- Complete **EITHER 1a.2, 1a.3 or 1a.4** as applicable for the type of measure and evidence.
- For composite performance measures:
 - A separate evidence form is required for each component measure unless several components were studied together.
 - If a component measure is submitted as an individual performance measure, attach the evidence form to the individual measure submission.
- All information needed to demonstrate meeting the evidence subcriterion (1a) must be in this form. An appendix of supplemental materials may be submitted, but there is no guarantee it will be reviewed.
- If you are unable to check a box, please highlight or shade the box for your response.
- Contact NQF staff regarding questions. Check for resources at [Submitting Standards webpage](#).

Note: The information provided in this form is intended to aid the Standing Committee and other stakeholders in understanding to what degree the evidence for this measure meets NQF's evaluation criteria.

1a. Evidence to Support the Measure Focus

The measure focus is evidence-based, demonstrated as follows:

- Outcome: ³ 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.
- Intermediate clinical outcome: a systematic assessment and grading of the quantity, quality, and consistency of the body of evidence ⁴ that the measured intermediate clinical outcome leads to a desired health outcome.

- **Process:** ⁵ a systematic assessment and grading of the quantity, quality, and consistency of the body of evidence ⁴ that the measured process leads to a desired health outcome.
- **Structure:** a systematic assessment and grading of the quantity, quality, and consistency of the body of evidence ⁴ that the measured structure leads to a desired health outcome.
- **Efficiency:** ⁶ evidence not required for the resource use component.
- For measures derived from patient reports, evidence should demonstrate that the target population values the measured outcome, process, or structure and finds it meaningful.
- **Process measures incorporating Appropriate Use Criteria:** See NQF's guidance for evidence for measures, in general; guidance for measures specifically based on clinical practice guidelines apply as well.

Notes

3. Generally, rare event outcomes do not provide adequate information for improvement or discrimination; however, serious reportable events that are compared to zero are appropriate outcomes for public reporting and quality improvement.

4. The preferred systems for grading the evidence are the Grading of Recommendations, Assessment, Development and Evaluation ([GRADE guidelines](#)) and/or modified GRADE.

5. 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 multistep process, the step with the strongest evidence for the link to the desired outcome should be selected as the focus of measurement. Note: A measure focused only on collecting PROM data is not a PRO-PM.

6. Measures of efficiency combine the concepts of resource use and quality (see NQF's [Measurement Framework: Evaluating Efficiency Across Episodes of Care](#); [AQA Principles of Efficiency Measures](#)).

1a.1. This is a measure of: (should be consistent with type of measure entered in De.1)

Outcome

☒ Outcome: Two domains of outcomes are measured: 1. Operative Mortality (death during the same hospitalization as surgery or within 30 days of the procedure), and 2. Presence of at least one of these major complications: pneumonia, acute respiratory distress syndrome, bronchopleural fistula, pulmonary embolus, initial ventilator support greater than 48 hours, reintubation/respiratory failure, tracheostomy, myocardial infarction, or unexpected return to the operating room.

☐ Patient-reported outcome (PRO):

PROs include HRQoL/functional status, symptom/symptom burden, experience with care, health-related behaviors. (A PRO-based performance measure is not a survey instrument. Data may be collected using a survey instrument to construct a PRO measure.)

☐ Intermediate clinical outcome (e.g., lab value):

☐ Process:

☐ Appropriate use measure:

☐ Structure:

☐ Composite:

1a.2 LOGIC MODEL Diagram or briefly describe the steps between the healthcare structures and processes (e.g., interventions, or services) and the patient's health outcome(s). The relationships in the diagram should be easily understood by general, non-technical audiences. Indicate the structure, process or outcome being measured.

Postoperative complications and operative mortality are important negative outcomes associated with lung cancer resection surgery, including lobectomy, the most frequently performed lung resection procedure. The STS lung cancer resection risk model (Fernandez et al, 2016) identifies predictors of these outcomes, including patient age, smoking status, comorbid medical conditions, and other patient characteristics, as well as operative approach and the extent of pulmonary resection. Knowledge of these predictors informs clinical decision making by enabling physicians and

patients to understand the associations between individual patient characteristics and outcomes and – with continuous feedback of performance data over time – fosters quality improvement.

Fernandez FG, Kosinski AS, Burfeind W, et al. The Society of Thoracic Surgeons lung cancer resection risk model: higher quality data and superior outcomes. *Ann Thorac Surg* 2016;102:370-7.

1a.3 Value and Meaningfulness: IF this measure is derived from patient report, provide evidence that the target population values the measured **outcome, process, or structure** and finds it meaningful. (Describe how and from whom their input was obtained.)

n/a

****RESPOND TO ONLY ONE SECTION BELOW -EITHER 1a.2, 1a.3 or 1a.4) ****

1a.2 FOR OUTCOME MEASURES including PATIENT REPORTED OUTCOMES - Provide empirical data demonstrating the relationship between the outcome (or PRO) to at least one healthcare structure, process, intervention, or service.

Data in the STS General Thoracic Surgery Database (GTSD) have demonstrated a reduction in perioperative morbidity and equivalent long-term survival when minimally invasive approaches for lobectomy are used instead of a standard thoracotomy. Specifically, STS data have shown that minimally invasive lung cancer resection has a 50% reduction in major complications compared with a thoracotomy approach, adjusted for age, sex, and comorbidities. There is a general consensus among STS surgeons and the STS GTSD task force that stage I lung cancer is usually resectable with a minimally invasive approach. Because many patients desire a minimally invasive approach, and STS data and other published data demonstrate improved risk-adjusted outcomes, the STS considers it appropriate to include the percent of minimally invasive lobectomies for stage I lung cancer as a process measure on STS biannual reports to GTSD participants.

Kozower BD, O'Brien SM, Kosinski AS, et al. The Society of Thoracic Surgeons composite score for rating program performance for lobectomy for lung cancer. *Ann Thorac Surg* 2016;101:1379-87.

1a.3. SYSTEMATIC REVIEW(SR) OF THE EVIDENCE (for INTERMEDIATE OUTCOME, PROCESS, OR STRUCTURE PERFORMANCE MEASURES, INCLUDING THOSE THAT ARE INSTRUMENT-BASED) If the evidence is not based on a systematic review go to section 1a.4) If you wish to include more than one systematic review, add additional tables.

What is the source of the systematic review of the body of evidence that supports the performance measure? A systematic review is a scientific investigation that focuses on a specific question and uses explicit, prespecified 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)

☐ Clinical Practice Guideline recommendation (with evidence review)

☐ US Preventive Services Task Force Recommendation

☐ Other systematic review and grading of the body of evidence (e.g., *Cochrane Collaboration, AHRQ Evidence Practice Center*)

☐ Other

Source of Systematic Review: <ul style="list-style-type: none"> • Title • Author • Date • Citation, including page number • URL 	
Quote the guideline or recommendation verbatim about the process, structure or intermediate outcome being measured. If not a guideline, summarize the conclusions from the SR.	
Grade assigned to the evidence associated with the recommendation with the definition of the grade	
Provide all other grades and definitions from the evidence grading system	
Grade assigned to the recommendation with definition of the grade	
Provide all other grades and definitions from the recommendation grading system	
Body of evidence: <ul style="list-style-type: none"> • Quantity – how many studies? • Quality – what type of studies? 	
Estimates of benefit and consistency across studies	
What harms were identified?	
Identify any new studies conducted since the SR. Do the new studies change the conclusions from the SR?	

1a.4 OTHER SOURCE OF EVIDENCE

If source of evidence is NOT from a clinical practice guideline, USPSTF, or systematic review, please describe the evidence on which you are basing the performance measure.

1a.4.1 Briefly SYNTHESIZE the evidence that supports the measure. A list of references without a summary is not acceptable.

1a.4.2 What process was used to identify the evidence?

1a.4.3. Provide the citation(s) for the evidence.

1b. Performance Gap

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.

1b.1. Briefly explain the rationale for this measure (e.g., how the measure will improve the quality of care, the benefits or improvements in quality envisioned by use of this measure)

If a COMPOSITE (e.g., combination of component measure scores, all-or-none, any-or-none), SKIP this question and answer the composite questions.

n/a

1b.2. Provide performance scores on the measure as specified (current and over time) at the specified level of analysis. (*This is required for maintenance of endorsement.* Include mean, std dev, min, max, interquartile range, scores

by decile. Describe the data source including number of measured entities; number of patients; dates of data; if a sample, characteristics of the entities include.) This information also will be used to address the sub-criterion on improvement (4b1) under Usability and Use.

The measure was calculated in two overlapping 3-year time periods, January 1, 2014 – December 31, 2016 and January 1, 2013 – December 31, 2015. For each time period, we provide the number of measured entities (No. of participants), the number of eligible patient records (No. of operations), and the distribution of composite score estimates by percentiles and geographic region. We present results for all the participants and for the subset of participants with at least 30 eligible cases.

	January 1, 2013 – December 31, 2015 All participants =30 cases		January 1, 2014 – December 31, 2016 All participants =30 cases	
No. of participants	242	185	233	186
No. of operations	23594	22752	24912	24318
Mean	0.972	0.972	0.973	0.974
SD	0.007	0.008	0.006	0.007
IQR	0.008	0.009	0.007	0.009
Minimum	0.945	0.945	0.953	0.953
10%	0.961	0.96	0.965	0.965
20%	0.967	0.967	0.969	0.968
30%	0.97	0.969	0.971	0.971
40%	0.971	0.971	0.973	0.973
50%	0.973	0.973	0.974	0.975
60%	0.974	0.975	0.976	0.976
70%	0.975	0.976	0.977	0.977
80%	0.977	0.978	0.979	0.979
90%	0.979	0.98	0.981	0.982
Maximum	0.988	0.988	0.987	0.987
Midwest	49	38	48	38
Northeast	71	52	67	51
South	82	65	82	67
West	40	30	36	30

1b.3. If no or limited performance data on the measure as specified is reported in 1b2, then provide a summary of data from the literature that indicates opportunity for improvement or overall less than optimal performance on the specific focus of measurement.

n/a (see data reported in 1b2)

1b.4. Provide disparities data from the measure as specified (current and over time) by population group, e.g., by race/ethnicity, gender, age, insurance status, socioeconomic status, and/or disability. *(This is required for maintenance of endorsement. Describe the data source including number of measured entities; number of patients; dates of data; if a*

sample, characteristics of the entities included.) For measures that show high levels of performance, i.e., “topped out”, disparities data may demonstrate an opportunity for improvement/gap in care for certain sub-populations. This information also will be used to address the sub-criterion on improvement (4b1) under Usability and Use.

DATES: Jan. 1, 2014 - Dec. 31, 2016

INCIDENCE N= 33,326

DEMOGRAPHICS

Age (years)

Mean	65.7
Median	67.0
25th Percentile	59.0
75th Percentile	73.0

Gender, Female 54.7%

Race

Caucasian	84.9%
Black	8.8%
Asian	2.9%
Native American	0.3%
Native Hawaiian/Pac Islander	0.2%
Other	2.5%
Multiple Races	0.7%
Missing/unknown/not documented	1.1%

1b.5. If no or limited data on disparities from the measure as specified is reported in 1b.4, then provide a summary of data from the literature that addresses disparities in care on the specific focus of measurement. Include citations. Not necessary if performance data provided in 1b.4

n/a (see data reported in 1b.4)

1c. Composite Quality Construct and Rationale

1c.1. A composite performance measure is a combination of two or more component measures, each of which individually reflects quality of care, into a single performance measure with a single score.

For purposes of NQF measure submission, evaluation, and endorsement, the following will be considered composites:

- Measures with two or more individual performance measure scores combined into one score for an accountable entity.
- Measures with two or more individual component measures assessed separately for each patient and then aggregated into one score for an accountable entity:
 - all-or-none measures (e.g., all essential care processes received, or outcomes experienced, by each patient);

1c.1. Please identify the composite measure construction: two or more individual performance measure scores combined into one score

1c.2. Describe the quality construct, including:

- the overall area of quality
- included component measures and
- the relationship of the component measures to the overall composite and to each other.

The STS Lobectomy Composite Score measures surgical performance for patients treated with lobectomy for lung cancer. Similar to other STS composite measures, this measure is based on a combination of an operative mortality outcome measure and the risk-adjusted occurrence of any of several major complications. To assess overall quality, the composite comprises the following two domains:

1. Operative Mortality (death during the same hospitalization as surgery or within 30 days of the procedure)
2. Presence of at least one of these major complications: pneumonia, acute respiratory distress syndrome, bronchopleural fistula, pulmonary embolus, initial ventilator support greater than 48 hours, reintubation/respiratory failure, tracheostomy, myocardial infarction, or unexpected return to the operating room.

Participants receive a score for each of the two domains, plus an overall composite score. The overall composite score was created by a weighted combination of the above two domains. In addition to receiving a numeric score, participants are assigned to rating categories designated by one to three stars:

1 star: lower-than-expected performance

2 stars: as-expected-performance

3 stars: higher-than-expected-performance

1c.3. Describe the rationale for constructing a composite measure, including how the composite provides a distinctive or additive value over the component measures individually.

Risk-adjusted mortality has historically been the dominant outcomes metric for thoracic surgery, but in an era when the average mortality rates for these procedures have declined to very low levels, it can be difficult to differentiate performance based on mortality alone. Specifically, mortality alone fails to take into account the fact that not all operative survivors received equal quality care, e.g., patients who survive surgery but are debilitated by a major postoperative complication. Calculating a composite score from a weighted combination of operative mortality and major complications provides a more comprehensive measure of overall surgical quality.

1c.4. Describe how the aggregation and weighting of the component measures are consistent with the stated quality construct and rationale.

The composite score is created by a weighted combination of two domains (operative mortality and major complications) resulting in a single composite score. Operative mortality is weighted approximately four times that of a major complication in the composite, consistent with the STS adult cardiac surgery quality measures. The STS General Thoracic Surgery Database working group believes this is an improvement from its previous lung cancer resection model in which mortality and major morbidity were weighted equally.

For more information on the STS composite methodology, please see the attachment:

Kozower BD, O'Brien SM, Kosinski AS, et al. The Society of Thoracic Surgeons composite score for rating program performance for lobectomy for lung cancer. *Ann Thorac Surg* 2016;101:1379-87.

Composite Measure Testing (subcriteria 2a2, 2b1-2b6)

Measure Number (if previously endorsed):

Composite Measure Title: STS Lobectomy for Lung Cancer Composite Score

Date of Submission: 11/15/2017

Composite Construction:

- ☒ Two or more individual performance measure scores combined into one score
- ☐ All-or-none measures (e.g., all essential care processes received or outcomes experienced by each patient)

Instructions: Please contact NQF staff before you begin.

- If a component measure is submitted as an individual performance measure, the non-composite measure testing form must also be completed and attached to the individual measure submission.
- Measures must be tested for all the data sources and levels of analyses that are specified. ***If there is more than one set of data specifications or more than one level of analysis, contact NQF staff*** about how to present all the testing information in one form.
- **Sections 1, 2a2, 2b1, 2b2, and 2b4 must be completed.**
- **For composites with outcome and resource use measures**, section **2b3** also must be completed.
- If specified for **multiple data sources/sets of specifications** (e.g., claims and EHRs), section **2b5** also must be completed.
- Respond to all questions as instructed with answers immediately following the question. All information on testing to demonstrate meeting the subcriteria for reliability (2a2) and validity (2b1-2b6) and composites (2c) must be in this form. An appendix for *supplemental* materials may be submitted, but there is no guarantee it will be reviewed.
- If you are unable to check a box, please highlight or shade the box for your response.
- Maximum of 25 pages (*including questions/instructions*; minimum font size 11 pt; do not change margins). **Contact NQF staff if more pages are needed.**
- Contact NQF staff regarding questions. Check for resources at [Submitting Standards webpage](#).
- For information on the most updated guidance on how to address social risk **factors variables** and testing in this form refer to the release notes for version 7.1 of the Measure Testing Attachment. and the 2017 Measure Evaluation Criteria and Guidance.

Note: The information provided in this form is intended to aid the Standing Committee and other stakeholders in understanding to what degree the testing results for this measure meet NQF's evaluation criteria for testing.

2a2. Reliability testing ¹⁰ demonstrates the measure data elements are repeatable, producing the same results a high proportion of the time when assessed in the same population in the same time period and/or that the measure score is precise. For **instrument-based measures** (including **PRO-PMs**) and **composite performance measures**, reliability should be demonstrated for the computed performance score.

2b1. Validity testing ¹¹ demonstrates that the measure data elements are correct and/or the measure score correctly reflects the quality of care provided, adequately identifying differences in quality. For **instrument based measures (including PRO-PMs) and composite performance measures**, validity should be demonstrated for the computed performance score.

2b2. Exclusions are supported by the clinical evidence and are of sufficient frequency to warrant inclusion in the specifications of the measure; ¹²

AND

If patient preference (e.g., informed decisionmaking) is a basis for exclusion, there must be evidence that the exclusion impacts performance on the measure; in such cases, the measure must be specified so that the information about patient preference and the effect on the measure is transparent (e.g., numerator category computed separately, denominator exclusion category computed separately). ¹³

2b3. For outcome measures and other measures when indicated (e.g., resource use):

- **an evidence-based risk-adjustment strategy** (e.g., risk models, risk stratification) is specified; is based on patient factors (including clinical and social risk factors) that influence the measured outcome and are present at start of care; ^{14,15} and has demonstrated adequate discrimination and calibration

OR

- rationale/data support no risk adjustment/ stratification.

2b4. Data analysis of computed measure scores demonstrates that methods for scoring and analysis of the specified measure allow for **identification of statistically significant and practically/clinically meaningful** ¹⁶ **differences in performance;**

OR

there is evidence of overall less-than-optimal performance.

2b5. If multiple data sources/methods are specified, there is demonstration they produce comparable results.

2b6. Analyses identify the extent and distribution of **missing data** (or nonresponse) and demonstrate that performance results are not biased due to systematic missing data (or differences between responders and nonresponders) and how the specified handling of missing data minimizes bias.

2c. For composite performance measures, empirical analyses support the composite construction approach and demonstrate that:

2c1. the component measures fit the quality construct and add value to the overall composite while achieving the related objective of parsimony to the extent possible; and

2c2. the aggregation and weighting rules are consistent with the quality construct and rationale while achieving the related objective of simplicity to the extent possible.

(if not conducted or results not adequate, justification must be submitted and accepted)

Notes

10. Reliability testing applies to both the data elements and computed measure score. Examples of reliability testing for data elements 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 of the measure score addresses precision of measurement (e.g., signal-to-noise).

11. Validity testing applies to both the data elements and computed measure score. Validity testing of data elements typically analyzes agreement with another authoritative source of the same information. Examples of validity testing of the measure score include, but are not limited to: testing hypotheses that the measures scores indicate quality of care, e.g., measure scores are different for groups known to have differences in quality assessed by another valid quality measure or method; correlation of measure scores with another valid indicator of quality for the specific topic; or relationship to conceptually related measures (e.g., scores on process measures to scores on outcome measures). Face validity of the measure score as a quality indicator may be adequate if accomplished through a systematic and transparent process, by identified experts, and explicitly addresses whether performance scores resulting from the measure as specified can be used to distinguish good from poor quality. The degree of consensus and any areas of disagreement must be provided/discussed.

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

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

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

15. 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 percent v. 75 percent) 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 less-than-optimal performance may not demonstrate much variability across providers.

1. DATA/SAMPLE USED FOR ALL TESTING OF THIS MEASURE

Often the same data are used for all aspects of measure testing. In an effort to eliminate duplication, the first five questions apply to all measure testing. If there are differences by aspect of testing, (e.g., reliability vs. validity) be sure to indicate the specific differences in question 1.7.

1.1. What type of data was used for testing? (Check all the sources of data identified in the measure specifications and data used for testing the measure. Testing must be provided for all the sources of data specified and intended for measure implementation. **If different data sources are used for different components in the composite, indicate the component after the checkbox. If different data sources are used for the numerator and denominator, indicate N [numerator] or D [denominator] after the checkbox.**)

Measure Specified to Use Data From: (must be consistent with data sources entered in S.17)	Measure Tested with Data From:
<input type="checkbox"/> abstracted from paper record	<input type="checkbox"/> abstracted from paper record
<input type="checkbox"/> claims	<input type="checkbox"/> claims
<input checked="" type="checkbox"/> registry	<input checked="" type="checkbox"/> registry
<input type="checkbox"/> abstracted from electronic health record	<input type="checkbox"/> abstracted from electronic health record
<input type="checkbox"/> eMeasure (HQMF) implemented in EHRs	<input type="checkbox"/> eMeasure (HQMF) implemented in EHRs
<input type="checkbox"/> other:	<input type="checkbox"/> other:

1.2. If an existing dataset was used, identify the specific dataset (the dataset used for testing must be consistent with the measure specifications for target population and healthcare entities being measured; e.g., Medicare Part A claims, Medicaid claims, other commercial insurance, nursing home MDS, home health OASIS, clinical registry).

STS General Thoracic Surgery Database, Version 2.3

1.3. What are the dates of the data used in testing? 01/01/2014 – 12/31/2016

1.4. What levels of analysis were tested? (testing must be provided for all the levels specified and intended for measure implementation, e.g., individual clinician, hospital, health plan)

Measure Specified to Measure Performance of: (must be consistent with levels entered in item S.20)	Measure Tested at Level of:
<input type="checkbox"/> individual clinician	<input type="checkbox"/> individual clinician
<input checked="" type="checkbox"/> group/practice	<input checked="" type="checkbox"/> group/practice
<input checked="" type="checkbox"/> hospital/facility/agency	<input checked="" type="checkbox"/> hospital/facility/agency
<input type="checkbox"/> health plan	<input type="checkbox"/> health plan
<input type="checkbox"/> other:	<input type="checkbox"/> other:

1.5. How many and which measured entities were included in the testing and analysis (by level of analysis and data source)? (identify the number and descriptive characteristics of measured entities included in the analysis (e.g., size, location, type); if a sample was used, describe how entities were selected for inclusion in the sample)

The analysis population consisted of all STS records for patients meeting measure inclusion criteria who had their surgery during January 1, 2014 through December 31, 2016. The population included 24,912 patient records from 233 hospitals.

1.6. How many and which patients were included in the testing and analysis (by level of analysis and data source)? (identify the number and descriptive characteristics of patients included in the analysis (e.g., age, sex, race, diagnosis); if a sample was used, describe how patients were selected for inclusion in the sample)

Includes 24,912 eligible patients. Patient characteristics are below.

Age (years), mean (SD)	67.3 (9.5)
Male	44.6%
Body Mass Index (kg/m ²), mean, (SD)	27.6 (6.1)
Hypertension	62.0%
Steroid therapy	3.0%
Congestive heart failure	2.5%
Coronary artery disease	20.6%
Peripheral vascular disease	8.9%
Reoperation	5.5%
Preoperative chemotherapy within 6 months	6.5%
Cerebrovascular disease	7.6%
Diabetes mellitus	18.7%
Renal failure	1.1%
Dialysis	0.5%
Cigarette smoking	
Never smoked	15.3%
Past smoker	61.7%
Current smoker	23.0%
Forced expiratory volume in 1 second percent of predicted	84.5 (19.7)
Zubrod score	
0	45.9%
1	50.2%
2	3.2%
3	0.6%
4	0.1%
5	<0.1%
ASA Class	
0	0.2%
2	15.2%
3	76.3%
4	8.3%
5	<0.1%
Pathologic stage	
0	71.0%
I	17.1%
II	10.4%
IV	1.5%
Year of operation	
2014	32.1%
2015	34.1%
2016	33.8%

1.7. If there are differences in the data or sample used for different aspects of testing (e.g., reliability, validity, exclusions, risk adjustment), identify how the data or sample are different for each aspect of testing reported below.

The STS tests reliability based on three years of data in the General Thoracic Surgery Database (see 1.5 above). Validity testing is conducted on an annual basis through the audit of data completeness and accuracy in randomly-selected surgical records at randomly-selected GTSD participant sites (see 2b1.2 below).

1.8 What were the social risk factors that were available and analyzed? For example, patient-reported data (e.g., income, education, language), proxy variables when social risk data are not collected from each patient (e.g. census tract), or patient community characteristics (e.g. percent vacant housing, crime rate) which do not have to be a proxy for patient-level data.

Patient social risk data are not collected in the General Thoracic Surgery Database. Through the collection of insurance information, information on dual Medicare/Medicaid eligibility is available from the database, which can serve as a proxy for low income and patient vulnerability. However, this information is not presently included in STS data analysis nor as a basis for stratification in STS measures.

2a2. RELIABILITY TESTING

Note: If accuracy/correctness (validity) of data elements was empirically tested, separate reliability testing of data elements is not required – in 2a2.1 check critical data elements; in 2a2.2 enter “see section 2b2 for validity testing of data elements”; and skip 2a2.3 and 2a2.4.

2a2.1. What level of reliability testing was conducted? (may be one or both levels)

Note: Current guidance for composite measure evaluation states that reliability must be demonstrated for the composite performance measure score.

☒ **Performance measure score** (e.g., signal-to-noise analysis)

2a2.2. Describe the method of reliability testing and what it tests (describe the steps—do not just name a method; what type of error does it test; what statistical analysis was used)

Reliability is conventionally defined as the proportion of variation in a measure that is due to true between-unit differences (i.e., signal) as opposed to random statistical fluctuations (i.e., noise). Equivalently, it is the squared correlation between a measurement and the true value. Accordingly, reliability was defined as the square of the Pearson correlation coefficient (ρ^2) between the set of participant-specific estimates

$\hat{\theta}_1, \dots, \hat{\theta}_N$ and the corresponding unknown true values, $\theta_1, \dots, \theta_N$, that is:

$$\rho^2 = \frac{\sum_{j=1}^N (\hat{\theta}_j - \frac{1}{N} \sum_{h=1}^N \hat{\theta}_h) (\theta_j - \frac{1}{N} \sum_{h=1}^N \theta_h)}{\sum_{j=1}^N (\hat{\theta}_j - \frac{1}{N} \sum_{h=1}^N \hat{\theta}_h)^2 \sum_{j=1}^N (\theta_j - \frac{1}{N} \sum_{h=1}^N \theta_h)^2}$$

The quantity ρ^2 was estimated by its posterior mean, namely,

$$\hat{\rho}^2 = \frac{1}{5000} \sum_{l=1}^{5000} \rho_{(l)}^2$$

where

$$\rho_{(l)}^2 = \frac{\sum_{j=1}^N (\hat{\theta}_j - \frac{1}{N} \sum_{h=1}^N \hat{\theta}_h) (\theta_j^{(l)} - \frac{1}{N} \sum_{h=1}^N \theta_h^{(l)})}{\sum_{j=1}^N (\hat{\theta}_j - \frac{1}{N} \sum_{h=1}^N \hat{\theta}_h)^2 \sum_{j=1}^N (\theta_j^{(l)} - \frac{1}{N} \sum_{h=1}^N \theta_h^{(l)})^2}$$

with $\theta_h^{(l)}$ denoting the value of θ_j on the l -th MCMC sample $\hat{\theta}_j = \sum_{l=1}^{5000} \theta_j^{(l)} / 5000$ denoting the posterior mean of θ_j . A 95% credible interval for ρ^2 was obtained by calculating the 125th smallest and 125th largest values of $\rho_{(l)}^2$ across the 5,000 MCMC samples.

2a2.3. What were the statistical results from reliability testing? (e.g., percent agreement and kappa for the critical data elements; distribution of reliability statistics from a signal-to-noise analysis)

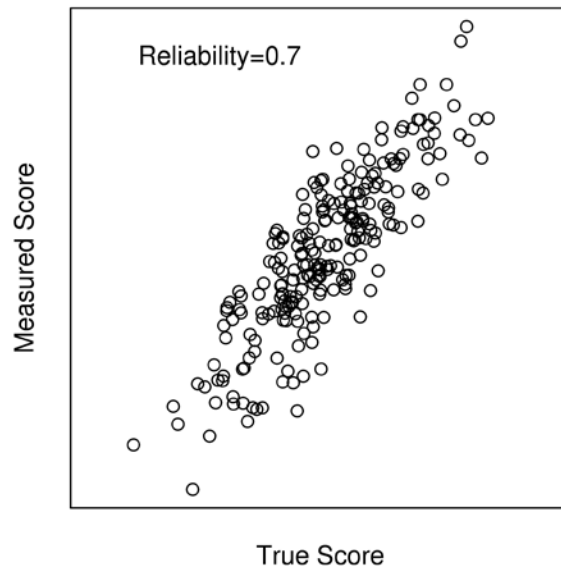
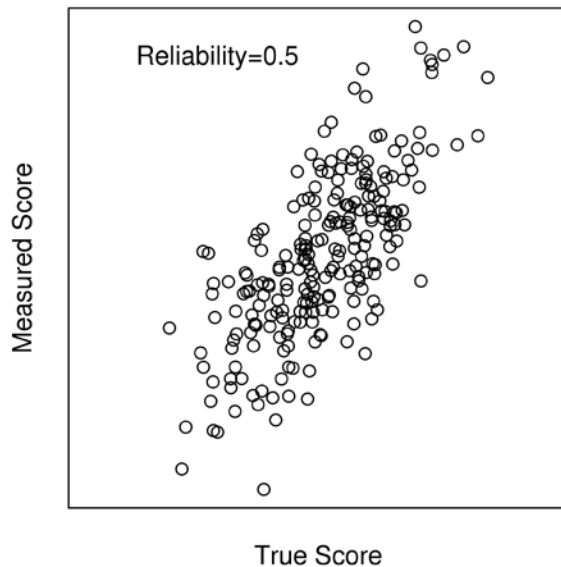
Based on all the 233 participants the reliability (proportion of signal variation) is 44.6%, 95% credible interval [CrI] (34.6%, 54.1%). Reliability increases when considering participants with a particular minimum number of cases within the time window as displayed below.

	No Minimum	≥30 cases	≥50 cases	≥100 cases	≥150 cases
No. of participants	233	186	156	101	53
Reliability	44.6%	51.7%	56.1%	60.9%	68.0%
95% CrI	(34.6%-54.1%)	(41.3%-61.4%)	(45.2%-65.6%)	(49.0%-71.2%)	(53.6%-79.7%)

2a2.4 What is your interpretation of the results in terms of demonstrating reliability? (i.e., what do the results mean and what are the norms for the test conducted?)

Reliability increases when considering participants with increasing minimum number of cases. Starting with participants with at least 30 cases, there is a moderate reliability of 0.517 (51.7%), and reliability is 0.68 (68%) when only large-volume participants (at least 150 cases) are considered. The increase in reliability is the result of a more precise estimation of a participant's measure value; in other words with the same between-participants variability, the reliability increases when the participant measurement error decreases with more cases per participant.

To visualize this effect of a decreasing measurement error on reliability, while keeping the same between-participant variability, we created two figures illustrating the accuracy of the measured scores when the true reliability is 0.50 and 0.70. Because the true score for the composite measure is unknown, we used simulated data with formula $\text{Measured Score}_i = \text{True Score}_i + e_i$ where $i = 1, 2, \dots, 233$ indicates the 233 participants and where True Score_i and e_i both follow normal distributions. The standard deviations of the normal distributions were chosen such that the measure (score) has a reliability of 0.50 on the left figure and reliability of 0.70 on the right figure. Each figure has true score along the x-axis, and the estimated (measured) value of this true score along the y-axis. With a decreasing measurement error of the score (as is the case with increase in the number of cases per participant), the correlation between the true and measured values of the score increases, and thus also, equivalently, the reliability increases because reliability can be expressed as a square of this correlation (Pearson correlation). Although a high reliability of 0.70 shows a very close correlation between true and measured scores, a more moderate reliability of 0.50 still visualizes a strong association (correlation) between the true and measured values of the score.



2b1. VALIDITY TESTING

Note: Current guidance for composite measure evaluation states that validity should be demonstrated for the composite performance measure score. If not feasible for initial endorsement, acceptable alternatives include assessment of content or face validity of the composite OR demonstration of validity for each component. Empirical validity testing of the composite measure score is expected by the time of endorsement maintenance.

2b1.1. What level of validity testing was conducted?

☐ Critical data elements (data element validity must address ALL critical data elements)

☐ Composite performance measure score

☐ Empirical validity testing

☐ Systematic assessment of face validity of performance measure score as an indicator of quality or resource use (i.e., is an accurate reflection of performance on quality or resource use and can distinguish good from poor performance) **NOTE:** Empirical validity testing is expected at time of maintenance review; if not possible, justification is required.

☒ Validity testing for component measures (check all that apply)

Note: applies to ALL component measures, unless already endorsed or are being submitted for individual endorsement.

☐ Endorsed (or submitted) as individual performance measures

☒ Critical data elements (data element validity must address ALL critical data elements)

☐ Empirical validity testing of the component measure score(s)

☐ Systematic assessment of face validity of component measure score(s) as an indicator of quality or resource use (i.e., is an accurate reflection of performance on quality or resource use and can distinguish good from poor performance)

2b1.2. For each level of testing checked above, describe the method of validity testing and what it tests (describe the steps—do not just name a method; what was tested, e.g., accuracy of data elements compared to authoritative source, relationship to another measure as expected; what statistical analysis was used)

When data arrive at the data warehouse, they are checked carefully for logical inconsistencies, missing required fields, and parent/child variable relationship violations. Any inconsistencies or violations are communicated to participants in the detailed Data Quality Report that is generated automatically following each harvest file submission. Upon receipt of

the Data Quality Report, participants are given an opportunity to correct the data, which substantially improves the quality and completeness of the data submitted for analysis. If the data inconsistencies are not changed by the participant prior to harvest close, the data warehouse performs consistency edits and/or parent/child edits on the data in order for them to be analyzable. Participants are informed of such edits to their data in the Data Quality Report.

Since 2010, the STS has contracted with Telligen (formerly IFMC) and, most recently, Cardiac Registry Support, LLC (CRS) to conduct audits of the STS General Thoracic Surgery Database on the Society's behalf to evaluate the accuracy, consistency and comprehensiveness of data collection, which has validated the integrity of the data. Currently, auditors validate case inclusion and 15 lobectomy and 5 esophagectomy cancer cases are randomly chosen for review of 39 individual data elements. The auditors abstract each designated medical record to validate data elements previously submitted to the STS data warehouse. Agreement rates are calculated for each of the 39 elements as well as for an overall agreement rate. Five sites were randomly selected for the first audit, which took place in 2010. In 2016, 25 sites were audited.

2b1.3. What were the statistical results from validity testing? (e.g., correlation; t-test)

STS audited 10% of participants in the General Thoracic Surgery Database in 2016 using an independent auditing firm (CRS). The sites were randomly selected and audited for data completeness and accuracy. Auditors compared case logs at each facility and cases submitted to the STS GTSD to assess completeness of data submission. There was consistent agreement across all participants for data completeness. Data accuracy was assessed by reabstraction of 15 randomly chosen lobectomy cancer cases and 5 esophagectomy cancer cases, comparing 39 data elements in the medical chart with the data file submitted to the STS GTSD. The agreement rate was 96.78% for overall data accuracy in 2016, with a range in agreement from 94.3% to 99.0%.

For comparison, the overall agreement rates in 2010 and 2011 were 89.9% and 94.6%, respectively (across the 33 data elements reviewed at that time). The range in agreement was from 76.5% to 95.5% in 2010, and from 88.8% to 97.5% in 2011.

Aggregate agreement rates from the 2016 audit for each of the 39 variables (data elements) and for each of the variable categories are displayed in the table below. The STS does not have access to audit results at the level of individual surgical cases; we are therefore unable to provide the kappa statistic.

CATEGORY	FIELD_NAME	NUM	DEN	Agreement Rate
PRE-OPERATIVE EVALUATION	OVERALL_ALL_FIELDS	6455	6738	95.80%
PRE-OPERATIVE EVALUATION	Admission Date	497	500	99.40%
PRE-OPERATIVE EVALUATION	Prior Cardiothoracic Surgery	488	500	97.60%
PRE-OPERATIVE EVALUATION	Pre-Op Chemo-Current Malignancy	489	500	97.80%
PRE-OPERATIVE EVALUATION	Pre-Op Thoracic Radiation Therapy	489	500	97.80%
PRE-OPERATIVE EVALUATION	Diabetes	413	423	97.64%
PRE-OPERATIVE EVALUATION	Diabetes Therapy	68	82	82.93%
PRE-OPERATIVE EVALUATION	Cigarette Smoking	489	500	97.80%
PRE-OPERATIVE EVALUATION	Pulmonary Function Tests Performed	419	423	99.05%
PRE-OPERATIVE EVALUATION	FEV1 Predicted	316	414	76.33%
PRE-OPERATIVE EVALUATION	Zubrod Score	491	500	98.20%
PRE-OPERATIVE EVALUATION	Lung Cancer	420	423	99.29%
PRE-OPERATIVE EVALUATION	Clinical Staging Method- Lung- EBUS	408	419	97.37%

CATEGORY	FIELD_NAME	NUM	DEN	Agreement Rate
PRE-OPERATIVE EVALUATION	Clinical Staging Method-Lung-PET or PET/CT	397	419	94.75%
PRE-OPERATIVE EVALUATION	Lung Cancer Tumor Size-T	377	419	89.98%
PRE-OPERATIVE EVALUATION	Lung Cancer Nodes-N	409	419	97.61%
PRE-OPERATIVE EVALUATION	Esophageal Cancer	77	77	100.00%
PRE-OPERATIVE EVALUATION	Clinical Staging Method- Esophageal-EUS	69	75	92.00%
PRE-OPERATIVE EVALUATION	Esophageal Cancer Tumor-T	68	72	94.44%
PRE-OPERATIVE EVALUATION	Clinical Diagnosis of Nodal Involvement	71	73	97.26%
DIAGNOSIS AND PROCEDURES	OVERALL_ALL FIELDS	4842	4978	97.27%
DIAGNOSIS AND PROCEDURES	Category of Disease-Primary	479	499	95.99%
DIAGNOSIS AND PROCEDURES	Date of Surgery	498	500	99.60%
DIAGNOSIS AND PROCEDURES	Procedure Start Time	493	500	98.60%
DIAGNOSIS AND PROCEDURES	Procedure End Time	482	500	96.40%
DIAGNOSIS AND PROCEDURES	ASA Classification	487	500	97.40%
DIAGNOSIS AND PROCEDURES	Procedure	500	500	100.00%
DIAGNOSIS AND PROCEDURES	Patient Disposition	491	500	98.20%
DIAGNOSIS AND PROCEDURES	Pathologic Staging-Lung Cancer-T	405	419	96.66%
DIAGNOSIS AND PROCEDURES	Pathologic Staging-Lung Cancer-N	411	419	98.09%
DIAGNOSIS AND PROCEDURES	Lung Cancer-Number of Nodes	385	419	91.89%
DIAGNOSIS AND PROCEDURES	Pathologic Staging-Esophageal Cancer-T	69	74	93.24%
DIAGNOSIS AND PROCEDURES	Pathologic Staging-Esophageal Cancer-N	73	74	98.65%
DIAGNOSIS AND PROCEDURES	Esophageal Cancer-Number of Nodes	69	74	93.24%
POST-OPERATIVE EVENTS	OVERALL_ALL FIELDS	1487	1500	99.13%
POST-OPERATIVE EVENTS	Unexpected Return to OR	493	500	98.60%
POST-OPERATIVE EVENTS	Pneumonia	494	500	98.80%
POST-OPERATIVE EVENTS	Initial Vent Support >48 Hours	500	500	100.00%
DISCHARGE	OVERALL_ALL FIELDS	1935	1993	97.09%
DISCHARGE	Discharge Date	499	500	99.80%
DISCHARGE	Discharge Status	490	500	98.00%
DISCHARGE	Readmission within 30 Days of Discharge	484	493	98.17%
DISCHARGE	Status 30 Days After Surgery	462	500	92.40%
	OVERALL_ALL FIELDS	14719	15209	96.78%

2b1.4. What is your interpretation of the results in terms of demonstrating validity? (i.e., what do the results mean and what are the norms for the test conducted?)

The most recent audits of the General Thoracic Surgery Database have demonstrated a high degree of data validity. Overall data accuracy rates have increased substantially since audits of the GTSD were first conducted in 2010; agreement ranges have also narrowed, indicating greater consistency in data accuracy among audited sites.

2b2. EXCLUSIONS ANALYSIS

Note: Applies to the composite performance measure, as well all component measures unless they are already endorsed or are being submitted for individual endorsement.

NA ☐ no exclusions — skip to section [2b4](#)

2b2.1. Describe the method of testing exclusions and what it tests (describe the steps—do not just name a method; what was tested, e.g., whether exclusions affect overall performance scores; what statistical analysis was used)

We excluded patients with missing data for age, sex, or discharge mortality status. In addition we excluded patients with non-elective status, occult or stage 0 tumors, or American Society of Anesthesiologists class VI. We believe these are clinically appropriate exclusions and are necessary to make the measure a consistent performance measure for the comparison across participants. The exclusions are precisely defined and specified.

2b2.2. What were the statistical results from testing exclusions? (include overall number and percentage of individuals excluded, frequency distribution of exclusions across measured entities, and impact on performance measure scores)

There were 183 (0.7%) occult or stage 0 tumors, 8 (0.03%) ASA VI, and 337 (1.3%) non-elective status patients, resulting in the overall exclusion of 2.1% (528 of 24,912 patient records). Impact of these exclusions on the performance measure is negligible due to the small proportion of cases excluded.

2b2.3. What is your interpretation of the results in terms of demonstrating that exclusions are needed to prevent unfair distortion of performance results? (i.e., the value outweighs the burden of increased data collection and analysis.

Note: If patient preference is an exclusion, the measure must be specified so that the effect on the performance score is transparent, e.g., scores with and without exclusion)

For the measure to consistently quantify the surgical quality of lobectomy for lung cancer per its definition (outcome domains of operative mortality and major complications), it is necessary and clinically appropriate to exclude cases with non-elective status, occult or stage 0 tumors, or American Society of Anesthesiologists class VI.

2b3. RISK ADJUSTMENT/STRATIFICATION FOR OUTCOME OR RESOURCE USE MEASURES

Note: Applies to all outcome or resource use component measures, unless already endorsed or are being submitted for individual endorsement.

If not an intermediate or health outcome, or PRO-PM, or resource use measure, skip to section [2b4](#).

2b3.1. What method of controlling for differences in case mix is used? (check all that apply)

☐ Endorsed (or submitted) as individual performance measures

☐ No risk adjustment or stratification

☒ Statistical risk model with risk factors

☐ Stratification by risk categories

☐ Other,

2b3.1.1 If using statistical risk models, provide detailed risk model specifications, including the risk model method, risk factors, coefficients, equations, codes with descriptors, and definitions.

Participant-specific risk-adjusted operative mortality and major complication rates were estimated using a bivariate random-effects logistic regression model. The term bivariate refers to the fact that both operative mortality and major complications were analyzed together in a single model,

not estimated one at a time in separate models. Random-effects refers to the assumption that the provider-specific parameters of interest are assumed to arise from a specified distribution defined by parameters that are also estimated

in the modelling process. Detailed description is provided in published statistical appendix; a copy is appended to the end of this document. Risk factors in the model were: age, sex, year of operation, body mass index, hypertension, steroid therapy, congestive heart failure, coronary artery disease, peripheral vascular disease, reoperation, preoperative chemotherapy within 6 months, cerebrovascular disease, diabetes mellitus, renal failure, dialysis, past smoker, current smoker, forced expiratory volume in 1 second percent of predicted, Zubrod score (linear plus quadratic), American Society of Anesthesiologists class (linear plus quadratic), and pathologic stage.

2b3.2. If an outcome or resource use component measure is not risk adjusted or stratified, provide rationale and analyses to demonstrate that controlling for differences in patient characteristics (case mix) is not needed to achieve fair comparisons across measured entities.

n/a

2b3.3a. Describe the conceptual/clinical and statistical methods and criteria used to select patient factors (clinical factors or social risk factors) used in the statistical risk model or for stratification by risk (e.g., potential factors identified in the literature and/or expert panel; regression analysis; statistical significance of $p < 0.10$; correlation of x or higher; patient factors should be present at the start of care) Also discuss any “ordering” of risk factor inclusion; for example, are social risk factors added after all clinical factors?

Covariates in this model were selected a priori based on a combination of literature review and expert group consensus, and as described in Kozower, et al. (2016). All covariates were retained in the model and were not added or removed based on a statistical variable selection algorithm.

No social risk factors were used in the statistical risk model or for stratification.

Kozower BD, O’Brien SM, Kosinski AS, et al. The Society of Thoracic Surgeons composite score for rating program performance for lobectomy for lung cancer. *Ann Thorac Surg* 2016;101:1379-87.

2b3.3b. How was the conceptual model of how social risk impacts this outcome developed? Please check all that apply:

- ☒ Published literature
- ☐ Internal data analysis
- ☒ Other (please describe)

Expert group consensus

2b3.4a. What were the statistical results of the analyses used to select risk factors?

Estimated odds ratios are summarized in the table below.

Variable	Operative Mortality		Major Morbidity	
	OR (95% CI)	p-value	OR (95% CI)	p-value
Age, yrs, (per 1 yr increase)	1.043 (1.027, 1.058)	<.0001	1.011 (1.006, 1.017)	<.0001
Male	1.372 (1.081, 1.743)	0.0094	1.377 (1.252, 1.514)	<.0001
Body Mass Index (kg/m ²), (per 1 unit increase)	0.958 (0.937, 0.98)	0.0002	0.986 (0.978, 0.994)	0.0007
Hypertension	1.471 (1.106, 1.955)	0.0079	0.986 (0.889, 1.095)	0.7936
Steroid therapy	1.419 (0.844, 2.387)	0.1866	1.027 (0.797, 1.322)	0.839
Congestive heart failure	1.611 (1.004, 2.585)	0.0483	1.202 (0.942, 1.535)	0.1395
Coronary artery disease	1.308 (1.007, 1.698)	0.0443	1.286 (1.150, 1.438)	<.0001
Peripheral vascular disease	1.738 (1.298, 2.328)	0.0002	1.248 (1.085, 1.435)	0.0019
Reoperation	1.328 (0.894, 1.975)	0.1604	1.110 (0.926, 1.331)	0.2583
Preoperative chemotherapy within 6 months	1.229 (0.791, 1.911)	0.3592	1.268 (1.065, 1.509)	0.0075
Cerebrovascular disease	1.062 (0.744, 1.514)	0.7409	1.116 (0.955, 1.304)	0.1674

	Operative Mortality		Major Morbidity	
Variable	OR (95% CI)	p-value	OR (95% CI)	p-value
Diabetes mellitus	1.026 (0.775, 1.358)	0.8591	0.968 (0.858, 1.091)	0.5888
Renal failure	1.695 (0.873, 3.29)	0.119	1.387 (0.986, 1.95)	0.0604
Dialysis	4.110 (1.761, 9.596)	0.0011	1.005 (0.535, 1.888)	0.9885
Past smoker	1.172 (0.774, 1.776)	0.4533	1.522 (1.272, 1.821)	<.0001
Current smoker	1.411 (0.889, 2.238)	0.1441	2.168 (1.790, 2.627)	<.0001
FEV in 1 second percent of predicted (per 1 unit increase)	0.991 (0.985, 0.997)	0.0028	0.987 (0.985, 0.99)	<.0001
Zubrod score (per 1 unit increase)	1.233 (0.895, 1.699)	0.1997	1.182 (1.030, 1.355)	0.0172
Squared Zubrod score (per 1 unit increase)	1.039 (0.922, 1.17)	0.5295	1.021 (0.962, 1.083)	0.5003
ASA Class (per 1 unit increase)	2.160 (0.383, 12.181)	0.3828	1.127 (0.595, 2.137)	0.7139
Squared ASA Class (per 1 unit increase)	0.909 (0.691, 1.196)	0.4952	1.032 (0.931, 1.144)	0.5532
Pathologic stage I	1.216 (0.910, 1.626)	0.1867	1.200 (1.068, 1.349)	0.0022
Pathologic stage II	1.660 (1.199, 2.298)	0.0022	1.142 (0.984, 1.325)	0.0797
Pathologic stage IV	1.575 (0.686, 3.615)	0.2841	1.222 (0.862, 1.733)	0.2593
Year of operation (per 1 yr increase)	0.916 (0.797, 1.053)	0.2188	0.925 (0.874, 0.978)	0.0065

2b3.4b. Describe the analyses and interpretation resulting in the decision to select social risk factors (e.g. prevalence of the factor across measured entities, empirical association with the outcome, contribution of unique variation in the outcome, assessment of between-unit effects and within-unit effects.) Also describe the impact of adjusting for social risk (or not) on providers at high or low extremes of risk.

All covariates were retained in the model and were not added or removed based on a statistical variable selection algorithm.

As noted in 1.8 above, patient social risk data are not collected in the General Thoracic Surgery Database.

2b3.5. Describe the method of testing/analysis used to develop and validate the adequacy of the statistical model or stratification approach (describe the steps—do not just name a method; what statistical analysis was used)

Continuous variables were evaluated with respect to linearity of effect and needed transformations were considered resulting in addition of squared ASA class and Zubrod score. The calibration of the model was assessed with the Hosmer-Lemeshow statistic. The discrimination of the model was assessed with the C-statistic.

Provide the statistical results from testing the approach to controlling for differences in patient characteristics (case mix) below.

If stratified, skip to 2b3.9

2b3.6. Statistical Risk Model Discrimination Statistics (e.g., c-statistic, R-squared):

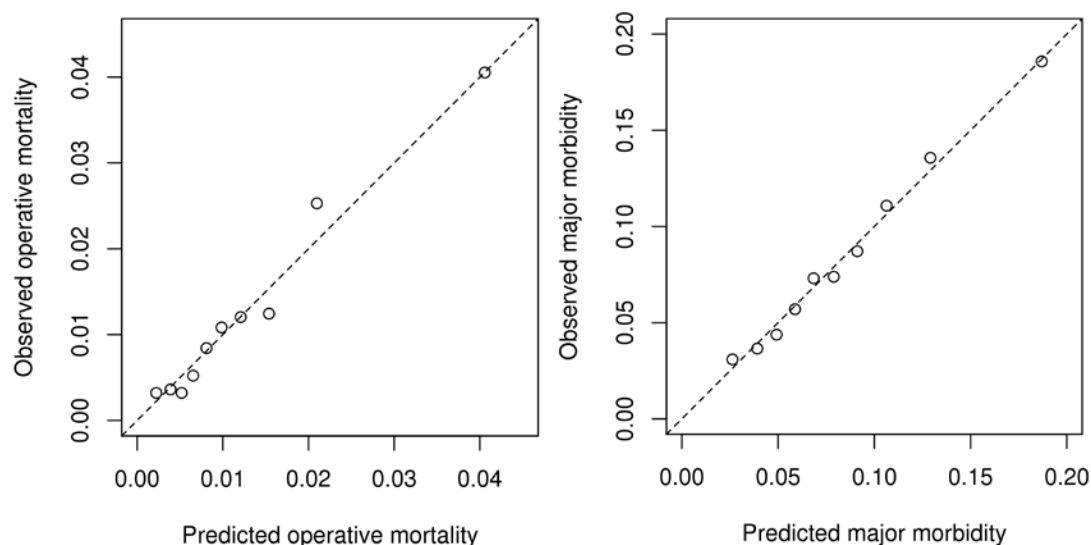
Operative mortality model: C-statistic is 0.731. Major morbidity model: C-statistic is 0.667.

2b3.7. Statistical Risk Model Calibration Statistics (e.g., Hosmer-Lemeshow statistic):

Operative mortality model: Hosmer and Lemeshow Goodness-of-Fit Test p-value=0.47 (Chi-Square=7.65, df=8). Major morbidity model: Hosmer and Lemeshow Goodness-of-Fit Test p-value=0.44 (Chi-Square=7.95, df=8).

2b3.8. Statistical Risk Model Calibration – Risk decile plots or calibration curves:

Risk decile plots below show good alignment of predicted and observed probabilities of outcome (operative mortality and major morbidity) within deciles of predicted values.



2b3.9. Results of Risk Stratification Analysis:

n/a

2b3.10. What is your interpretation of the results in terms of demonstrating adequacy of controlling for differences in patient characteristics (case mix)? (i.e., what do the results mean and what are the norms for the test conducted)

The results demonstrated that the STS lobectomy risk models are well calibrated and have good discrimination power. They are suitable for controlling for differences in case-mix between centers.

2b3.11. Optional Additional Testing for Risk Adjustment (*not required*, but would provide additional support of adequacy of risk model, e.g., testing of risk model in another data set; sensitivity analysis for missing data; other methods that were assessed)

n/a

2b4. IDENTIFICATION OF STATISTICALLY SIGNIFICANT & MEANINGFUL DIFFERENCES IN PERFORMANCE

Note: Applies to the composite performance measure.

2b4.1. Describe the method for determining if statistically significant and clinically/practically meaningful differences in performance measure scores among the measured entities can be identified (describe the steps—do not just name a method; what statistical analysis was used? Do not just repeat the information provided related to performance gap in 1b)

The degree of uncertainty surrounding an STS participant's composite measure estimate is indicated by calculating 95% Bayesian credible intervals (CI's) which are similar to conventional confidence intervals. Point estimates and CI's for an individual STS participant are reported along with a comparison to the overall average STS composite score. In addition, the composite measure result is converted into categories labeled as 1 to 3 stars. An STS participant receives 2 stars if the Bayesian credible interval surrounding their composite score overlaps the overall STS average. This rating implies that the STS participant's performance was not statistically different from the overall STS national average. If the Bayesian CI falls entirely above the STS national average, the participant receives 3 stars (higher-than-expected performance). If the Bayesian CI falls entirely below the STS national average, the participant receives 1 star (lower-than-expected performance).

2b4.2. What were the statistical results from testing the ability to identify statistically significant and/or clinically/practically meaningful differences in performance measure scores across measured entities? (e.g., number and percentage of entities with scores that were statistically significantly different from mean or some benchmark, different from expected; how was meaningful difference defined)

Among participants with at least 30 cases over 3 years, 93.1% of participants have received 2 stars, and the remaining participants have received either 1 or 3 stars.

January 1, 2014 through December 31, 2016

	All Participants	Participants N ≥ 30
Category	Number of Participants, %	Number of Participants, %
1-star	6, 2.6%	6, 3.2%
2-star	217, 93.1%	170, 91.4%
3-star	10, 4.3%	10, 5.4%

2b4.3. What is your interpretation of the results in terms of demonstrating the ability to identify statistically significant and/or clinically/practically meaningful differences in performance across measured entities? (i.e., what do the results mean in terms of statistical and meaningful differences?)

The Bayesian methodology allows direct probability interpretation of the results. The identified differences in performance are both statistically significant and clinically meaningful. The surgeon panel and users are satisfied with the distribution of participants across performance categories.

2b5. COMPARABILITY OF PERFORMANCE SCORES WHEN MORE THAN ONE SET OF SPECIFICATIONS

Note: Applies to all component measures, unless already endorsed or are being submitted for individual endorsement.

If only one set of specifications, this section can be skipped.

Note: This item is directed to measures that are risk-adjusted (with or without social risk factors) **OR** to measures with more than one set of specifications/instructions (e.g., one set of specifications for how to identify and compute the measure from medical record abstraction and a different set of specifications for claims or eMeasures). It does not apply to measures that use more than one source of data in one set of specifications/instructions (e.g., claims data to identify the denominator and medical record abstraction for the numerator). **Comparability is not required when comparing performance scores with and without social risk factors in the risk adjustment model. However, if comparability is not demonstrated for measures with more than one set of specifications/instructions, the different specifications (e.g., for medical records vs. claims) should be submitted as separate measures.**

2b5.1. Describe the method of testing conducted to compare performance scores for the same entities across the different data sources/specifications (describe the steps—do not just name a method; what statistical analysis was used)

n/a

2b5.2. What were the statistical results from testing comparability of performance scores for the same entities when using different data sources/specifications? (e.g., correlation, rank order)

n/a

2b5.3. What is your interpretation of the results in terms of the differences in performance measure scores for the same entities across the different data sources/specifications? (i.e., what do the results mean and what are the norms for the test conducted?)

n/a

2b6. MISSING DATA ANALYSIS AND MINIMIZING BIAS

Note: Applies to the overall composite measure.

2b6.1. Describe the method of testing conducted to identify the extent and distribution of missing data (or nonresponse) and demonstrate that performance results are not biased due to systematic missing data (or differences between responders and nonresponders) and how the specified handling of missing data minimizes bias (*describe the steps—do not just name a method; what statistical analysis was used*)

The quality of data in the STS General Thoracic Surgery Database has been improving. We managed the missing data with imputation. Missing body mass index (BMI) values (1%) were imputed utilizing the median of the observed BMI values. Missing FEV1 (3.4%) was imputed to the median within the smoking status categories. Missing pathologic stage (3.1%) was imputed to its mode (stage I). For binary risk factors, missing values were considered as indicating absence of the risk factor.

2b6.2. What is the overall frequency of missing data, the distribution of missing data across providers, and the results from testing related to missing data? (*e.g., results of sensitivity analysis of the effect of various rules for missing data/nonresponse; if no empirical sensitivity analysis, identify the approaches for handling missing data that were considered and pros and cons of each*)

To maximize use of available data, when encountering records with missing values of model covariates (with the exception of age and gender), the missing values were imputed. Patient records missing age or gender were excluded. Variables FEV1, steroid use, dialysis, and pathologic stage were each missing for approximately 3% of patients. Remaining variables had less than 1% of missing values.

2b6.3. What is your interpretation of the results in terms of demonstrating that performance results are not biased due to systematic missing data (or differences between responders and nonresponders) and how the specified handling of missing data minimizes bias? (*i.e., what do the results mean in terms of supporting the selected approach for missing data and what are the norms for the test conducted; if no empirical analysis, provide rationale for the selected approach for missing data*)

The rates of missing data were low and are getting lower. We therefore concluded that systematic missing data did not lead to bias in our measure.

2c. EMPIRICAL ANALYSIS TO SUPPORT COMPOSITE CONSTRUCTION APPROACH

Note: *If empirical analyses do not provide adequate results—or are not conducted—justification must be provided and accepted in order to meet the must-pass criterion of Scientific Acceptability of Measure Properties. Each of the following questions has instructions if there is no empirical analysis.*

2d1. Empirical analysis demonstrating that the component measures fit the quality construct, add value to the overall composite, and achieve the object of parsimony to the extent possible.

2d1.1 Describe the method used (*describe the steps—do not just name a method; what statistical analysis was used; if no empirical analysis, provide justification*)

To verify that each domain contributes statistical information, we calculated the operative mortality and major complication rates across program star ratings among 186 hospitals with at least 30 lobectomies within three years.

2d1.2. What were the statistical results obtained from the analysis of the components? (*e.g., correlations, contribution of each component to the composite score, etc.; if no empirical analysis, identify the components that were considered and the pros and cons of each*)

The table below demonstrates that the mortality and major complication rates decrease monotonically from one-star (below average) to three-star (above average) participants.

Operative Mortality and Major Complication Rates Across Star Ratings

	One star	Two Star	Three Star	All Programs
Operative mortality (95% CI)	2.1% (1.4%, 3.2%)	1.3% (1.1%, 1.4%)	0.4% (0.2%, 0.7%)	1.2% (1.1%, 1.4%)
Major complication (95% CI)	16.2% (14.1%, 18.6%)	8.4% (8.0%, 8.8%)	3.2% (2.5%, 4.1%)	8.3% (8.0%, 8.7%)

Among 186 hospitals with at least 30 lobectomies.

2d1.3. What is your interpretation of the results in terms of demonstrating that the components included in the composite are consistent with the described quality construct and add value to the overall composite? (i.e., what do the results mean in terms of supporting inclusion of the components; if no empirical analysis, provide rationale for the components that were selected)

Although risk-adjusted morbidity explains more of the variation in the overall composite score, it does not dominate. Both domains contribute statistical information.

2d2. Empirical analysis demonstrating that the aggregations and weighting rules are consistent with the quality construct and achieve the objective of simplicity to the extent possible

2d2.1 Describe the method used (describe the steps—do not just name a method; what statistical analysis was used; if no empirical analysis, provide justification)

To form the composite, we rescaled the morbidity and mortality domains by dividing by their respective standard deviations across STS participants and then added the two domains together. This weighting was then assessed by an expert panel to determine if it provided an appropriate reflection of the relative importance of the two domains.

2d2.2. What were the statistical results obtained from the analysis of the aggregation and weighting rules? (e.g., results of sensitivity analysis of effect of different aggregations and/or weighting rules; if no empirical analysis, identify the aggregation and weighting rules that were considered and the pros and cons of each)

After rescaling, the relative weights in the final composite of risk-standardized mortality and risk-standardized major morbidity were 0.827 and 0.173, respectively. An implication of this weighting is that a 1 percentage point change in a participant's risk-adjusted mortality rate has the same impact as a 4.8 percentage point change in the site's risk-adjusted morbidity rate.

2d2.3. What is your interpretation of the results in terms of demonstrating the aggregation and weighting rules are consistent with the described quality construct? (i.e., what do the results mean in terms of supporting the selected rules for aggregation and weighting; if no empirical analysis, provide rationale for the selected rules for aggregation and weighting)

This weighting was consistent with our expert panel's clinical assessment of each domain's relative importance.

Statistical Model

For the i -th of n_j patients at the j -th participant ($j=1, 2, \dots, N$), let Y_{1ji} be a binary indicator of operative mortality status (0=alive, 1=dead), let Y_{2ji} be an indicator of major complications (0=none, 1=at least one), and let $x_{ji} = (x_{1ji}, x_{2ji}, \dots, x_{qji})$ be a set of numerically encoded patient baseline characteristics (e.g. age in years; binary risk factors coded as 0=absent, 1=present, etc.). Let $\pi_{kji} = \Pr(Y_{kji} = 1 | x_{ji})$ denote the probability of the occurrence of the k -th endpoint where $k=1$ refers to mortality and $k=2$ refers to complications. The associations of x_{ji} with Y_{1ji} and Y_{2ji} are assumed to be described by a bivariate random effects logistic regression model with normally distributed hospital-specific random intercept parameters. In particular, we assume:

$$\begin{aligned} \text{(operative mortality)} \quad & \log \left(\frac{\pi_{1ji}}{1-\pi_{1ji}} \right) = \alpha_{1j} + x'_{ji}\beta_1 \\ \text{(major complication)} \quad & \log \left(\frac{\pi_{2ji}}{1-\pi_{2ji}} \right) = \alpha_{2j} + x'_{ji}\beta_2 \\ \text{(random effects)} \quad & (\alpha_{1j}, \alpha_{2j}) \stackrel{\text{iid}}{\sim} N(\mu, \Sigma), \end{aligned}$$

where $\beta_1 = (\beta_{11}, \beta_{12}, \dots, \beta_{1q})$ denotes a set of unknown regression coefficients relating covariates to mortality, $\beta_2 = (\beta_{21}, \beta_{22}, \dots, \beta_{2q})$ denotes a set of unknown regression coefficients relating covariates to major complication, $(\alpha_{1j}, \alpha_{2j})$ denote a set of normally distributed hospital-specific random effect parameters, and $N(\mu, \Sigma)$ denotes a bivariate normal distribution with mean $\mu = (\mu_1, \mu_2)$ and covariance $\Sigma = (\sigma_{11}, \sigma_{12}, \sigma_{21}, \sigma_{22})$. Conditional on π_{1ji} and π_{2ji} , the variables Y_{1ji} and Y_{2ji} are assumed to be distributed as two independent Bernoulli variables with parameters π_{1ji} and π_{2ji} , respectively. That is:

$$\Pr(Y_{1ji} = y_{1ji}, Y_{2ji} = y_{2ji} | \pi_{1ji}, \pi_{2ji}) = \prod_{k=1}^2 \pi_{kji}^{y_{kji}} (1 - \pi_{kji})^{1-y_{kji}}.$$

Outcomes of patients at different participants are assumed to be statistically independent, and outcomes of patients at the same participant are assumed to be conditionally independent given $(\alpha_{1j}, \alpha_{2j})$. The assumption that Y_{1ji} and Y_{2ji} are conditionally independent given π_{1ji} and π_{2ji} is likely to be violated in practice but is made in order to facilitate computation. Although the model assumes *conditional* independence between Y_{1ji} and Y_{2ji} , the model does not assume *marginal* independence between these two variables, as the underlying probabilities π_{1ji} and π_{2ji} depend on random effects parameters which account for within-hospital correlation.

Definition of Risk-Adjusted Rates

Based on this model, the j -th participant's risk-adjusted rates of operative mortality and major complications were defined as

$$\begin{aligned} \text{(operative mortality)} \quad & \theta_{1j} = \frac{\sum_{i=1}^{n_j} \text{expit}(\alpha_{1j} + x'_{ji}\beta_1)}{\sum_{i=1}^{n_j} \text{expit}(\mu_1 + x'_{ji}\beta_1)} \times \bar{Y}_1 \\ \text{(major complication)} \quad & \theta_{2j} = \frac{\sum_{i=1}^{n_j} \text{expit}(\alpha_{2j} + x'_{ji}\beta_2)}{\sum_{i=1}^{n_j} \text{expit}(\mu_2 + x'_{ji}\beta_2)} \times \bar{Y}_2 \end{aligned}$$

where \bar{Y}_1 denotes the overall aggregate observed rate of operative mortality in the study sample and \bar{Y}_2 denotes the overall aggregate observed rate of major complication in the study sample.

Definition of Composite Score

The overall composite of the j -th participant was defined as

$$\theta_j = \omega(1 - \theta_{1j}) + (1 - \omega)(1 - \theta_{2j})$$

Where $\omega = (1/\sigma_1)/(1/\sigma_1 + 1/\sigma_2)$ and σ_k denotes the standard deviation of θ_{kj} 's across participants, $k=1, 2$.

Estimation

Model parameters were estimated in a Bayesian framework by specifying a prior probability distribution for the unknown model parameters β_1, β_2, μ , and Σ . Because our prior knowledge was limited, we specified a vague proper prior distribution that consisted of independent normal distributions for the elements of β_1, β_2 , and μ , and an inverse Wishart distribution for Σ . Posterior means and credible intervals were calculated using Markov Chain Monte Carlo (MCMC) simulations as implemented in OpenBUGS version 3.2.2 software. Posterior summaries were calculated by generating 50,000 sets of simulated parameter values after a long burn-in period to ensure convergence and then thinning the sample to arrive at a final set of 5,000 iterations. The parameter θ_j was estimated as $\hat{\theta}_j = \sum_{l=1}^{5000} \theta_j^{(l)} / 5000$, where $\theta_j^{(l)}$ denotes the simulated values of θ_j at the l -th iteration of the MCMC procedure. A 95% Bayesian credible interval was obtained by calculating the 125th lowest and 125th highest values of θ_j across the 5000 simulated values.

2. Reliability and Validity—Scientific Acceptability of Measure Properties

Extent to which the measure, as specified, produces consistent (reliable) and credible (valid) results about the quality of care when implemented. **Measures must be judged to meet the sub criteria for both reliability and validity to pass this criterion and be evaluated against the remaining criteria.**

2a.1. Specifications The measure is well defined and precisely specified so it can be implemented consistently within and across organizations and allows for comparability. eMeasures should be specified in the Health Quality Measures Format (HQMF) and the Quality Data Model (QDM).

De.5. Subject/Topic Area (check all the areas that apply):

De.6. Non-Condition Specific(check all the areas that apply):

De.7. Target Population Category (Check all the populations for which the measure is specified and tested if any):

S.1. Measure-specific Web Page (Provide a URL link to a web page specific for this measure that contains current detailed specifications including code lists, risk model details, and supplemental materials. Do not enter a URL linking to a home page or to general information.)

http://www.sts.org/sites/default/files/documents/STSThoracicDataSpecsV2_3.pdf

S.2a. If this is an eMeasure, HQMF specifications must be attached. Attach the zipped output from the eMeasure authoring tool (MAT) - if the MAT was not used, contact staff. (Use the specification fields in this online form for the plain-language description of the specifications)

This is not an eMeasure Attachment:

S.2b. Data Dictionary, Code Table, or Value Sets (and risk model codes and coefficients when applicable) must be attached. (Excel or csv file in the suggested format preferred - if not, contact staff)

Attachment **Attachment:** STSThoracicDataSpecsV2_3-636463839166691726.pdf

S.2c. Is this an instrument-based measure (i.e., data collected via instruments, surveys, tools, questionnaires, scales, etc.)? Attach copy of instrument if available.

No, this is not an instrument-based measure Attachment:

S.2d. Is this an instrument-based measure (i.e., data collected via instruments, surveys, tools, questionnaires, scales, etc.)? Attach copy of instrument if available.

Not an instrument-based measure

S.3.1. For maintenance of endorsement: Are there changes to the specifications since the last updates/submission. If yes, update the specifications for S1-2 and S4-22 and explain reasons for the changes in S3.2.

No

S.3.2. For maintenance of endorsement, please briefly describe any important changes to the measure specifications since last measure update and explain the reasons.

n/a

S.4. Numerator Statement (Brief, narrative description of the measure focus or what is being measured about the target population, i.e., cases from the target population with the target process, condition, event, or outcome) **DO NOT** include the rationale for the measure.

IF an OUTCOME MEASURE, state the outcome being measured. Calculation of the risk-adjusted outcome should be described in the calculation algorithm (S.14).

The STS Lobectomy Composite Score comprises two domains:

1. Operative Mortality (death during the same hospitalization as surgery or within 30 days of the procedure)

2. Presence of at least one of these major complications: pneumonia, acute respiratory distress syndrome, bronchopleural fistula, pulmonary embolus, initial ventilator support greater than 48 hours, reintubation/respiratory failure, tracheostomy, myocardial infarction, or unexpected return to the operating room.

The composite score is created by a weighted combination of the above two domains resulting in a single composite score. Operative mortality and major complications were weighted inversely by their respective standard deviations across participants. This procedure is equivalent to first rescaling mortality and complications by their respective standard deviations and then assigning equal weighting to the rescaled mortality rate and rescaled complication rate. This is the same methodology used for other STS composite measures.

In addition to receiving a numeric score, participants are assigned to rating categories designated by the following:

1 star: lower-than expected performance

2 stars: as-expected-performance

3 star: higher-than-expected-performance

Patient Population: The STS GTSD was queried for all patients treated with lobectomy for lung cancer between January 1, 2014, and December 31, 2016. We excluded patients with non-elective status, occult or stage 0 tumors, American Society of Anesthesiologists class VI, and with missing data for age, sex, or discharge mortality status.

Time Window: 01/01/2014 - 12/31/2016

Model variables: Variables in the model: age, sex, year of operation, body mass index, hypertension, steroid therapy, congestive heart failure, coronary artery disease, peripheral vascular disease, reoperation, preoperative chemotherapy within 6 months, cerebrovascular disease, diabetes mellitus, renal failure, dialysis, past smoker, current smoker, forced expiratory volume in 1 second percent of predicted, Zubrod score (linear plus quadratic), American Society of Anesthesiologists class (linear plus quadratic), and pathologic stage.

S.5. Numerator Details *(All information required to identify and calculate the cases from the target population with the target process, condition, event, or outcome such as definitions, time period for data collection, specific data collection items/responses, code/value sets – Note: lists of individual codes with descriptors that exceed 1 page should be provided in an Excel or csv file in required format at S.2b)*

IF an OUTCOME MEASURE, describe how the observed outcome is identified/counted. Calculation of the risk-adjusted outcome should be described in the calculation algorithm (S.14).

Number of patients undergoing elective lobectomy for lung cancer for whom:

1. Postoperative events (POEvents - STS GTS Database, v 2.2, sequence number 1710) is marked “Yes” and one of the following items is marked:
 - a. Reintubation (Reintube - STS GTS Database, v 2.2, sequence number 1850)
 - b. Need for tracheostomy (Trach - STS GTS Database, v 2.2, sequence number 1860)
 - c. Initial ventilator support > 48 hours (Vent- STS GTS Database, v 2.2, sequence number 1840)
 - d. Acute Respiratory Distress Syndrome (ARDS - STS GTS Database, v 2.2, sequence number 1790)
 - e. Pneumonia (Pneumonia - STS GTS Database, v 2.2, sequence number 1780)
 - f. Pulmonary Embolus (PE - STS GTS Database, v 2.2, sequence number 1820)
 - g. Bronchopleural Fistula (Bronchopleural - STS GTS Database, v 2.2, sequence number 1810)
 - h. Myocardial infarction (MI - STS GTS Database, v 2.2, sequence number 1900)

Or

2. Unexpected return to the operating room (ReturnOR - STS GTS Database, Version 2.2, sequence number 1720) is marked “yes”

Or

3. One of the following fields is marked “dead”
 - a. Discharge status (MtDCStat - STS GTS Database, Version 2.2, sequence number 2200);
 - b. Status at 30 days after surgery (Mt30Stat - STS GTS Database, Version 2.2, sequence number 2240)

Please see STS General Thoracic Surgery Database Data Collection Form, Version 2.3-

http://www.sts.org/sites/default/files/documents/STSThoracicDCF_V2_3_MajorProc_Annotated.pdf

S.6. Denominator Statement *(Brief, narrative description of the target population being measured)*

Number of patients greater than or equal to 18 years of age undergoing elective lobectomy for lung cancer

S.7. Denominator Details *(All information required to identify and calculate the target population/denominator such as definitions, time period for data collection, specific data collection items/responses, code/value sets – Note: lists of individual codes with descriptors that exceed 1 page should be provided in an Excel or csv file in required format at S.2b.)*

IF an OUTCOME MEASURE, describe how the target population is identified. Calculation of the risk-adjusted outcome should be described in the calculation algorithm (S.14).

1. Lung cancer (LungCancer - STS GTS Database, v 2.2, sequence number 830) is marked “yes” and Category of Disease – Primary (CategoryPrim - STS GTS Database, v 2.2, sequence number 1300) is marked as one of the following: (ICD-9, ICD-10)

Lung cancer, main bronchus, carina (162.2, C34.00)

Lung cancer, upper lobe (162.3, C34.10)

Lung cancer, middle lobe (162.4, C34.2)

Lung cancer, lower lobe (162.5, C34.30)

Lung cancer, location unspecified (162.9, C34.90)

2. Patient has lung cancer (as defined in #1 above) and primary procedure is one of the following CPT codes:

Thoracoscopy, surgical; with lobectomy (32663)

Removal of lung, single lobe (lobectomy) (32480)

3. Status of Operation (Status - STS General Thoracic Surgery Database, Version 2.2, sequence number 1420) is marked as “Elective”

4. Only analyze the first operation of the hospitalization meeting criteria 1-3

S.8. Denominator Exclusions *(Brief narrative description of exclusions from the target population)*

Patients were excluded with non-elective status, occult or stage 0 tumors, American Society of Anesthesiologists class VI, and with missing data for age, sex, or discharge mortality status.

S.9. Denominator Exclusion Details *(All information required to identify and calculate exclusions from the denominator such as definitions, time period for data collection, specific data collection items/responses, code/value sets – Note: lists of individual codes with descriptors that exceed 1 page should be provided in an Excel or csv file in required format at S.2b.)*

Cases removed from calculations if Emergent, Urgent, or Palliative is checked under "Status of Operation"

OR if T0 is checked under Pathological Staging of the Lung / Lung Tumor: PathStageLungT(1540)

OR if VI is checked under ASA Classification: ASA (1470)

Only general thoracic procedures coded as primary lung or primary esophageal cancer are included in measure calculations, so occult carcinoma is effectively excluded.

S.10. Stratification Information *(Provide all information required to stratify the measure results, if necessary, including the stratification variables, definitions, specific data collection items/responses, code/value sets, and the risk-model covariates and coefficients for the clinically-adjusted version of the measure when appropriate – Note: lists of individual codes with descriptors that exceed 1 page should be provided in an Excel or csv file in required format with at S.2b.)*

n/a

S.11. Risk Adjustment Type (Select type. Provide specifications for risk stratification in measure testing attachment)

Statistical risk model

If other:

S.12. Type of score:

Rate/proportion

If other:

S.13. Interpretation of Score (*Classifies interpretation of score according to whether better quality is associated with a higher score, a lower score, a score falling within a defined interval, or a passing score*)

Better quality = Lower score

S.14. Calculation Algorithm/Measure Logic (*Diagram or describe the calculation of the measure score as an ordered sequence of steps including identifying the target population; exclusions; cases meeting the target process, condition, event, or outcome; time period for data, aggregating data; risk adjustment; etc.*)

Target population is patients treated with lobectomy for lung cancer. Patients were excluded with non-elective status, occult or stage 0 tumors, American Society of Anesthesiologists class VI, and with missing data for age, sex, or discharge mortality status. Outcomes were measured in two domains:

1. Operative Mortality (death during the same hospitalization as surgery or within 30 days of the procedure)
2. Presence of at least one of these major complications: pneumonia, acute respiratory distress syndrome, bronchopleural fistula, pulmonary embolus, initial ventilator support greater than 48 hours, reintubation/respiratory failure, tracheostomy, myocardial infarction, or unexpected return to the operating room.

Time window for analysis was between 01/01/2014 and 12/31/2016.

Analysis considered 24,912 patient records across 233 participant sites.

To form the composite, we rescaled the major complication and operative mortality domains by dividing by their respective standard deviations across STS participants and then added the two domains together. This weighting was then assessed by an expert panel to determine if it provided an appropriate reflection of the relative importance of the two domains.

After rescaling, the relative weights in the final composite of risk-standardized mortality and risk-standardized major morbidity were 0.827 and 0.173, respectively. An implication of this weighting is that a 1 percentage point change in a participant's risk-adjusted mortality rate has the same impact as a 4.8 percentage point change in the site's risk-adjusted morbidity rate. Our expert panel concurred that this weighting was consistent with their clinical assessment of each domain's relative importance.

S.15. Sampling (*If measure is based on a sample, provide instructions for obtaining the sample and guidance on minimum sample size.*)

IF an instrument-based performance measure (e.g., PRO-PM), identify whether (and how) proxy responses are allowed.

n/a

S.16. Survey/Patient-reported data (*If measure is based on a survey or instrument, provide instructions for data collection and guidance on minimum response rate.*)

Specify calculation of response rates to be reported with performance measure results.

n/a

S.17. Data Source (*Check ONLY the sources for which the measure is SPECIFIED AND TESTED*).

If other, please describe in S.18.

Other, Registry Data

S.18. Data Source or Collection Instrument (*Identify the specific data source/data collection instrument (e.g. name of database, clinical registry, collection instrument, etc., and describe how data are collected.)*)

IF instrument-based, identify the specific instrument(s) and standard methods, modes, and languages of administration.

STS General Thoracic Surgery Database, Version 2.3

S.19. Data Source or Collection Instrument (available at measure-specific Web page URL identified in S.1 OR in attached appendix at A.1)

Available at measure-specific web page URL identified in S.1

S.20. Level of Analysis (Check ONLY the levels of analysis for which the measure is SPECIFIED AND TESTED)

Facility

S.21. Care Setting (Check ONLY the settings for which the measure is SPECIFIED AND TESTED)

Inpatient/Hospital

If other:

S.22. COMPOSITE Performance Measure - Additional Specifications (Use this section as needed for aggregation and weighting rules, or calculation of individual performance measures if not individually endorsed.)

2. Validity – See attached Measure Testing Submission Form

2.1 For maintenance of endorsement

Reliability testing: If testing of reliability of the measure score was not presented in prior submission(s), has reliability testing of the measure score been conducted? If yes, please provide results in the Testing attachment. Please use the most current version of the testing attachment (v7.1). Include information on all testing conducted (prior testing as well as any new testing); use red font to indicate updated testing.

2.2 For maintenance of endorsement

Has additional empirical validity testing of the measure score been conducted? If yes, please provide results in the Testing attachment. Please use the most current version of the testing attachment (v7.1). Include information on all testing conducted (prior testing as well as any new testing); use red font to indicate updated testing.

2.3 For maintenance of endorsement

Risk adjustment: For outcome, resource use, cost, and some process measures, risk-adjustment that includes social risk factors is not prohibited at present. Please update sections 1.8, 2a2, 2b1,2b4.3 and 2b5 in the Testing attachment and S.140 and S.11 in the online submission form. NOTE: These sections must be updated even if social risk factors are not included in the risk-adjustment strategy. You MUST use the most current version of the Testing Attachment (v7.1) -- older versions of the form will not have all required questions.

3. Feasibility

Extent to which the specifications including measure logic, require data that are readily available or could be captured without undue burden and can be implemented for performance measurement.

3a. Byproduct of Care Processes

For clinical measures, the required data elements are routinely generated and used during care delivery (e.g., blood pressure, lab test, diagnosis, medication order).

3a.1. Data Elements Generated as Byproduct of Care Processes.

Generated or collected by and used by healthcare personnel during the provision of care (e.g., blood pressure, lab value, diagnosis, depression score), Coded by someone other than person obtaining original information (e.g., DRG, ICD-9 codes on claims), Abstracted from a record by someone other than person obtaining original information (e.g., chart abstraction for quality measure or registry)

If other:

3b. Electronic Sources

The required data elements are available in electronic health records or other electronic sources. If the required data are not in electronic health records or existing electronic sources, a credible, near-term path to electronic collection is specified.

3b.1. To what extent are the specified data elements available electronically in defined fields (i.e., data elements that are needed to compute the performance measure score are in defined, computer-readable fields) Update this field for maintenance of endorsement.

ALL data elements are in defined fields in a combination of electronic sources

3b.2. If ALL the data elements needed to compute the performance measure score are not from electronic sources, specify a credible, near-term path to electronic capture, OR provide a rationale for using other than electronic sources. For maintenance of endorsement, if this measure is not an eMeasure (eCQM), please describe any efforts to develop an eMeasure (eCQM).

n/a

3b.3. If this is an eMeasure, provide a summary of the feasibility assessment in an attached file or make available at a measure-specific URL. Please also complete and attach the NQF Feasibility Score Card.

Attachment:

3c. Data Collection Strategy

Demonstration that the data collection strategy (e.g., source, timing, frequency, sampling, patient confidentiality, costs associated with fees/licensing of proprietary measures) can be implemented (e.g., already in operational use, or testing demonstrates that it is ready to put into operational use). For eMeasures, a feasibility assessment addresses the data elements and measure logic and demonstrates the eMeasure can be implemented or feasibility concerns can be adequately addressed.

3c.1. Required for maintenance of endorsement. Describe difficulties (as a result of testing and/or operational use of the measure) regarding data collection, availability of data, missing data, timing and frequency of data collection, sampling, patient confidentiality, time and cost of data collection, other feasibility/implementation issues.

IF instrument-based, consider implications for both individuals providing data (patients, service recipients, respondents) and those whose performance is being measured.

Missing data are sought by the DCRI from participants when the data are initially sent to DCRI for analysis.

Data are collected continuously by the participating sites and harvested by the DCRI twice yearly. Reports are then sent back to the sites about 3 months after a harvest.

No individual patient identifiers are collected by the DCRI.

Data Collection:

Participants of the STS General Thoracic Surgery Database generally have data managers on staff to collect these data. Costs to develop the measure included volunteer thoracic surgeons' time, STS staff time, and DCRI statistician and project management time.

Other fees:

STS General Thoracic Surgery Database participant surgeons pay an annual participant fee of \$550 or \$700, depending on whether the participant is an STS member or not. STS membership thus provides surgeons with a 21% discount on the non-member database participation fee.

3c.2. Describe any fees, licensing, or other requirements to use any aspect of the measure as specified (e.g., value/code set, risk model, programming code, algorithm).

See 3c.1

4. Usability and Use

Extent to which potential audiences (e.g., consumers, purchasers, providers, policy makers) 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.

4a. 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 (or the data on performance results are available). If not in use at the time of initial endorsement, then a credible plan for implementation within the specified timeframes is provided.

4.1. Current and Planned Use

NQF-endorsed measures are expected to be used in at least one accountability application within 3 years and publicly reported within 6 years of initial endorsement in addition to performance improvement.

Specific Plan for Use	Current Use (for current use provide URL)
	Public Reporting STS General Thoracic Surgery Database http://publicreporting.sts.org/gtsd Quality Improvement (external benchmarking to organizations) STS General Thoracic Surgery Database http://publicreporting.sts.org/gtsd Quality Improvement (Internal to the specific organization) STS General Thoracic Surgery Database http://publicreporting.sts.org/gtsd

4a1.1 For each CURRENT use, checked above (update for maintenance of endorsement), provide:

- Name of program and sponsor
- Purpose
- Geographic area and number and percentage of accountable entities and patients included
- Level of measurement and setting

See 4a1.2

4a1.2. If not currently publicly reported OR used in at least one other accountability application (e.g., payment program, certification, licensing) what are the reasons? (e.g., Do policies or actions of the developer/steward or accountable entities restrict access to performance results or impede implementation?)

STS is actively promoting public reporting of the STS adult cardiac, congenital heart, and general thoracic surgery performance measures. This is consistent with the explicitly stated STS philosophy that "As a national leader in health care transparency and accountability, The Society of Thoracic Surgeons believes that the public has a right to know the quality of surgical outcomes." (<http://www.sts.org/registries-research-center/sts-public-reporting>) In our efforts to operationalize public reporting, the STS Public Reporting Task Force has and will continue to develop public report cards that are consumer centric. Public reporting remains a top priority for the Society, and STS is striving for even stronger involvement among Database participants.

Currently, more than 650 Adult Cardiac Surgery Database (ACSD) participants voluntarily consent to be a part of the STS Public Reporting and more than 550 ACSD participants have consented to report publicly via the Consumer Reports public reporting initiative. Additionally, more than 100 Congenital Heart Surgery Database (CHSD) participants are currently enrolled in STS Public Reporting.

As of July 2017, General Thoracic Surgery Database (GTSD) participants were included in the Public Reporting initiative and more than 250 participants currently consent to report outcomes publicly on the STS website. This includes discharge mortality rate and median postoperative length of stay for lobectomy procedures for lung cancer, including scores and star ratings for the Lobectomy for Lung Cancer Composite Measure in addition to its domains of 1) absence

of mortality, and 2) absence of major complication. Participant outcomes are published alongside GTSD overall outcomes and National Inpatient Sample (NIS) outcomes.

-ACSD public reporting online may be found here: <http://publicreporting.sts.org/acsd>

-CHSD public reporting online may be found here: <http://publicreporting.sts.org/chsd>

-GHSD public reporting online may be found here: <http://publicreporting.sts.org/gtsd>

4a1.3. If not currently publicly reported OR used in at least one other accountability application, provide a credible plan for implementation within the expected timeframes -- any accountability application within 3 years and publicly reported within 6 years of initial endorsement. (*Credible plan includes the specific program, purpose, intended audience, and timeline for implementing the measure within the specified timeframes. A plan for accountability applications addresses mechanisms for data aggregation and reporting.*)

n/a

4a2.1.1. Describe how performance results, data, and assistance with interpretation have been provided to those being measured or other users during development or implementation.

How many and which types of measured entities and/or others were included? If only a sample of measured entities were included, describe the full population and how the sample was selected.

STS's combined mortality and morbidity model for pulmonary resection for lung cancer is important and appropriate for public reporting for the following reasons:

1.) within the broad category of lung cancer resections, lobectomy is the single most common major procedure that a thoracic surgeon performs;

2.) these procedures are therefore useful and appropriate to use as a benchmark for performance by general thoracic surgery programs. By providing surgeons and teams with risk-adjusted results, they can identify how they are performing compared with other programs in the STS General Thoracic Database, which generally includes the top thoracic programs in the nation. This will assist them in focusing performance improvement efforts. Also, when publicly reported, the outcomes for these common procedures provide patients and their families with comparative performance information to aid in selection of a provider;

3.) major morbidity is relatively common after lung resection; however, although mortality is rare, it should be captured as well in an outcome measure, thereby identifying ALL adverse events after lung resection;

4.) this measure is reported in an easy to understand format which summarizes the results of all participants who were included in the analysis. The participant's score is illustrated graphically in relation to the 25th, 50th and 75th percentiles of the distribution across participants, and is accompanied by the 95% Bayesian credible interval. Surgeons easily grasp this result and the visual display powerfully shows them just where they perform compared to their peers on a bi-annual basis. In addition, these risk-adjusted results allow surgeons to benchmark their program and initiate QI efforts, as needed. In providing transparency through public reporting of this measure, surgeons can better compare their patients' outcomes with national benchmarks and patients will be better informed consumers of health care.

4a2.1.2. Describe the process(es) involved, including when/how often results were provided, what data were provided, what educational/explanatory efforts were made, etc.

See 4a2.1.1

4a2.2.1. Summarize the feedback on measure performance and implementation from the measured entities and others described in 4d.1.

Describe how feedback was obtained.

The general thoracic surgeons from across the U.S. who comprise the STS General Thoracic Surgery Task Force meet periodically to discuss the participant reports and to consider potential enhancements to the GTSD.

Additions/clarifications to the data collection form and to the content/format of the participant reports are discussed and implemented as appropriate.

Most recently, STS surgeon members have expressed interest in real-time, online data updates, which has led to the development of dashboard-type reporting on STS.org. The general thoracic dashboard is scheduled for launch in 2018.

Also, general thoracic public reporting was initiated in the summer of 2017 (<http://publicreporting.sts.org/gtsd>), making star ratings for consenting participant groups available to participants as well as the public.

4a2.2.2. Summarize the feedback obtained from those being measured.

See 4a2.2.1

4a2.2.3. Summarize the feedback obtained from other users

Given the very recent launch of general thoracic public reporting, the STS has not yet received sufficient feedback from non-participants to be able to assess the impact of the public reporting initiative.

4a2.3. Describe how the feedback described in 4a2.2.1 has been considered when developing or revising the measure specifications or implementation, including whether the measure was modified and why or why not.

n/a

Improvement

Progress toward achieving the goal of high-quality, efficient healthcare for individuals or populations is demonstrated. If not in use for performance improvement at the time of initial endorsement, then a credible rationale describes how the performance results could be used to further the goal of high-quality, efficient healthcare for individuals or populations.

4b1. Refer to data provided in 1b but do not repeat here. Discuss any progress on improvement (trends in performance results, number and percentage of people receiving high-quality healthcare; Geographic area and number and percentage of accountable entities and patients included.)

If no improvement was demonstrated, what are the reasons? If not in use for performance improvement at the time of initial endorsement, provide a credible rationale that describes how the performance results could be used to further the goal of high-quality, efficient healthcare for individuals or populations.

Operative mortality in the STS General Thoracic Surgery Database has decreased from 2.2% in the years 2002 to 2008 to 1.4% from 2012 to 2014. These data represent the highest quality lung cancer surgery in the United States. It is important to recognize that a large proportion of the general thoracic surgery in the US is not performed by general thoracic surgeons certified by the American Board of Thoracic Surgery. Results by STS General Thoracic Database participants, who are almost all ABTS certified, are generally superior to those of surgeons performing these procedures who do not participate in the GTSD, and who are often not ABTS certified.

Kozower and colleagues (Ann Thorac Surg 2010) have previously demonstrated that compared with the Nationwide Inpatient Sample database, from 2002 to 2008, patients in the GTSD had lower unadjusted discharge mortality rates, median length of stay, and pulmonary complication rates for lobectomy.

The major morbidity rate has increased from 8.6% to 9.1% during the same time. A potential explanation for this observation is more complete coding of complications by data abstractors as the result of education efforts from STS, as well as inclusion of unexpected return to the operating room for any reason instead of only for bleeding.

Fernandez FG, Kosinski AS, Burfeind W, Park B, DeCamp MM, Seder C, Marshall B, Magee MJ, Wright CD, Kozower BD. The Society of Thoracic Surgeons Lung Cancer Resection Risk Model: Higher Quality Data and Superior Outcomes. Ann Thorac Surg. 2016 Aug;102(2):370-7.

Kozower BD, Sheng S, O'Brien SM, et al. STS database risk models: predictors of mortality and major morbidity for lung cancer resection. Ann Thorac Surg 2010;90:875–83.

4b2. Unintended Consequences

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).

4b2.1. Please explain any unexpected findings (positive or negative) during implementation of this measure including unintended impacts on patients.

We are not aware of any unexpected findings associated with implementation of this measure.

4b2.2. Please explain any unexpected benefits from implementation of this measure.

n/a

5. Comparison to Related or Competing Measures

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

5. Relation to Other NQF-endorsed Measures

Are there related measures (conceptually, either same measure focus or target population) or competing measures (conceptually both the same measure focus and same target population)? If yes, list the NQF # and title of all related and/or competing measures.

Yes

5.1a. List of related or competing measures (selected from NQF-endorsed measures)

1790 : Risk-Adjusted Morbidity and Mortality for Lung Resection for Lung Cancer

5.1b. If related or competing measures are not NQF endorsed please indicate measure title and steward.

n/a (measure #1790 is NQF endorsed, eligible for endorsement maintenance in this Surgery Project cycle)

5a. Harmonization of Related Measures

The measure specifications are harmonized with related measures;

OR

The differences in specifications are justified

5a.1. If this measure conceptually addresses EITHER the same measure focus OR the same target population as NQF-endorsed measure(s):

Are the measure specifications harmonized to the extent possible?

Yes

5a.2. If the measure specifications are not completely harmonized, identify the differences, rationale, and impact on interpretability and data collection burden.

Measure #1790 includes a broader range of lung resection procedures than the Lobectomy Composite, and therefore includes a larger number of cases and potentially provides performance data to more general thoracic surgeons. Of the two measures, only the Lobectomy Composite is currently publicly reported.

5b. Competing Measures

The measure is superior to competing measures (e.g., is a more valid or efficient way to measure);

OR

Multiple measures are justified.

5b.1. If this measure conceptually addresses both the same measure focus and the same target population as NQF-endorsed measure(s):

Describe why this measure is superior to competing measures (e.g., a more valid or efficient way to measure quality); OR provide a rationale for the additive value of endorsing an additional measure. (Provide analyses when possible.)

n/a

Appendix

A.1 Supplemental materials may be provided in an appendix. All supplemental materials (such as data collection instrument or methodology reports) should be organized in one file with a table of contents or bookmarks. If material pertains to a specific submission form number, that should be indicated. Requested information should be provided in the submission form and required attachments. There is no guarantee that supplemental materials will be reviewed.

Attachment Attachment: KozowerOBrienKosinski-et-al-2016-PIIS0003497515017531.pdf

Contact Information

Co.1 Measure Steward (Intellectual Property Owner): The Society of Thoracic Surgeons

Co.2 Point of Contact: Mark, Antman, mantman@sts.org, 312-202-5856-

Co.3 Measure Developer if different from Measure Steward: The Society of Thoracic Surgeons

Co.4 Point of Contact: Mark, Antman, mantman@sts.org, 312-202-5856-

Additional Information

Ad.1 Workgroup/Expert Panel involved in measure development

Provide a list of sponsoring organizations and workgroup/panel members' names and organizations. Describe the members' role in measure development.

Members of the STS Task Force on Quality Initiatives provide surgical expertise as needed. The STS Workforce on National Databases meets at the STS Annual Meeting and reviews the measures on a yearly basis. Changes or updates to the measure will be at the recommendation of the Workforce.

Measure Developer/Steward Updates and Ongoing Maintenance

Ad.2 Year the measure was first released: 2016

Ad.3 Month and Year of most recent revision: 01, 2016

Ad.4 What is your frequency for review/update of this measure? annually

Ad.5 When is the next scheduled review/update for this measure? 01, 2018

Ad.6 Copyright statement:

Ad.7 Disclaimers:

Ad.8 Additional Information/Comments: