

NATIONAL QUALITY FORUM

TO: NQF Members and Public

FR: NQF Staff

RE: Pre-voting review for *National Voluntary Consensus Standards for Cost and Resource Use (Cycle 2): A Consensus Report*

DA: October 20, 2011

Resource use measures count the frequency of defined health system resources, are broadly applicable and comparable measures of health services counts that are applied to a population or event. This project seeks to endorse cost and resource use measures, which will serve as building blocks for efficiency of care measures and signal the measure development industry of the urgent need to develop measures of efficiency that integrate quality domains with cost and resource use measures. This is NQF's first effort focused on endorsing cost and resource use measures.

Eleven measures were evaluated for suitability as voluntary consensus standards for accountability and performance improvement during review cycle two. Four condition-focused Technical Advisory Panels (TAPs) for pulmonary, cardiovascular and diabetes, bone and joint, and cancer conditions were convened to assist the project's 23-member Steering Committee in making recommendations. To date, the Steering Committee has recommended a total of eight cost and resource use measures for endorsement; four of which are a result of the second review cycle. For one measure (1595: ETG based diabetes cost of care measure), the Steering Committee vote was split and there was no consensus. The Committee seeks public and member comment on this measure to be considered in subsequent Committee discussion to determine the final recommendation for the measure.

The draft document, *National Voluntary Consensus Standards for Cost and Resource Use (Cycle 2): A Consensus Report*, is posted on the NQF website along with the following additional information:

- [measure submission forms](#); and
- [meeting and call summaries](#) from the TAP and Steering Committee's discussions.

This report builds upon the commenting draft report posted for cycle 1. Elements of the report that have been enhanced and added to the draft are reflected in **blue text**. Following the measure summaries, a new section has been added to the report entitled "Additional Considerations". In addition to commenting on the measures, we request that comments also be submitted on the "Applying Resource Use Measure Evaluation Criteria" and "Additional Considerations" sections of the report. Please submit your comments on these sections in the "General Comment" section of the commenting tool, including a reference to the page and line number in the report.

Pursuant to section II.A of the Consensus Development Process v. 1.8, this draft document, along with the accompanying material, is being provided to you at this time for purposes of review and comment only and is not intended to be used for voting purposes. You may post your comments and view the comments of others on the [NQF website](#).

NATIONAL QUALITY FORUM

**NQF Member comments must be submitted no later than 6:00 pm ET, November 21.
Public comments must be submitted no later than 6:00 pm ET, November 14.**

Thank you for your interest in NQF's work. We look forward to your review and comments.

NATIONAL QUALITY FORUM

NATIONAL VOLUNTARY CONSENSUS STANDARDS FOR COST AND RESOURCE USE (CYCLE 2): A CONSENSUS REPORT

DRAFT REPORT FOR COMMENTING

OCTOBER 20, 2011

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

NATIONAL VOLUNTARY CONSENSUS STANDARDS FOR COST AND RESOURCE USE (CYCLE 2): A CONSENSUS REPORT

TABLE OF CONTENTS

1		
2		
3		
4		
5	<u>EXECUTIVE SUMMARY</u>	4
6	<u>STRATEGIC DIRECTIONS FOR NQF</u>	8
7	<u>RELATED NQF WORK</u>	10
8	<u>RESOURCE USE MEASURES IN CONTEXT</u>	10
9	<u>NQF'S CONSENSUS DEVELOPMENT PROCESS</u>	12
10	<u>Evaluating Potential Consensus Standards</u>	13
11	<i>Principles for Resource Use Measure Evaluation</i>	13
12	<u>Applying the Resource Use Measure Evaluation Criteria</u>	14
13	<i>Importance</i>	15
14	<i>Scientific Acceptability</i>	15
15	<i>Usability</i>	27
16	<i>Feasibility</i>	28
17	<u>Harmonization & Best-in-Class</u>	29
18	<u>RECOMMENDATIONS FOR ENDORSEMENT</u>	30
19	<u>1560: Relative Resource Use for People with Asthma (NCQA)</u>	33
20	<u>1561: Relative Resource Use for People with COPD (NCQA)</u>	35
21	<u>1611: ETG-Based Pneumonia Cost of Care Measure (Ingenix)</u>	37
22	<u>1609: ETG/PEG-Based Hip/Knee Replacement Cost of Care Measure (Ingenix)</u>	39
23	<u>Candidate Consensus Standards Not Recommended for Endorsement</u>	42
24	<u>1591: ETG-Based Congestive Heart Failure (CHF) Cost of Care Measure (Ingenix)</u>	42
25	<u>1594 ETG-Based Coronary Artery Disease (CAD) Cost of Care Measure (Ingenix)</u>	45
26	<u>1599: ETG-Based Non-Condition Specific Cost of Care Measure (Ingenix)</u>	48
27	<u>1603: ETG/PEG-Based Hip Fracture Cost of Care Measure (Ingenix)</u>	51
28	<u>1605: ETG-Based Asthma Cost of Care Measure (Ingenix)</u>	53
29	<u>1608: ETG-Based Chronic Obstructive Pulmonary Disease Cost of Care Measure (COPD)</u>	
30	<u>(Ingenix)</u>	56

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

31	<u>Candidate Consensus Standards with No Committee Consensus</u>	58
32	<u>1595: ETG-Based Diabetes Cost of Care Measure (Ingenix)</u>	58
33	<u>ADDITIONAL RECOMMENDATIONS</u>	61
34	NEXT STEPS	74
35	NOTES	74
36	APPENDIX A—SPECIFICATIONS FOR COST AND RESOURCE USE MEASURES	
37	2011 (Cycle 2)	77
38	1560: Relative Resource Use for People with Asthma	77
39	1561: Relative Resource Use for People with COPD	79
40	1609: ETG-Based Hip/Knee Replacement Cost of Care Measure	81
41	1611: ETG-Based Pneumonia Cost of Care Measure	83
42	APPENDIX B—STEERING COMMITTEE	85
43	APPENDIX C—TECHNICAL ADVISORY PANELS	87
44	Cardiovascular/Diabetes Technical Advisory Panel	87
45	Pulmonary Technical Advisory Panel	88
46	Bone/Joint Technical Advisory Panel	89
47	Cancer Technical Advisory Panel	90
48	APPENDIX D—RESOURCE USE MEASUREMENT TERMS	91
49		
50		
51		
52		
53		
54		
55		
56		
57		

NATIONAL QUALITY FORUM

58 NATIONAL VOLUNTARY CONSENSUS STANDARDS FOR COST AND RESOURCE 59 USE (CYCLE 2): A CONSENSUS REPORT

60

61 EXECUTIVE SUMMARY

62 As current health reform efforts focus on expanding coverage, increasing access to care, and
63 reducing costs, it is important to understand how the system uses resources in the context of
64 health outcomes. Combining resource use (or cost) and quality data will enable the system to
65 better evaluate efficiency of care. Understanding resource use measurement as a building block
66 of efficiency is a first step toward this goal. For the purposes of this project, resource use
67 measures are defined as broadly applicable and comparable measures of health services counts,
68 in terms of units or dollars applied to a population or event (e.g., diagnoses, procedures, or
69 encounters). A resource use measure counts the frequency of defined health system resources;
70 some may further apply a dollar amount (e.g., allowable charges, paid amounts, or standardized
71 prices) to each unit of resource use.

72

73 This Consensus Development Process (CDP) project will endorse resource use (or cost)
74 measures that will serve as building blocks for efficiency of care measures and signal the
75 measure development industry of the urgent need to develop measures of resource use and
76 efficiency that integrate quality domains with resource use measures. In applying the Resource
77 Use Measure Evaluation Criteria for the first time, the Technical Advisory Panels (TAPs) and
78 Steering Committee encountered several overarching issues during their discussions and
79 evaluations of the measures. Some issues varied by developer as each developer submitted
80 measures with very distinct approaches. [This report reflects the discussion of those issues as well](#)
81 [as the measure-specific evaluation summaries for 11 measures reviewed during the second](#)
82 [review cycle.](#)

83

84 [In the second cycle of the project, four additional measures have been recommended for](#)
85 [endorsement as voluntary consensus standards suitable for accountability and performance](#)
86 [improvement:](#)

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

87

88

- (1560) Relative resource use for people with asthma (NCQA)

89

- (1561) Relative resource use for people with COPD (NCQA)

90

- (1609) ETG based hip/knee replacement cost of care measure (Ingenix)

91

- (1611) ETG based pneumonia cost of care (Ingenix)

92

93

94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

115 NATIONAL VOLUNTARY CONSENSUS STANDARDS FOR COST AND RESOURCE 116 USE (CYCLE 2): A CONSENSUS REPORT

117

118 BACKGROUND

119

120 The United States' healthcare expenditures are unmatched by any country in the world.¹ This
121 spending, however, has not resulted in better health for Americans. In fact, higher spending has
122 not decreased mortality, increased patient satisfaction, or led to improvements in access or higher
123 quality of care.^{2,3,4} This phenomenon of high spending with disproportionate outcomes points to
124 a system laden with waste. The contributing factors to this alarming trend are as complex as the
125 health care system itself, with physician practice patterns, regional market influences, and access
126 to care as major players. Meanwhile, the United States' healthcare spending continues to
127 increase at a rate of seven percent per year and is largely focused on treating acute and chronic
128 illness rather than preventive care.⁵

129

130 As ongoing health reform efforts focus on expanding coverage, increasing access to care, and
131 reducing costs, it is important to understand how resources are currently being used in the system
132 in the context of quality, preferably related to health outcomes. Linking resource use (or cost)
133 and quality measures will enable the system to better evaluate efficiency of care. Several
134 provisions in the Affordable Care Act (ACA), slated to be implemented over the next five years,
135 require using resource use data to further support efforts to move toward a value-based
136 purchasing (VBP) payment model. One such provision requires the Secretary of Health and
137 Human Services to develop an episode grouper that combines separate but clinically related
138 items and services into an episode of care for an individual.⁶ Additionally, resource use data will
139 be included on the physician compare website, as well as a physician value modifier that will be
140 used to adjust fee-for-service (FFS) payments by combining physician performance on quality
141 and resources use. While the ACA legislation is focused on the Medicare population,
142 understanding resource use measurement as a building block of efficiency, even in the context of
143 commercial-based measures, is a first step toward meeting these goals.

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

144 For the purposes of this project, resource use measures are defined as broadly applicable and
145 comparable measures of health services counts (in terms of units or dollars) that are applied to a
146 population or event (broadly defined to include diagnoses, procedures, or encounters). A
147 resource use measure counts the frequency of defined health system resources; some may further
148 apply a dollar amount (e.g., allowable charges, paid amounts, or standardized prices) to each unit
149 of resource use. Current approaches for measuring resource use range from broadly focused
150 measures, such as per capita measures, which address total healthcare spending (or resource use)
151 per person, to those with a more narrow focus, such as measures dealing with the healthcare
152 spending or resource use of an individual procedure (e.g., a hip replacement).

153 This second phase of a two-phase effort will endorse resource use measures through the
154 Consensus Development process (CDP). These measures will serve as building blocks for
155 efficiency of care measures and signal to the measure development industry the urgent need to
156 develop resource use and efficiency measures that integrate quality domains. Phase one, which
157 began in 2009, was aimed at understanding resource use measures and identifying the important
158 attributes to consider in their evaluation. During this phase, the current NQF Measure Evaluation
159 Criteria used to evaluate quality measures was reviewed and refined by the Resource Use
160 Steering Committee to address the unique aspects of resource use measures, resulting in the [NQF
161 Resource Use Measure Evaluation Criteria](#). A single Steering Committee was used across both
162 phases of work, with the addition of four Technical Advisory Panels (TAPs) in Phase two to
163 assist the Committee in evaluating the measures' clinical and methodological aspects. The CDP
164 project was divided into two review cycles, between which 14 focus areas were assigned:

165

166

167

168

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

Cycle 1

Cardiovascular

- Congestive heart failure (CHF)
- Coronary artery disease (CAD)
- Acute myocardial infarction (AMI)

Stroke

Diabetes

Non-condition specific (e.g. per capita-population)

Cycle 2

Pulmonary

- Chronic obstructive pulmonary disease (COPD)
- Asthma
- Pneumonia

Cancer

- Breast cancer
- Colorectal cancer

Bone/Joint

- Hip or knee replacement
- Hip or pelvic fracture

- Low back pain

169

170 This report reflects the discussion and overarching issues the Committee identified while
171 evaluating cost and resource use measures submitted to the project; measure-specific evaluation
172 summaries are provided for 11 measures reviewed during Cycles 1 and 2.

173

174 **STRATEGIC DIRECTIONS FOR NQF**

175 NQF's mission includes three parts: 1) building consensus on national priorities and goals for
176 performance improvement and working in partnership to achieve them; 2) endorsing national
177 consensus standards for measuring and publicly reporting on performance; and 3) promoting the
178 attainment of national goals through education and outreach programs. As greater numbers of
179 quality measures are developed and brought to NQF for consideration of endorsement, NQF
180 must assist stakeholders in measuring "what makes a difference" and addressing what is
181 important to achieve the best outcomes for patients and populations.

182

NATIONAL QUALITY FORUM

183 Several strategic issues have been identified to guide consideration of candidate consensus
184 standards:

185 **DRIVE TOWARD HIGH PERFORMANCE.** Over time, the bar of performance expectations
186 should be raised to encourage achievement of higher levels of system performance.

187 **EMPHASIZE COMPOSITES.** Composite measures provide much-needed summary information
188 pertaining to multiple dimensions of performance and are more comprehensible to patients and
189 consumers.

190 **MOVE TOWARD OUTCOME MEASUREMENT.** Outcome measures provide information of
191 keen interest to consumers and purchasers, and when coupled with healthcare process measures,
192 they provide useful and actionable information to providers. Outcome measures also focus
193 attention on much-needed system-level improvements because achieving the best patient
194 outcomes often requires a carefully designed care process, teamwork, and coordinated action on
195 the part of many providers.

196 **CONSIDER DISPARITIES IN ALL WE DO.** Some of the greatest performance gaps relate to
197 care of minority populations. Particular attention should be focused on identifying disparities-
198 sensitive performance measures and on identifying the most relevant
199 race/ethnicity/language/socioeconomic strata for reporting purposes.

200

201 **NATIONAL PRIORITIES PARTNERSHIP AND THE NATIONAL QUALITY** 202 **STRATEGY**

203 The [National Priorities Partnership](#), a multi-stakeholder collaborative of 48 organizations
204 convened by NQF, plays a key role in identifying strategies for achieving national goals for
205 quality healthcare and facilitating coordinated, multi-stakeholder action. The Department of
206 Health and Human Services has asked the Partnership for its collective, multi-stakeholder input
207 on the [National Quality Strategy](#) (NQS) framework, which includes three inextricably linked
208 domains—better care, affordable care, and healthy people/healthy communities—around which

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

209 priorities, goals, measures, and strategic opportunities for improvement are to be identified or
210 refined.

211
212 When the NQS was announced in March 2011, one of the priorities it identified was [making](#)
213 [quality care more affordable](#). The resource use measure endorsement process is an important step
214 toward measuring affordable care by evaluating resource use and cost measures. These measures
215 can identify opportunities to reduce the rate of growth in healthcare spending, and when paired
216 with quality measures, can help evaluate the efficiency of the healthcare system.

217

218 **RELATED NQF WORK**

219 This project is NQF's first effort focused on evaluating and endorsing cost and resource use
220 measures. In 2009, NQF completed a measurement framework for evaluating efficiency across
221 patient-focused episodes of care. This report, [NQF Measurement Framework: Evaluating](#)
222 [Efficiency across Patient-Focused Episodes of Care](#), presents the NQF-endorsed[®] measurement
223 framework for assessing efficiency, and ultimately value, associated with the care over the
224 course of an episode of illness and sets forth a vision to guide ongoing and future efforts.

225

226 **RESOURCE USE MEASURES IN CONTEXT**

227 This consensus development process seeks to endorse resource use (or cost) measures as
228 building blocks toward measuring efficiency of care. Efficiency can be defined broadly as the
229 resource use (or cost) associated with a specific level of performance with respect to the other
230 five Institute of Medicine (IOM) aims of quality: safety, timeliness, effectiveness, equity, and
231 patient-centeredness.⁷ Resource use measures can also be used to assess value by integrating
232 preference-weighted assessments of the quality and cost performance of a specified stakeholder,
233 such as an individual patient, consumer organization, payer, provider, government, or society.⁸

234

235 As a building block in understanding efficiency and value, NQF supports using and reporting of
236 resource use measures in the context of quality performance, preferably outcome measures.

237 Using resource use measures independent of quality measures does not provide an accurate

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

238 assessment of efficiency or value and may lead to adverse unintended consequences in the
239 healthcare system.

240
241 Resource use measures used to assess efficiency and value should be important to measure, have
242 scientifically acceptable properties, and be usable and feasible. Those resource use measures
243 under evaluation in this process should independently meet these endorsement standards. Future
244 efforts will need to evaluate how resource use measures can be paired with appropriate quality
245 measures to assess the healthcare system's efficiency. These efforts should consider quality and
246 resource measure alignment of the underlying population, exclusions, and risk adjustment,
247 among other measure properties.

248
249 Given the diverse perspectives on cost and resource use measurement in healthcare, it is
250 important to articulate, in the context of this project and the measures submitted, the
251 terminology, purpose, and perspectives these measures represented. Recognizing this is NQF's
252 first project in the resource use measurement arena, there is a clear gap in the NQF portfolio for
253 these types of measures. NQF also recognizes that while the measure submission process is open
254 to any entity wishing to submit measures for evaluation, the measures submitted and evaluated in
255 this process are not representative of all approaches to measuring healthcare costs and resources
256 that exist in the market today. This report is a reflection of the evaluation process of the
257 measurement approaches submitted to this project for review.

258
259 Each of the measurement approaches submitted for review calculate the use of various resources
260 using administrative claims data, categorize them by type of resource [e.g., pharmacy, durable
261 medical equipment, evaluation and management (E&M) visits] and apply a costing methodology
262 (either actual prices paid or standardized prices). When developers further apply a dollar value to
263 utilization counts, the dollar value serves as a weight for each resource. Due to the limitations in
264 the data types available for measuring resource use in healthcare, administrative claims data are
265 the primary source of this information for the measures submitted to this project. Further

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

266 discussion of costing approaches and the use of administrative claims data are addressed later in
267 the report.

268
269 Also important to understand in the context of this report is the way in which the terms “cost,”
270 “resource use,” and “prices” are used. The term “cost” can represent very different constructs to
271 various stakeholders. In the context of this report, cost (or cost of care measures) reflects the
272 actual prices *paid* by health plans for health plan member for utilization; resource use or
273 “resource use measures” further apply standardized prices to utilization counts. Prices charged
274 by providers in healthcare, by many accounts, is not a good measure of utilization as prices
275 charged can be a reflection of the negotiating position of health plans vis-à-vis providers in a
276 given market. Prices paid is generally a reflection of the cost the health plan incurs to cover the
277 claims submitted for its members; some measures also report a member (consumer) cost based
278 on member co-pays. For a provider, (e.g., a physician or nurse practitioner) a cost of care
279 measure would reflect the payment the provider received from the health plan for care provided.
280 For a purchaser, a resource use measure can be used to assess the utilization of healthcare
281 services across health plans, while a cost of care measure can be used to assess how well a health
282 plan is managing charges and utilization of providers within the health plan’s network. Given the
283 other types of costs attributed to healthcare, it is important to note that these measures do not
284 capture or represent production costs (fixed or any other costs to the provider to deliver care),
285 administrative costs, government funding to support healthcare delivery, or societal costs (e.g.,
286 lost wages, sick days).

287

288 **NQF’S CONSENSUS DEVELOPMENT PROCESS**

289 NQF’s National Voluntary Consensus Standards for Cost and Resource Use project seeks to
290 endorse resource use and cost measures for performance improvement and accountability in the
291 context of quality measures.

292

293

NATIONAL QUALITY FORUM

294 **Evaluating Potential Consensus Standards**

295 Candidate consensus standards were solicited through a Call for Measures on January 31, 2011.
296 Within the Cycle 2 condition areas, 19 measures were submitted and evaluated for suitability as
297 voluntary consensus standards for accountability; 12 of these were withdrawn by the developer.
298 The measures were evaluated using NQF Resource Use Measure Evaluation Criteria. Four
299 condition-focused TAPs for pulmonary, cardiovascular and diabetes, bone and joint, and cancer
300 conditions rated each candidate consensus standard according to the subcriteria and identified
301 strengths and weaknesses to assist the Committee in making recommendations. The 23-member,
302 multi-stakeholder Committee evaluated the subcriteria of the non-condition specific measures,
303 provided final evaluations of the four main criteria—importance to measure and report, scientific
304 acceptability of the measure properties, usability, and feasibility—and made endorsement
305 recommendations for all measures. Measure developers were available during TAP and
306 Committee discussions to respond to questions and clarify any issues or concerns.

307 ***Principles for Resource Use Measure Evaluation***

308 In Phase One of this project, the Committee defined resource use measures and their constructs
309 to better understand how to evaluate these measures. For the purposes of this project, resource
310 use measures are defined as broadly applicable and comparable measures of health services
311 counts (units or dollars) applied to a population or event (diagnoses, procedures, or encounters).
312 Resource use measure scores may be expressed as counts, dollars, or even observed-to-expected
313 ratios. The Committee developed the following principles to frame its subsequent efforts to
314 refine the evaluation criteria for resource use measures and evaluate resource use measures for
315 endorsement:

- 316 1. Efficiency is one of the Institute of Medicine (IOM) five quality aims and is a function
317 of resource use and health outcomes: *Efficiency = fx(resource use, health outcomes)*
- 318 2. Resource use measures are the amount of resources used per population, episode, or
319 procedure.

NATIONAL QUALITY FORUM

- 320 3. Resource use measures are an important building block for measures of efficiency of
321 care; future measurement efforts should integrate and explicitly incorporate measures of quality,
322 health outcomes, or appropriateness.
- 323 4. The justification for and intended purpose of resource use measures is to examine,
324 understand, and ultimately reduce unnecessary costs in care.
- 325 5. There is a continuum of resource use measures (i.e., per capita to per procedure); all types
326 under consideration for endorsement must meet NQF evaluation criteria for such measures.
- 327 6. The resource use measure specification and calculation must be explicitly stated and
328 transparent so the approach can be deconstructed and implemented in a standard manner.
- 329 7. Comprehensive measures are preferable, even if combining multiple service categories
330 into one resource use estimate increases complexity; using methodologically sound methods is of
331 paramount importance.
- 332 8. The final resource use measure result should be clear and understandable for all
333 stakeholders to interpret.
- 334 9. Methods for combining the component scores influence the interpretation of the measure
335 results and must be justified (e.g., averaging across all component scores may obscure low or
336 high scores of individual components).
- 337 10. While resource use measure developers may have fundamental differences in approach,
338 these principles should apply across all types and approaches.
- 339 11. NQF considers transparency as key to ensuring the intended audiences understand the
340 results and can use them for decision making. Resource use measures are often highly complex,
341 with lengthy algorithm decision trees that can make clarity difficult, particularly when some
342 components may be only partially transparent to the user.

343

344 **Applying the Resource Use Measure Evaluation Criteria**

345 With a working definition of resource use measures and guiding principles in place, the
346 Committee completed a detailed review of the standard NQF Measure Evaluation Criteria. This
347 review resulted in the NQF Resource Use Measure Evaluation Criteria, based on the same four
348 major criteria used to evaluate quality measures—importance, scientific acceptability, usability,

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

349 and feasibility—with targeted changes to the subcriteria to address the unique attributes of
350 resource use measures.

351
352 In applying the Resource Use Measure Evaluation Criteria for the first time, the TAPs and
353 Committee encountered several overarching issues during their discussions and evaluations of
354 the measures. Some issues varied by developer, as each developer submitted measures with very
355 distinct approaches. The Committee factored these issues into its ratings and recommendations
356 for multiple measures, recognizing the need to balance the quantity and specificity of
357 information required to evaluate adequately the measure and the burden on the developer to
358 provide this information. These issues are included below in the discussion of each criterion, in
359 addition to the summary provided of each individual measure in the evaluation summary table.
360

361 ***Importance***

362 The importance criterion for resource use measures, like that for quality measures, is aimed at
363 determining the extent to which the measure’s focus (e.g., hip fractures, coronary artery disease)
364 is important to measure and report. For resource use measures, the developers were asked to
365 demonstrate high impact by showing there is variation and opportunities for improvement in the
366 delivery of care for the identified condition. The TAP concluded that the measures submitted
367 were broad and inclusive of high-impact conditions. Additional subcriteria were tailored
368 specifically for resource use measures. These subcriteria included an evaluation of whether the
369 intent of the measure had been clearly described and whether the resource use service categories
370 selected to measure costs accurately reflected the intent and focus of the measure. All measure
371 submissions were found to be important.

372

373 ***Scientific Acceptability***

374 Similar to quality measures, evaluating the scientific acceptability of resource use measures
375 includes reviewing the measure’s specifications, reliability and validity testing, and approach to
376 addressing disparities. The completeness, repeatability of the specifications, and the adequacy of

NATIONAL QUALITY FORUM

377 the reliability testing methodology and results are evaluated within the reliability criterion.
378 Applying the validity criteria, the Committee was asked to determine whether the specifications
379 reflected the intent of the measure and address those areas where there was variation, as
380 demonstrated in importance. The validity criterion also includes an assessment of the adequacy
381 of validity testing, exclusions, risk-adjustment, and the identification of meaningful differences.

382

383 ***Resource Use Specification Modules***

384 The resource use measure specifications were delineated by five main modules, including: 1)
385 data protocol, 2) measure clinical logic, 3) measure construction logic, 4) adjustments for
386 comparability, and 5) measure reporting. To allow for user flexibility, the developers were
387 permitted to submit measurement steps in the data protocol and reporting modules as
388 specifications or guidelines, or to not submit instructions at all. Specifications are inherent
389 measure characteristics that must be fully implemented to obtain valid measure results.
390 Guidelines, on the other hand, are suggested approaches from the developer on possible ways to
391 implement these steps. Evaluation of resource use measure specifications proved to be the most
392 intensive effort in the review process. The issues identified within each of the specification
393 modules have been outlined below.

394

395 *Data protocol*

396 The data protocol module allows developers to submit instructions and analytic steps for
397 cleaning or aggregating relevant data necessary to implement the specifications and produce
398 valid results. Measure developers submitted the following data protocol information: data
399 preparation, data inclusion criteria, data exclusion criteria, and considerations for missing data.
400 Recognizing that not all developers create specifications around these steps, the Committee
401 concluded these items could be submitted as specifications or guidelines, or not submitted at all.

402

403 All of the measures submitted use administrative claims as the data source. Administrative
404 claims offer the benefit of reduced administrative burden for providers and measure
405 implementers in collecting and reporting data elements. However, variation in coding practices

NATIONAL QUALITY FORUM

406 has the potential to affect the reliability and validity of any measure that relies on administrative
407 and claims data alone, including resource use measures. This may be particularly true for entities
408 providing care under capitated financial arrangements that may capture fewer diagnostic and
409 procedural codes per record than those operating under traditional FFS arrangements.

410

411 Accountable entities may outsource services through pharmacy benefit managers (PBMs) or
412 behavioral/mental health carve-outs, which may result in incomplete or missing pharmacy or
413 behavioral/mental health data. These entities can outsource administration of outpatient
414 prescription drug benefits to PBMs.⁹ Carve-out arrangements allow accountable entities to
415 separate behavioral/mental health insurance benefits by contracting with a third party to manage
416 care or the insurance risk for patients requiring these services.¹⁰ The Committee agreed that total
417 resource use for entities that do not receive member claim information from carve-out pharmacy
418 and behavioral/mental health services may not be comparable to resource use for those that do
419 not outsource these services. In this instance, interpreting the overall costs for a patient across
420 health plans with and without carve-out arrangements would be misleading.

421

422 However, entities without member claims data from their carve-out arrangements can be flagged
423 for comparison with entities with similar missing benefit information. Because resource use
424 measures allow claims to be assigned to resource use categories (i.e., laboratory and imaging),
425 these categories can be used to compare costs across entities, even when outsourcing
426 arrangements are present. For example, comparing laboratory costs or imaging costs across
427 entities within a total per-capita resource use measure would be informative even when
428 pharmacy data are not available.

429

430 *Clinical logic*

431 Evaluation of the measure clinical logic included steps to identify the condition or event of
432 interest and any clustering of diagnoses or procedures. This evaluation included examining the
433 clinical topic area and determining whether or not the measure accounts for comorbid conditions,

NATIONAL QUALITY FORUM

434 disease interactions, clinical hierarchies, clinical severity levels, and concurrency of clinical
435 events.

436

437 The complexity of the submitted measure specifications made evaluating the measure's clinical
438 logic challenging. For example, measure developers designed various methodologies to assign
439 patients to a severity level; however, due to complex algorithms, specific details and code lists
440 used to determine the assignment of patients to severity categories were difficult to interpret.

441

442 Exclusions were a focus during evaluation of the resource use measure's clinical logic. Although
443 the creation of homogenous populations enables comparability, measure developers should
444 ensure that measure exclusions do not allow for complications from poor care to drive patients
445 out of the episode, thus rewarding entities that provide inadequate care. For example, a biased
446 measure score may be created by excluding patients with acute myocardial infarction (AMI) who
447 are discharged from a skilled nursing facility or excluding patients who are not discharged alive.

448

449 Finally, resource use measures that seek to create more homogenous patient populations often are
450 limited by the ability of administrative claims data to assess patient health status and severity
451 accurately. For example, measures submitted were unable to differentiate between community-
452 acquired and healthcare-acquired pneumonia. Measures submitted also were unable to identify
453 staging information to assess the severity of a cancer diagnosis.

454

455 *Construction logic*

456 The measure construction logic evaluation included a review of the steps used to cluster, group,
457 or assign claims beyond those associated with the measure's clinical logic and an assessment of
458 how the various components of the measure (episode logic, clinical logic, risk adjustment) work
459 together. Measures were evaluated to determine if the temporal parameters including trigger and
460 termination rules are appropriate for the clinical logic specified within the measure. For example,
461 the Committee evaluated the post-hospitalization period in an episode of AMI to ensure it was

NATIONAL QUALITY FORUM

462 appropriate for the measure’s intent, level of analysis, attribution approach, and statistical
463 properties.

464
465 The Committee evaluated the validity of the measures by examining the interaction of the
466 measure components including the specified level of analysis and the risk adjustment approach.
467 There is a need for nationally endorsed measures at the individual clinician level of measurement
468 and the experts encourage development of measures at this level. However, the Committee
469 expected developers to demonstrate statistical differences at sample sizes that would be observed
470 in the level of analysis specified. Further, attribution of the measure to the individual or group
471 practice level was discussed at length, focusing on the appropriateness and generalizability.
472 While sample size and attribution could be submitted as guidelines, the Committee agreed these
473 testing results contribute to the measure’s scientific acceptability.

474
475 Measures submitted as a part of an episode grouper were challenging to evaluate because the
476 assignment of claims into the episode, comorbidities and interactions, clinical hierarchies, and
477 the handling of concurrent of clinical events included lengthy algorithm decision trees that were
478 at times unclear and only partially transparent to the reviewers. Measures submitted to this
479 project were evaluated as standalone measures of resource use; however, the construction logic
480 within episode grouper-based approaches include claim assignment decisions, or tie-breaker
481 logic, which were not clearly explained in the evaluation of single resource use measures. Tie-
482 breaker logic is a mechanism to determine how a claim or record is assigned to an episode if it is
483 eligible for assignment to multiple episodes. For example, if a patient fills a prescription that
484 could be mapped to multiple open episodes, tie-breaking logic could be used to determine how
485 this cost would be assigned. The Committee expected developers to provide a clear and
486 transparent explanation of this tie-breaker logic, how claims would be assigned to episodes, and
487 how various open episodes interact with each other. While resource use measures are complex,
488 developers have a responsibility to provide an explanation of the construction logic within the
489 grouper; however the explanations submitted were often insufficient.

490

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

491 *Adjustments for comparability*

492 A measure's result can be influenced by confounding external factors that can affect the measure
493 score. Measure developers submitted steps for adjusting the measure to increase comparability.
494 These adjustments include risk adjustment, stratification approach, and the costing method used
495 within the measure.

496

497 Risk-adjustment methodologies varied considerably across measure developers. A combination
498 of complexity and a varying degree of transparency of the risk-adjustment approach made
499 evaluating the methods challenging. The experts agreed that the details on the performance of
500 risk models were vital to determining the model's adequacy—specifically, how the presence of
501 certain claims drives categorization into different risk categories and the goodness of fit of the
502 risk model. Of the various methodologies reviewed, none was considered to be superior. A
503 [Society of Actuaries report](#) shared with the Committee comparing various risk-adjustment
504 methodologies [e.g., Hierarchical Clinical Categories (HCC), Adjusted Clinical Groups (ACG),
505 Episode-risk-group (ERG)] was informative; however, more research and guidance on the
506 appropriateness of the models for specific applications are needed, as the Committee deemed this
507 report to be an inadequate analysis of the risk-adjustment models for the purposes of this project.
508 For example, the Committee asserted that risk-adjustment models be tested and may need to be
509 recalibrated based on the measure's target population. Guidance presented in the SOA report was
510 insufficient in assisting the Committee's assessment of risk-adjustment model performance
511 across various datasets, across various homogenous populations (including Medicaid or
512 Medicare), or the credibility of risk-adjustment models across various population sizes. The
513 Committee agreed that submissions lacking the necessary information to evaluate the risk model
514 fully should not be considered in future efforts to evaluate resource use measures. Descriptions
515 of the risk models should include model calibration statistics (i.e., the R-squared value), a
516 discussion of how variables were selected (i.e., based on statistical significance or clinical
517 indicators), and sensitivity analyses.

518

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

519 Stratification can be a mechanism to create homogenous risk populations; however, similar to the
520 concern that exclusions may remove patients out of an episode inappropriately, measure
521 developers need to ensure that the risk stratification approach does not allow for complications
522 from poor care to drive patients into a higher risk stratum, thus rewarding entities who provide
523 inadequate care. For example, for patients with coronary artery disease (CAD), creating risk
524 strata based on subsequent revascularization has this potential for adverse consequences.

525
526 The developers were asked to specify a costing method to apply to the measure. For the
527 measures submitted, the costing approaches were either specified for the actual prices paid (i.e.,
528 cost of care measures) or for standardized prices (i.e., resource use measure). Standardized
529 pricing allows users to compare the use and intensity of health services while holding actual paid
530 amounts constant. Resource use measures that apply standardized prices allow for comparison of
531 resource use units across regions and markets, while actual prices allow for comparison of prices
532 paid. The Committee agreed that both approaches could be appropriate for different applications;
533 however a measure used as a national consensus standard must select a single costing approach.
534 Including both costing approaches within the same measure could reduce comparability and limit
535 the user's ability to identify the source of variation. For this reason, developers that submitted a
536 single measure with an option for the user to determine which costing method to apply were
537 asked either to split the submission into two separate measures or select one of the approaches to
538 apply to a single measure submission. At the Committee's request, measures that were
539 unknowingly evaluated and voted on with optional costing approaches were re-voted during the
540 Cycle 2 Committee meeting based on developer selection of a single costing approach to be
541 applied (actual prices paid) to all of their measures.

542
543 Subsequent Committee discussions on applying an actual price approach for national
544 comparisons at an individual provider level identified additional concerns. Specifically, the
545 Committee noted the potential for misinterpreting physician resource use in national reporting.
546 This pricing approach includes environmental factors (i.e., local facility and labor costs) that may
547 be outside of an individual provider's control. The Committee agreed that when actual prices

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

548 paid are reported, utilization counts should be reported as well. The concern over the use of
549 actual prices also was considered in the measure’s usability. However, there was agreement that
550 actual prices paid by health plans to providers is important to measure and report; for example,
551 regional comparisons at the individual provider level where environmental factors may not be as
552 prominent, or nationally at higher levels of measurement (i.e. health plan level). Measures based
553 on actual prices paid are encouraged for endorsement, noting that the validity will be examined
554 through the interaction of the measure’s specified level of analysis, risk adjustment model, and
555 attribution approach.

556
557 Finally, measures submitted to this project spanned various levels of measurement analysis, from
558 regional, to health plan, to individual provider. Measures specified at a higher level of
559 measurement (i.e., health plan or regional) allowed for a comprehensive view of health service
560 resource use by measuring all costs for a person across settings and providers. While the
561 Committee encouraged measurement at the individual and group practice level, measures
562 submitted to this project had difficulty demonstrating reliability and validity at this level. *Across*
563 *all levels of measurement, the Committee engaged in a detailed evaluation of the risk adjustment*
564 *approach and minimum sample size to ensure that the measures produced a valid and reliable*
565 *score.*

566 567 *Reporting*

568 The reporting module includes steps for attribution, peer grouping, defining outliers and
569 thresholds, sample size requirements, and benchmarking. These reporting steps could be
570 submitted as measure specifications or guidelines, or could be left to the user’s discretion.
571 Specifications limit user options and flexibility and must be strictly adhered to, whereas
572 guidelines are well thought-out guidance to users, allowing for user flexibility.

573
574 While sample size considerations could be submitted as guidelines or specifications in the
575 reporting module, the Committee found that sample size was also relevant to the discussion of
576 other modules and reliability and validity testing. To evaluate the number of patients required for

NATIONAL QUALITY FORUM

577 a measure to demonstrate meaningful and statistically significant differences, the Committee
578 encouraged measure developers to provide simulations and sensitivity analyses during the
579 evaluation. When measures are specified at the individual provider level, confidence intervals
580 need to be presented, especially when displaying information with small sample sizes. Using
581 confidence intervals allows the user to assess the estimated range of the measure score and true
582 differences in provider performance.

583
584 Outliers were handled at both the episode and the claim level. During data preparation, high
585 outlier claims were generally subject to a statistical technique used to limit the effect of extreme
586 values and the effect of spurious outliers, known as *winsorization*.¹¹ Low cost claims were either
587 winsorized or, more typically, were removed from measure analysis. Winsorization often sets
588 outliers to a percentile of data; for example, all outliers above the 95th percentile are set to the
589 value at the 95th percentile. Developers who chose to remove low-cost episodes indicated they
590 took this approach because these episodes were likely to be incomplete and thus have the
591 potential to skew the results. The Committee requested additional details from the developers on
592 the effect of the winsorization and exclusion at the claim and episode level on the measure score.
593 The experts noted that detailed listing and analysis of high-cost outliers could be useful for
594 targeted improvement activities.

595
596 As part of the reporting module, the attribution approach could also be submitted as measure
597 guidelines or specifications or left to the user to define. The attribution approach is distinct from
598 the level of analysis in that the level of analysis is the unit in which the measure has been tested
599 and specified, while the attribution approach determines how the costs or resources are assigned
600 to a provider, group of providers, health plan, or region. Regardless of the approach submitted,
601 the Committee agreed that it should reasonably allow for the accountable entity to affect the
602 resource use of the patient. For example, if the attribution approach assigns a patient to the
603 primary care provider (PCP) based on one evaluation and management (E/M) visit, the approach
604 should not assign all of the previous hospitalization costs during the measurement year before the
605 patient's first visit to this PCP. Proper consideration should be given to how the timing of patient

NATIONAL QUALITY FORUM

606 encounters affects the attribution rules and potential for unfair assignment of costs to providers.
607 Lack of consideration for these types of factors creates the potential for unintended consequences
608 of providers “gaming the system” to avoid attribution of extraneous costs to their profile for new
609 patients with whom they have had limited contact.

610

611 *Approach to disparities*

612 Identifying and measuring disparities in care delivery is critically important to understanding
613 variations in cost and improving quality. Gender and age were the most common factors
614 accounted for in the stratification for disparities in the measures reviewed. The lack of
615 information on race and ethnicity in commercial administrative data limited the ability of the
616 resource use measures under evaluation to reflect disparities accurately in the results. Additional
617 efforts should be pursued to capture this information more systematically. The Committee was
618 unable to assess the measure’s ability to identify disparities based on underlying limitations in
619 the data. Measures were evaluated based on their ability to stratify if the underlying data
620 included information on race and ethnicity.

621

622 *Reliability and Validity testing*

623 The next component to evaluating a measure’s scientific acceptability is determining whether the
624 measure testing approach and results demonstrate that the measure is reliable and valid.
625 Reliability testing should demonstrate that the measure results are repeatable, producing the
626 same results a high proportion of the time when assessed in the same population in the same time
627 period, or that the measure score is precise. Validity testing must demonstrate that the measure
628 data elements are correct or that the measure score correctly reflects the cost of care or resources
629 provided, adequately distinguishing high and low resource use. If face validity is the only
630 validity addressed, it must be assessed systematically. Reliability and validity testing can be
631 demonstrated at the measure score or the data element level.

632

633

634

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

635 *Data element reliability*

636 Discussion of data element reliability was driven by the fact that the submitted resource use
637 measures relied on administrative claims data. Administrative claims provide accessible
638 information on the processes of care and can generally be obtained as a byproduct of the care
639 process. While administrative claims data reduces measure error due to manual chart abstraction
640 and transcription, developers cannot rely on the administrative claims to capture patient clinical
641 characteristics accurately without proper data element validity testing. Claims data provide only
642 limited clinical information, lack detail in determining patient health severity, and are subject to
643 variation in coding processes by the accountable entities. The Committee agreed that these
644 concerns span measures of quality and resource use and are not limited to the measures currently
645 under evaluation.

646

647 *Measure score reliability*

648 Measure developers also performed varying levels of reliability assessments at the measure score
649 level. The Committee was interested in assessing the measure's precision or ability to detect
650 signal rather than noise. Measures demonstrated lower levels of measure score reliability
651 assessments including parallel development of episode grouper software and SAS using the exact
652 same specifications. While these tests demonstrated match rates of more than 99.9 percent, they
653 do not facilitate assessments of the measure score's precision. Further, developers whose
654 measures have been in use attempted to demonstrate the reliability of the observed/expected
655 results (O/E) over time; however, doing so does not provide an assessment of precision of the
656 measure score. The Committee suggested other robust methodologies that could be used to
657 demonstrate a high level of reliability, including signal-to-noise ratio analysis using Analysis of
658 Variance (ANOVA) or intra-class correlation coefficient to demonstrate measure score
659 reliability.

660

661 *Data element validity*

662 The validity testing submitted at the data element level was often weak because there were no
663 comparisons to other independent claims databases or other authoritative data sources (e.g., the

NATIONAL QUALITY FORUM

664 patient's medical record). In addition, a comparison of the distribution of important variables to
665 the literature would provide a more robust assessment of the validity of the data elements used.

666
667 With the exception of developers who require regular data audits to ensure data integrity, the
668 measure submissions generally contained weak evidence of data integrity checks (i.e., percentage
669 of missing values, missing diagnosis codes, or inconsistent dates). However, developers often
670 provided guidelines for data preparation and missing data in the data protocol module.

671
672 Most measures submitted to the project were tested in large administrative claims databases
673 representative of the target population. The Committee noted one exception in which a hip
674 fracture measure was tested in a population with an age distribution outside of the age range in
675 which the condition was most prevalent. The TAP agreed this testing approach calls to question
676 the validity (and in fact the importance) of the measure as it has been tested and used to measure
677 costs in a population where this condition is not high impact and has limited clinical relevance.

678
679 *Measure score validity*

680 Validity testing at the measure score level often relied on face validity that the measure score
681 was valid based on clinical review and empirical results. The measure score, however, was often
682 not validated by correlating measure scores with other valid indicators or by showing that the
683 score produces different results when applied to subgroups known to have differences in
684 resource use, as a more complex validity testing approach would demonstrate. Developers often
685 demonstrated face validity by describing the distribution of measure score results, outlier status,
686 and type of service. While the Committee accepted this as a minimum threshold for
687 demonstrating validity, they suggested more robust methods, including correlating the measure
688 score with other valid indicators, should be applied in future iterations and testing.

689

NATIONAL QUALITY FORUM

690 ***Usability***

691 The focus of the usability criteria is to determine whether the measure results are usable for the
692 intended audience. This includes an evaluation of whether the measure is currently in use and the
693 results are being reported for performance improvement and accountability purposes, and
694 whether the results are considered meaningful and useful. For resource use measures, usability
695 also includes the evaluation of whether it has been demonstrated that the measure construct and
696 its components (e.g., risk-adjustment methodology, clinical logic) can be deconstructed to enable
697 transparency and understanding.

698

699 Resource use measures presented some specific challenges to applying the concepts identified
700 within the usability criterion. For example, the issue of accountability is a charged one. No
701 consensus existed as to who the intended audience of these measures should be—purchasers, the
702 public at large (consumers), health plans, and health plan members, are all likely users of this
703 information. It was noted that for the public at large, extra effort would be required to make the
704 reporting of these measure results as clear as possible; ensuring clarity is the focus of consumer-
705 oriented organizations that share data such as these. There was agreement that these measures
706 should not be reported alone, but in the context of quality measures.

707

708 Another challenge the TAPs and Committees encountered was differentiating between usability
709 and usefulness and determining whether a measure is inherently usable because it is in use. For
710 measures not currently in use, they questioned how usefulness should be demonstrated since
711 there is a lack of knowledge of the practical application of the measure.

712

713 The Committee also questioned the usability of measures that are embedded in a complex
714 episode-grouper system in which each individual measure's logic is interwoven and tied to the
715 logic of another measure, which may not be under evaluation. They struggled with how to
716 evaluate the usability of a single measure without evaluating the entire grouper system.

717

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

718 The final overarching issue identified within the usability criteria relates to transparency. Many
719 of the TAP and Committee members expressed concern over the complexity of certain
720 methodologies used and questioned whether this complexity masks these measures' ability to be
721 transparent. Difficulty understanding how the risk adjustment, severity level assignments, and
722 episode logic work together in a measure may make it difficult for a physician, for example, to
723 understand completely which of his or her patients have been included in the costs attributed to
724 them and how the complexity of the patient population has been accounted for in those costs.
725 Some Committee members argued that this lack of transparency and understanding of the
726 construction logic affects the ability of the reported measure score to be used and may limit the
727 physician or health plan from identifying how and where to improve scores. Committee members
728 also questioned whether there should be an expectation that these complex measures would
729 require an investment of time to be interpreted and understood. It was pointed out, however, that
730 by using the resource use service categories identified within the measure, action could be taken
731 using the categories in which high costs were most evident (e.g., imaging, outpatient visits).

732 ***Feasibility***

733 The feasibility criterion focuses on the extent to which the measure can be implemented with
734 undue burden and identifies any barriers to implementation. The feasibility subcriteria used to
735 evaluate the resource use measures are identical to those used to evaluate quality measures.
736 Because all of the resource use measures submitted to this project solely rely on administrative
737 claims data, the subcriteria evaluating the availability of required data via electronic sources and
738 whether the data are routinely generated required very little discussion. The remaining feasibility
739 subcriteria, however, illuminated some important issues related to implementing resource use
740 measures, which often use very complex, sophisticated methodologies to adjust risk and
741 determine episode logic, for example. The TAPs and the Committee discussed this issue of
742 complexity for the implementer (and for the users of the results) during their evaluation of
743 susceptibility to errors and inaccuracies. Some members expressed concern that the complexity
744 of the methodologies lends itself to user error, most likely on behalf of the programmer who
745 would develop the code to run the measures. This issue may be mitigated by the purchase of a

NATIONAL QUALITY FORUM

746 product that is pre-programmed to implement the measure with imported data or the submission
747 of data to an organization that audits, computes the measure, and reports the information back to
748 the user.

749
750 Additionally, having been in use in the marketplace by health plans and purchasers for many
751 years, these measures often use some proprietary component or are imbedded in sophisticated
752 proprietary products. For product lines that include large episode-grouping tools encompassing
753 many conditions, a user would be required to purchase some or parts of a product suite to run a
754 single episode for diabetes, for example. For this reason, the feasibility of implementing an
755 individual clinical episode may be very limited. The Committee expressed concern that the
756 financial burden on a practice or system to purchase these products could be very significant,
757 thus creating a barrier to measuring resource use applying NQF-endorsed standards.

758

759 **Harmonization and Best-in-Class**

760 In Phase One of this resource use measurement project, the Committee agreed that because this
761 is NQF's first effort focused on evaluating resource use measures, identifying "best-in-class" and
762 requiring harmonization among resource use measures would be premature. In the context of
763 resource use measures, similar measures may share the same measure type (e.g., per episode, per
764 capita), or measure the same costs/resources (e.g., actual prices paid vs. standard prices, resource
765 service categories), or address the same population (e.g., people with diabetes). Competing
766 measures would share all of the characteristics previously listed. Among the eight measures
767 recommended for endorsement, there were no competing measures. Recommended measures
768 that were the same measure type were submitted from the same developer and were already
769 harmonized. With the exception of the two non-condition-specific total cost of care measures
770 (submitted by the same developer and recommended in Cycle 1), which employ different costing
771 methodologies, all recommended measures addressed different populations. Future resource use
772 measure endorsement efforts should explore the potential ways in which harmonization among
773 similar measures might be achieved. Specifically, identifying which measure constructs (e.g.,

NATIONAL QUALITY FORUM

774 condition-specific episode trigger and end mechanisms, age ranges), if any, could be harmonized
775 for standard measurement is needed in this measurement area. Also, exploring the implications
776 of harmonization for the resource use measure development community in which proprietary
777 measure components are common would be useful as the portfolio of endorsed resource use
778 measures expands.

779

780 **RECOMMENDATIONS FOR ENDORSEMENT**

781 This report presents the results from the evaluation of 11 measures considered during review
782 cycles one and two under NQF's CDP.

783

784 ***Evaluation of Measure Costing Approaches***

785 Early in the evaluation process, the Committee agreed that it was important to distinguish
786 measure results obtained using standardized prices and actual prices paid; dividing the costing
787 approaches into separate measures was determined to be the best approach to ensure this
788 distinction was made for standardized implementation and prevent inaccurate comparisons. As
789 such, developers that submitted a single measure with an option for the user to determine which
790 costing method to apply, were asked either to split the submission into two separate measures, or
791 select one of the approaches to apply to a single measure submission. This was requested of both
792 HealthPartners (in cycle one) and of Ingenix (in cycles one and two). HealthPartners
793 subsequently resubmitted two separate measures, one applying each costing approach; Ingenix
794 resubmitted all of their measures applying only actual prices paid.

795

796 During the initial evaluation and voting for recommendation of the Ingenix measures, there was
797 not a shared understanding among the Committee that the measures had been submitted with
798 flexibility in the costing approach. Ingenix chose to resubmit their measures using actual prices
799 paid. Once the measures were resubmitted to the Committee applying the single costing
800 approach, the Committee was given the opportunity to determine if the selection in the costing
801 approach warranted a re-vote. The Committee requested a revote since there was not a shared
802 understanding on the original costing approach by Ingenix, thus all Ingenix measures were re-

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

803 voted on during the Cycle 2 Committee meeting. The re-vote was for overall recommendation
804 for endorsement only. This is reflected as such in the measure evaluation summaries below.

805

806 ***Evaluation of Measurement Approaches***

807 The NQF measure evaluation process calls for each submitted measure to be evaluated
808 individually, based on its own merit. This was also the approach used in this project.

809 Additionally in this project, given the nature of the various of resource use measure developers,
810 measures developed by a single developer shared many common underlying measure constructs
811 and processes. By understanding the common constructs shared among a group of measures from
812 a developer (i.e. general methods), it lays the foundation for understanding the nuances specific
813 to each individual measure. During the measure evaluation process, the TAPs and Committees
814 often identified some recurring themes within the criteria discussions that applied across
815 measures from an individual developer, regardless of condition focus of the individual measure.
816 Some of these recurring themes have been captured in several of the measure evaluation
817 summaries and some have been identified below.

818

819 ***Ingenix Feasibility***

820 Each of the individual Ingenix measures [(Episode Treatment Groups (ETGs))] exist as part of a
821 larger grouper system, and requires the use of the entire grouper to produce results for the
822 individual ETGs. Because, each of the condition-specific ETGs submitted to this project require
823 the use of the Ingenix grouper product to implement the measures, the Committee's discussion of
824 the feasibility criterion for these measures was done for all of these measures at one time. As a
825 part of feasibility discussion, the Committee was provided with a pricing table for each of the
826 products required for implementation of these condition-specific ETGs.

827

828 Because these measures primarily use administrative claims data, all of the data required to
829 implement these measures is generated as a byproduct of care and is available electronically.
830 There was concern around the measure's susceptibility to inaccuracies as Ingenix does not have a
831 formal audit system to ensure that all of data is included and correct. In terms of barriers to use,

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

832 the purchase and implementation of this product could be cost prohibitive for some entities.
833 Annually, for physicians the cost to implement this project could range from of the small
834 package \$70,000 (for a group of less than 800 physicians) to \$110,000 (for over 2,000 physicians
835 in the group). For health plans, the annual cost could range from \$90,115 (for less than 400,000
836 covered lives) to is \$135,000 (for over a million covered lives). The Steering Committee
837 concluded that this cost is comparable to the cost of other proprietary fees associated with other
838 risk adjustment models of its caliber (e.g., ACGs used by HealthPartners). These prices include
839 costs associated with the licensure of the proprietary software and the cost of all of their
840 measures, over 558 ETGs, but not implementation. The Steering Committee acknowledged that
841 while the methodology is very complex, the system may be used without Ingenix’s technical
842 support, if the user spends time thoroughly reviewing the documentation.

843

844 **Candidate Consensus Standards Recommended for Endorsement**

845 Four measures are recommended for endorsement as voluntary consensus standards suitable for
846 accountability and performance improvement.

847

848 The evaluation summary tables follow the list of measures and summarize the results of the
849 TAP’s and Committee’s evaluation of and voting on the candidate consensus standards that were
850 recommended for endorsement. Hyperlinks are provided from each summary table to the
851 detailed measure specifications. To access the meeting transcripts and recordings in which these
852 measures are discussed, refer to the [project web page](#).

853

854 The Committee recommended the following candidate consensus standards for endorsement:

855 **Pulmonary**

856 (1560) Relative Resource Use for People with Asthma (NCQA).....32

857 (1561) Relative Resource Use for People with COPD (NCQA).....34

858 (1611) ETG-Based Pneumonia Cost of Care (Ingenix).....37

859 **Bone/Joint**

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

860 (1609) ETG/PEG-Based Hip/Knee Replacement Cost of Care Measure (Ingenix).....39

861

862

863

864

865 **Evaluation Summary—Candidate Consensus Standards Recommended for**
 866 **Endorsement**

<p>1560: Relative Resource Use for People with Asthma (NCQA)</p>
<p>Description: This measure addresses the resource use of members identified as having asthma. Both encounter and pharmacy data are used to identify members for inclusion in the eligible population, and the results are adjusted to account for age, gender, and HCC-RRU risk classifications that predict cost variability (Refer to Attachment S8_Clinical Logic for additional information).</p> <p>Resource Use Type: Per capita (population- or patient-based)</p> <p>Data Type: Administrative claims; Electronic Clinical Data : Electronic Health Record; Electronic Clinical Data : Imaging/Diagnostic Study; Electronic Clinical Data : Laboratory; Electronic Clinical Data : Pharmacy Paper Records</p> <p>Resource Use Service Categories: Inpatient services: Inpatient facility services, Inpatient services: Evaluation and management, Inpatient services: Procedures and surgeries, Inpatient services: Imaging and diagnostic, Inpatient services: Lab services, Inpatient services: Admissions/discharges, Ambulatory services: Outpatient facility services, Ambulatory services: Emergency Department, Ambulatory services: Pharmacy, Ambulatory services: Evaluation and management, Ambulatory services: Procedures and surgeries, Ambulatory services: Imaging and diagnostic, Ambulatory services: Lab services</p> <p>Care Setting: Ambulatory Care : Ambulatory Surgery Center (ASC), Ambulatory Care : Clinic/Urgent Care, Ambulatory Care : Clinician Office, Hospital/Acute Care Facility, Imaging Facility, Laboratory, Pharmacy, Post Acute/Long Term Care Facility: Nursing Home/Skilled Nursing Facility</p> <p>Level of Analysis: Clinician: Group/Practice, Health Plan, Integrated Delivery System, Population : National, Population: Regional</p> <p>Measure Developer: National Committee for Quality Assurance (NCQA), 1100 13th Street NW, STE 1000, Washington, District Of Columbia, 20005</p>
<p>Committee Recommendation for Endorsement: Y-13; N-0; Abstain-1</p>
<p>Conditions/Questions for Developer:</p> <ol style="list-style-type: none"> 1. Could this measure be improved by including other diagnostic criteria to ensure all appropriate asthma patients are captured? 2. How have you come up with the age strata in your risk-adjustment? 3. Can secondary diagnosis be taken into account within the measurement year? 4. Is cost during the measurement year part of the risk-adjustment strategy? 5. Are your measure results published publically? <p>Developer Response:</p> <ol style="list-style-type: none"> 1. Using asthma as a principal diagnosis will make it difficult to identify most patients, especially those who are acute and come into the ER and are diagnosed with bronchitis first, and then asthma. 2. The age strata for risk-adjustment are designed around known utilization patterns and clinical treatment patterns. 3. All costs for anyone with asthma are counted. 4. The HCC uses any services during the year to appropriately categorize patients into those 13 risk cohorts by severity of comorbidity. They also look at ICD-9 and procedural codes to categorize them and then go back and look at the number of times those services were offered to that population. Therefore, if a patient has multiple co-morbidities, that factors into the risk-adjustment, and will put a patient into a more severe risk-adjustment category. 5. Results are published through NCQA's Quality Compass module which contains the individual plan results by detailed service category along with a quality score.
<p>1. Importance to Measure and Report</p> <p>1a.High Impact: H-9; M-0; L-0; I-0</p> <p>TAP Discussion: The TAP agrees that asthma is an important area of healthcare to measure due to its high cost and the potential for improvements in care.</p> <p>1b. Resource use/cost problems: H-7; M-2; L-0; I-0</p>

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

<p>TAP Discussion: The TAP agrees that asthma represents a resource use problem and noted that there is a well-documented opportunity for improvement.</p> <p>1c. Purpose clearly described: H-9; M-0; L-0; I-0</p> <p>Discussion: The TAP believes the purpose and objective are clear; this subcriterion has been met.</p> <p>1d. Resource use service categories consistent and representative: H-9; M-0; L-0; I-0</p> <p>TAP Discussion: The TAP believes this subcriterion has been met; there were not issues raised.</p>
<p>Overall Importance: Y-16, N-0</p> <p>Committee Discussion: The Steering Committee agrees this criterion has been met.</p>
<p>2. Scientific Acceptability of Measure Properties:</p> <p>2a. Overall Reliability: H-8; M-1; L-0; I-0</p> <p>2a1. Measure well defined and precisely specified: H-9; M-0; L-0; I-0</p> <p>TAP Discussion: The TAP believes this subcriterion has been met.</p> <p>2a2. The results are repeatable: H-8; M-1; L-0; I-0</p> <p>TAP Discussion: There was general agreement from the TAP that following a methodology of including <i>all</i> costs avoids having to consider what costs should or should not be associated with asthma. The developer reaffirmed that the measures are valid for any health plan; they are population-based measures and have been tested and can be used in physician groups with a sufficient number of patients. A population of at least 400 members is needed for the methodology to be valid, so it consequently tends to be larger physician groups that can use the measures.</p> <p>2b. Overall Validity: H-5; M-4; L-0; I-0</p> <p>2b1. Evidence is consistent with intent: H-6; M-3; L-0; I-0</p> <p>TAP Discussion: The TAP agrees there is good overall evidence of face validity, but also a general desire to see more specific discussion around the face validity of the use of HCC's in this population.</p> <p>2b2. Score/Analysis: H-6; M-3; L-0; I-0</p> <p>TAP Discussion: The face validity of HCC's was found to be clear, but the logic behind the age stratification was unclear. 2b3. Exclusions: H-6; M-3; L-0; I-0</p> <p>TAP Discussion: The TAP had an in-depth discussion regarding measure exclusions. The measure developer explained that cardiovascular conditions are not specifically excluded, but are used in the risk adjustment model. Patients with COPD are excluded. Exclusions affect the denominator population over either year within the two-year criteria, which is similar to the HEDIS asthma measure. There was agreement that the exclusion of COPD (which resulted in 38% of the initial population being eliminated) seems appropriate, particularly in light of the age range increasing to 64. The TAP did express concern that excluding acute respiratory failure could exclude poorly managed asthma patients. However, NCOA noted that acute respiratory failure only accounted for 3% of the population, so it doesn't meet their 5% threshold of concern.</p> <p>2b4. Risk Adjustment: H-7; M-2; L-0; I-0</p> <p>TAP Discussion: The TAP believes the risk-adjustment strategy seems appropriate. Several strategies are tested by NCOA, and the same methodology is used for all of their measures. The developer stratifies the population by age and gender and uses HCC's to risk adjust the population.</p> <p>2b5. Identification of statistically significant/meaningful differences: H-8; M-1; L-0; I-0</p> <p>TAP Discussion: There was general agreement that the distribution of the scores' detail score was appropriate. There was concern regarding whether the measure score could differentiate statistically significant and clinically significant variation.</p> <p>2b6. Multiple data sources: N/A</p> <p>2c. Stratification for disparities: H-5; M-3; L-0; I-1</p> <p>TAP Discussion: The TAP believes stratification is needed although the data isn't available at this time.</p>
<p>Overall Scientifically Acceptable: Yes [Y-12; N-2 (Committee Vote)]</p> <p>Overall Reliability: H-12; M-3; L-0; I-0</p> <p>Overall Validity: H-4; M-9; L-1; I-0</p> <p>Committee Discussion: The Committee agreed with the TAP's analysis of reliability and raised no additional concerns. There was further discussion around missing pharmacy data, and confirmation that plans submit separate components (total medical, quality, and pharmacy, for example) to NCOA and are allowed to have a certain number of missing components. NCOA then holds the plans accountable for ensuring that they have the complete data required to report the measure, and any plans that are missing a major component of the measure specification would not end up in the NCOA reporting product. The Committee asked the developers to defend the measure's use of indirect standardization in creating standardized prices.</p>
<p>Usability:</p>

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

<p>3a. Measure performance results are publicly reported: H-8; M-1; L-0; I-0 <i>TAP Discussion:</i> The TAP was satisfied that NCOA publically reports measure results and provides support to enable understanding of those results. Purchasers are using this information, along with NCOA quality measures, to improve value for their employees. Asthma is a bit more difficult because there is only one NCOA quality measure to associate with this cost measure, however there are more quality measures in the pipeline.</p> <p>3b. Measure results are meaningful/useful for public reporting and quality improvement: H-6; M-3; L-0; I-0 <i>TAP Discussion:</i> The measure is straightforward and easy to interpret. NCQA uses standardized pricing tables, which are reviewed annually. Health plans are the main users for this data. However, purchasers and the large employers will also drive a need for this information. The TAP wondered how smaller businesses would implement this measure, and NCQA explained that they provide help through their annual conferences, webinar services and a dedicated webpage.</p> <p>3c. Data and results can be decomposed for transparency and understanding: H-8; M-1; L-0; I-0 <i>TAP Discussion:</i> The TAP believes the methodology was transparent and appropriate.</p> <p>3d. Harmonized or justification for differences: N/A</p>
<p>Overall Usability: H-9; M-5; L-0; I-0 Committee Discussion: The Steering Committee was concerned about the ability of small groups to implement this measure.</p>
<p>4. Feasibility:</p> <p>4a. Data elements routinely generated during care process: H-9; M-0; L-0; I-0 <i>TAP Discussion:</i> The TAP agrees this subcritierion has been met; the data is a byproduct of care.</p> <p>4b. Data elements available electronically: H-9; M-0; L-0; I-0 <i>TAP Discussion:</i> The TAP agrees this subcritierion has been met; the data is available electronically.</p> <p>4c. Susceptibility to inaccuracies/ unintended consequences identified: H-7; M-2; L-0; I-0 <i>TAP Discussion:</i> There was agreement that NCQA did a sufficient job recognizing where the challenges with data inaccuracies are and have adequately addressed these challenges.</p> <p>4d. Data collection strategy can be implemented: H-8; M-1; L-0; I-0 <i>TAP Discussion:</i> All the data submitted to NCQA must go through a certified auditor before it's reported to NCQA. As part of their annual analysis, NCQA reviews outliers, but currently the outliers are less than half a percent for this measure.</p>
<p>Overall Feasibility: H-10; M-4; L-0; I-0 Committee Discussion: No additional concerns were raised by the Steering Committee regarding feasibility.</p>

867

<p>1561: Relative Resource Use for People with COPD (NCQA)</p> <p>Description: This measure addresses the resource use of members identified with COPD. Clinical diagnosis of COPD during the measurement year is used to identify members for inclusion in the eligible population and the results are adjusted to account for age, gender, and HCC-RRU risk classifications that predict cost variability (Refer to Attachment S8_Clinical Logic for additional information).</p> <p>Resource Use Type: Per capita (population- or patient-based)</p> <p>Data Type: Administrative claims, Electronic Clinical Data: Electronic Health Record, Electronic Clinical Data: Imaging/Diagnostic Study, Electronic Clinical Data : Laboratory, Electronic Clinical Data : Pharmacy, Paper Records</p> <p>Resource Use Service Categories: Inpatient services: Inpatient facility services, Inpatient services: Evaluation and management, Inpatient services: Procedures and surgeries, Inpatient services: Imaging and diagnostic, Inpatient services: Lab services, Inpatient services: Admissions/discharges, Ambulatory services: Outpatient facility services, Ambulatory services: Emergency Department, Ambulatory services: Pharmacy, Ambulatory services: Evaluation and management, Ambulatory services: Procedures and surgeries, Ambulatory services: Imaging and diagnostic, Ambulatory services: Lab services</p> <p>Care Setting: Ambulatory Care : Ambulatory Surgery Center (ASC), Ambulatory Care : Clinic/Urgent Care, Ambulatory Care : Clinician Office, Hospital/Acute Care Facility, Imaging Facility, Laboratory, Pharmacy, Post Acute/Long Term Care Facility: Nursing Home/Skilled Nursing Facility</p> <p>Level of Analysis: Clinician : Group/Practice, Health Plan, Integrated Delivery System, Population: Community, Population: National, Population : Regional</p> <p>Measure Developer: National Committee for Quality Assurance (NCQA), 1100 13th street NW, STE 1000, Washington, District Of Columbia, 20005</p>
<p>Committee Recommendation for Endorsement: Y-13; N-0; Abstain-1</p>
<p>Conditions/Questions for Developer:</p>

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

1. If the goal is to eventually link these measures with quality measures and stratification is different, how will that be plausible?
2. What is the upper age limit to be included in this measure?
3. How do you ensure similar populations are compared?

Developer Response:

1. The resource use strata are different than they are for clinical quality strata, which are not risk-adjusted. As the quality measures further increase and perhaps in the future become risk-adjusted, there will be more room for comparability.
2. There is no upper age limit to this measure.
3. By risk adjusting to the specified level using the HCC's and the 13 different cohorts, NCQA end up comparing relatively similar plan populations. The quality index for this measure is use of diagnostic spirometer and exacerbations measures. There is no attribution of specific procedures to COPD yet.

1. Importance to Measure and Report

1a. High Impact: H-9; M-0; L-0; I-0

TAP Discussion: The TAP was in agreement that this is an important area of measurement.

1b. Resource use/cost problems: H-9; M-0; L-0; I-0

TAP Discussion: The TAP believes while there is variation in resource use was identified in other parts of the submission, the information submitted in the form for this item only discussed the variations in clinical care provided.

1c. Purpose clearly described: H-8; M-1; L-0; I-0

TAP Discussion: The TAP was concerned that the measure submission applied only to newly diagnosed patients. The developer clarified that it is supposed to apply to anyone with a diagnosis with COPD. Otherwise, the purpose of the measure is to evaluate the total cost of care for COPD patients within a 1 year timeframe was clear.

1d. Resource use service categories consistent and representative: H-9; M-0; L-0; I-0

TAP Discussion: The TAP believes this subcriterion has been met.

Overall Importance: Y-14, N-0

Committee Discussion: The Steering Committee agreed the measure focused on an important area of healthcare.

2. Scientific Acceptability of Measure Properties:

2a. Overall Reliability: H-7; M-2; L-0; I-0

2a1. Measure well defined and precisely specified: H-9; M-0; L-0; I-0

TAP Discussion: The TAP believes the specifications provided are clear and precise. The developer provided clarification on age stratification for resource use categories indicating that they are based on utilization patterns in the data-set, not clinical factors.

2a2. The results are repeatable: H-8; M-1; L-0; I-0

TAP Discussion: A similar methodology was used for this measure as for NCQA measure #1560, the primary difference being in the selection of the population. The TAP was concerned about the multiple populations being studied including commercial, Medicare, and Medicaid, due to the age range (unlike Measure 1560, where the age range cut off at 64). There was also concern that NCQA did not distinguish the fee-for-service versus the beneficiaries in Medicare Advantage plans.

2b. Overall Validity: H-4; M-5; L-0; I-0

2b1. Evidence is consistent with intent: H-8; M-1; L-0; I-0

TAP Discussion: The TAP believes the measure is clearly defined; however, one of the challenges will be the fact that COPD has multiple co-morbidities, particularly when compared to asthma. It will therefore be difficult to know if you are measuring exactly COPD. Specifications should be explored on how to develop disease severity; however, this is difficult to do with administrative datasets.

2b2. Score/Analysis: H-6; M-3; L-0; I-0

TAP Discussion: The TAP believes that overall the validity testing was appropriate. Outliers are identified by tagging O/E ratios below .3 or above 3.

2b3. Exclusions: H-4; M-5; L-0; I-0

TAP Discussion: The TAP agrees the exclusions are well stated and are similar to the asthma measure.

2b4. Risk Adjustment: H-6; M-3; L-0; I-0

TAP Discussion: Cardiovascular disease maybe a major driver of the severity of COPD.. The risk adjustment approach appears reasonable for the data available. The intent is to compare across populations.

2b5. Identification of statistically significant/meaningful differences: H-5; M-4; L-0; I-0

TAP Discussion: The TAP believes NCQA did a sufficient job presenting their data in a transparent manner.

2b6. Multiple data sources:

TAP Discussion: N/A (using all administrative data)

2c. Stratification for disparities: H-5; M-4; L-0; I-0

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

<p>TAP Discussion: Examining differences in racial disparities for this data set is not yet possible, but there is stratification by gender. Race is not a required field for most provider systems and is usually unavailable except in the Medicare population.</p>
<p>Overall Scientifically Acceptable: Yes [Y-13; N-1 (Committee Vote)] Overall Reliability: H-11; M-3; L-0; I-0 Overall Validity: H-4; M-10; L-0; I-0 Committee Discussion: The Steering Committee was satisfied by the appropriateness of the risk-adjustment methodology employed to address the multiple co-morbidities associated with COPD. They agreed with the TAP's assessment of Scientific Acceptability and raised no new concerns.</p>
<p>3. Usability: 3a. Measure performance results are publicly reported: H-9; M-0; L-0; I-0 TAP Discussion: The TAP believes this subcriterion has been met as NCQA does extensive audits of their material on a regular basis. 3b. Measure results are meaningful/useful for public reporting and quality improvement: H-5; M-4; L-0; I-0 TAP Discussion: The TAP feels the results are usable and understandable. 3c. Data and results can be decomposed for transparency and understanding: H-6; M-3; L-0; I-0 TAP Discussion: The TAP feels this subcriterion has been met as NCQA does extensive audits of their material on a regular basis, and the measure can be deconstructed to facilitate transparency. 3d. Harmonized or justification for differences: N/A</p>
<p>Overall Usability: H-7; M-7; L-0; I-0 Committee Discussion: The Steering Committee valued NCQA's rigorous auditing processes and the transparency with which the developers construct their measures. In addition to being used by health plans, the Committee acknowledged the usefulness of measures for purchasers/providers, giving them much more leverage during negotiations for their annual purchasing agreements.</p>
<p>4. Feasibility: 4a. Data elements routinely generated during care process: H-9; M-0; L-0; I-0 TAP Discussion: The TAP believes this subcriterion has been met as data is a byproduct of care. 4b. Data elements available electronically: H-9; M-0; L-0; I-0 TAP Discussion: The TAP believes this subcriterion has been met; all data is available electronically. 4c. Susceptibility to inaccuracies/ unintended consequences identified: H-6; M-3; L-0; I-0 TAP Discussion: The TAP believes this subcriterion has been met. 4d. Data collection strategy can be implemented: H-8; M-1; L-0; I-0 TAP Discussion: The TAP believes this subcriterion has been met.</p>
<p>Overall Feasibility: H-10; M-4; L-0; I-0 Committee Discussion: There were no new additional comments from the Steering Committee relating to feasibility of NCQA measures.</p>

868

<p>1611: ETG Based Pneumonia Cost of Care Measure (Ingenix)</p> <p>Description: The measure focuses on resources used to deliver episodes of care for patients with pneumonia. Pneumonia episodes are defined using the Episode Treatment Groups (ETG) methodology and describe the unique presence of the condition for a patient and the services involved in diagnosing, managing and treating pneumonia. A number of resource use measures are defined for pneumonia episodes, including overall cost of care, cost of care by type of service, and the utilization of specific types of services. Each resource use measure is expressed as a cost or a utilization count per episode and comparisons with internal and external benchmarks are made using risk adjustment to support valid comparisons. As requested by NQF, the focus of this submission is for pneumonia episodes and will cover both measures at the pneumonia base and severity level and also a pneumonia composite measure where pneumonia episode results are combined across pneumonia severity levels. At the most detailed level, the measure is defined as the base condition of pneumonia and an assigned level of severity (e.g., resources per episode for pneumonia, severity level 1 episodes). Composite measures can then be created using these measurement units to meet a specific need. For example, a composite measure for pneumonia is derived by combining pneumonia episode results across pneumonia severity levels. Appropriate risk adjustment is applied to support comparisons (e.g., for physician measurement, adjusting for a physician's mix of pneumonia episodes by severity level when supporting a pneumonia composite comparison). The focus of this measure is on pneumonia. However, pneumonia episode results could also be included in a "pulmonary" or other clinical composite for a physician, combining episodes in clinical areas similar to pneumonia. Further, an "overall" composite for a physician can be created, again by aggregating episode results</p>

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

<p>across appropriate conditions and severity levels and applying proper risk adjustment when making comparisons.</p> <p>Resource Use Type: Per episode</p> <p>Data Type: Administrative claims, Other</p> <p>Resource Use Service Categories: Inpatient services: Inpatient facility services, Inpatient services: Admissions/discharges, Ambulatory services: Outpatient facility services, Ambulatory services: Emergency Department, Ambulatory services: Pharmacy, Ambulatory services: Evaluation and management, Ambulatory services: Procedures and surgeries, Ambulatory services: Imaging and diagnostic, Ambulatory services: Lab services</p> <p>Care Setting: Ambulatory Care : Ambulatory Surgery Center (ASC), Ambulatory Care : Clinic/Urgent Care, Ambulatory Care : Clinician Office, Emergency Medical Services/Ambulance, Home Health, Hospice, Hospital/Acute Care Facility, Imaging Facility, Laboratory, Post Acute/Long Term Care Facility: Nursing Home/Skilled Nursing Facility, Post Acute/Long Term Care Facility: Rehabilitation</p> <p>Level of Analysis: Clinician : Group/Practice, Clinician : Individual, Clinician : Team, Facility, Health Plan, Integrated Delivery System, Population: Community, Population: County or City, Population : National, Population : Regional, Population: State</p> <p>Measure Developer: Ingenix, 950 Winter Street, Suite 3800, Waltham, Massachusetts, 02451</p>
<p>Committee Recommendation for Endorsement: Y-12; N-4; Abstain-0</p>
<p>Conditions/Questions for Developer:</p> <p>1. Would it be possible to break down the measure by bacterial versus non-bacterial to try to separate out pneumonia types?</p> <p>Developer Response:</p> <p>1. Yes, the measure is stratified. To the extent that administrative claims code the differences in pneumonia types, the measure can be stratified to evaluate resource use differences between pneumonia types.</p>
<p>1.Importance to Measure and Report</p> <p>1a.High Impact: H-8; M-0; L-0; I-0 <i>TAP Discussion:</i> The TAP agreed that pneumonia is a high impact and high cost area.</p> <p>1b. Resource use/cost problems: H-8; M-0; L-0; I-0 <i>TAP Discussion:</i> The TAP believes this subcriterion has been met.</p> <p>1c. Purpose clearly described: H-8; M-0; L-0; I-0 <i>TAP Discussion:</i> The TAP feel the purpose and objective are clear.</p> <p>1d. Resource use service categories consistent and representative: H-7; M-1; L-0; I-0 <i>TAP Discussion:</i> The TAP agrees the service categories are consistent and representative.</p>
<p>Overall Importance: Y-14, N-1 Committee Discussion: The Steering Committee deemed the measure to be important.</p>
<p>2.Scientific Acceptability of Measure Properties:</p> <p>2a. Overall Reliability: H-3; M-3; L-0; I-1</p> <p>2a1.Measure well defined and precisely specified: H-3; M-4; L-0; I-0 <i>TAP Discussion:</i> Several TAP members were uncomfortable with the lack of transparency in the risk adjustment specifications and felt that the severity weights, particularly for the elderly, were unclear. The panel also had a hard time identifying clean periods. There was a strong feeling that there should be some separation between community-acquired and healthcare-acquired pneumonia, as they represent very different clinical conditions.</p> <p>2a2. The results are repeatable: H-6; M-1; L-0; I-0 <i>TAP Discussion:</i> The TAP had concerns regarding the fact that that there is no way to ascertain how Ingenix came up with the specific weights assigned to comorbidities.</p> <p>2b. Overall Validity: H-0; M-7; L-0; I-0</p> <p>2b1. Evidence is consistent with intent: H-4; M-3; L-0; I-0 <i>TAP Discussion:</i> The panel again asked for clarification regarding why the measure has different weighted scores for the elderly.</p> <p>2b2.Score/Analysis: H-0; M-5; L-2; I-0 <i>TAP Discussion:</i> The TAP was concerned that they weren't provided with enough information to understand how Ingenix assigned risk scores. Questions regarding how diagnostic descriptions leads to increased utilization were raised. The TAP remained doubtful as to whether this measure should be counted as one distinct population.</p> <p>2b3. Exclusions: H-2; M-4; L-1; I-0 <i>TAP Discussion:</i> The TAP felt that more data around the impact of exclusions (e.g. sensitivity analysis) would be helpful. Ingenix confirmed that there are no clinical exclusions from the measure, only cost exclusions.</p> <p>2b4. Risk Adjustment: H-1; M-3; L-2; I-1 <i>TAP Discussion:</i> The TAP believed that the risk-adjustment methodology is not readily transparent. More information on how risk</p>

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

<p>scores are assigned was requested from the developers.</p> <p>2b5. Identification of statistically significant/meaningful differences: H-0; M-7; L-0; I-0 <i>TAP Discussion:</i> Data submitted does demonstrate variation in resource use. However, there was a general feeling that meaningfulness is questionable since types of pneumonia cannot be separated out.</p> <p>2b6. Multiple data sources: N/A (using all administrative data)</p> <p>2c. Stratification for disparities: H-2; M-5; L-0; I-0 <i>TAP Discussion:</i> Gender and age can be stratified, but race data is not available in administrative claims.</p>
<p>Overall Scientifically Acceptable: Yes [Y-13; N-3 (Committee Vote)] Overall Reliability: H-3; M-11; L-2; I-0 Overall Validity: H-1; M-13; L-2; I-0 Committee Discussion: The Steering Committee agreed that this measure would not be clinically relevant at the physician level due to its limited ability to differentiate between community and hospital acquired pneumonia. In general, the Committee also believed that the "start and stop rules" would be more readily apparent for acute procedure-oriented measures such as knee replacements, as compared with chronic illnesses, which has less clear cut start and stop dates. The Committee reiterated the TAP's concern that Ingenix specified the measure for use in patients over 65 using commercial data to calibrate the model. Commercial patients over 65 are not representative of the general over 65 population.</p>
<p>Usability:</p> <p>3a. Measure performance results are publicly reported: H-0; M-6; L-1; I-0 <i>TAP Discussion:</i> The TAP agrees that despite the fact that multiple care organizations are currently using this measure, the inability to distinguishing between types of pneumonia severely limits the usability of the measure. They concurred that for individual organizations this limitation might be acceptable, but the measure wouldn't be useful as a national consensus standard. . NQF clarified that the measure has a specified for particular levels of analysis, and the ratings need to be reflective of that specification.</p> <p>3b. Measure results are meaningful/useful for public reporting and quality improvement: H-1; M-5; L-1; I-0 <i>TAP Discussion:</i> The TAP agrees that this subcriterion has been met.</p> <p>3c. Data and results can be decomposed for transparency and understanding: H-1; M-5; L-1; I-0 <i>TAP Discussion:</i> The TAP feels the measure would be more transparent if more user-friendly detail were provided.</p> <p>3d. Harmonized or justification for differences: N/A</p>
<p>Overall Usability: H-3; M-11; L-1; I-1 Committee Discussion: There were no additional concerns identified by the Steering Committee for this criterion.</p>
<p>4. Feasibility:</p> <p>4a. Data elements routinely generated during care process: H-7; M-0; L-0; I-0 <i>TAP Discussion:</i> The TAP believes this subcriterion has been met; data is a byproduct of care.</p> <p>4b. Data elements available electronically: H-7; M-0; L-0; I-0 <i>TAP Discussion:</i> The TAP believes this subcriterion has been met; data available electronically.</p> <p>4c. Susceptibility to inaccuracies/ unintended consequences identified: H-1; M-5; L-0; I-1 <i>TAP Discussion:</i> The TAP concluded there was a lack of information in the submission regarding data cleaning and missing data to sufficiently understand those areas.</p> <p>4d. Data collection strategy can be implemented: H-5; M-2; L-0; I-0 <i>TAP Discussion:</i> The TAP agrees this subcriterion has been met.</p>
<p>Overall Feasibility: H-1; M-8; L-7; I-0 Committee Discussion: See Ingenix feasibility discussion above.</p>

869

<p>1609: ETG/PEG Based hip/knee replacement Cost of Care Measure (Ingenix)</p> <p>Description: The measure focuses on resources used to deliver episodes of care for patients who have undergone a Hip/Knee Replacement. Hip Replacement and Knee Replacement episodes are initially defined using the Episode Treatment Groups (ETG) methodology and describe the unique presence of the condition for a patient and the services involved in diagnosing, managing and treating the condition. The Procedure Episode Group (PEG) methodology uses the ETG results and further logic to creating a procedure episode that focuses on the Hip Replacement and Knee Replacement component of the care. Procedure episodes identify a unique procedure event as well as the related services performed before and after the procedure including workup and therapy prior to the procedure as well as post-op activities such as repeated surgery and patient follow-up. Together, the ETG and</p>

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

<p>PEG methodologies identify the services involved in diagnosing, managing and treating patients with Hip/Knee Replacements. A methodology to assign a severity level to each episode is employed to group Hip and Knee Replacement episodes by level of risk.</p> <p>Resource Use Type: Per episode</p> <p>Data Type: Administrative claims</p> <p>Resource Use Service Category: Inpatient services: Inpatient facility services, Inpatient services: Admissions/discharges, Ambulatory services: Outpatient facility services, Ambulatory services: Emergency Department, Ambulatory services: Pharmacy, Ambulatory services: Evaluation and management, Ambulatory services: Procedures and surgeries, Ambulatory services: Imaging and diagnostic, Ambulatory services: Lab services</p> <p>Care Setting: Ambulatory Care : Ambulatory Surgery Center (ASC), Ambulatory Care : Clinic/Urgent Care, Ambulatory Care : Clinician Office, Emergency Medical Services/Ambulance, Home Health, Hospice, Hospital/Acute Care Facility, Imaging Facility, Laboratory Post Acute/Long Term Care Facility : Nursing Home/Skilled Nursing Facility, Post Acute/Long Term Care Facility : Rehabilitation</p> <p>Level of Analysis: Clinician : Group/Practice, Clinician : Individual, Clinician : Team, Facility, Health Plan, Integrated Delivery System, Population : Community, Population : County or City, Population : National, Population : Regional, Population : State</p> <p>Measure Developer: Ingenix, 950 Winter Street, Suite 3800, Waltham, Massachusetts, 02451</p>
<p>Committee Recommendation for Endorsement: Y-9; N-7; Abstain-0</p>
<p>Conditions/Questions for Developer: N/A</p> <p>Developer Response: N/A</p>
<p>1.Importance to Measure and Report –</p> <p>1a.High Impact: H-6; M-1; L-0; I-0</p> <p>TAP Discussion: The TAP deemed this measure to be a high cost/high impact area.</p> <p>1b. Resource use/cost problems: H-0; M-2; L-5; I-0</p> <p>TAP Discussion: The TAP felt that the measure would be able to identify large variation in resource use and cost. However, the TAP felt that the developers could have provided more information specifically related to hip/knee replacement variation in resource use in the measure submission.</p> <p>1c. Purpose clearly described: H-0; M-5; L-1; I-1</p> <p>TAP Discussion: The TAP felt that the purpose was sufficiently described.</p> <p>1d. Resource use service categories consistent and representative: H-2; M-5; L-0; I-0</p> <p>TAP Discussion: The TAP felt that the resource use service categories were appropriate.</p>
<p>Overall Importance: Y-17, N-0</p> <p>Committee Discussion: The Steering Committee deemed this measure to be important.</p>
<p>2.Scientific Acceptability of Measure Properties:</p> <p>2a. Reliability:</p> <p>2a1.Measure well defined and precisely specified: H-0; M-3; L-4; I-0</p> <p>TAP Discussion: The TAP wanted more information on how the developers handled right and left hip/knee replacement since there is limited ability to distinguish between right/left surgery in the administrative data used. It is important to capture the rate of surgery at the provider level to ensure that the current measure construct does not penalize those providers who chose conservative treatment for low severity patients. The developer should provide more clear information on the clinical logic, including the specific codes that are used to create the episodes. Overall, the TAP wanted more clarity on the clinical construction logic of the episode such as severity level assignments, assignment of claims with two concurrent episodes (i.e. tie breaking logic). The TAP also wanted more information on the procedure definitions, handling of comorbidities and the weighting of multiple co-occurring comorbidities.</p> <p>2a2. The results are repeatable: H-2; M-5; L-0; I-0</p> <p>TAP Discussion: The TAP wanted additional information on how reliable the physician level scores were over time.</p> <p>Overall Reliability: H-2; M-4; L-0; I-0</p> <p>TAP Discussion:</p> <p>2b. Validity</p> <p>2b1. Evidence is consistent with intent: H-2; M-4; L-1; I-0</p> <p>TAP Discussion: The TAP felt that the evidence was consistent with the intent of the measure.</p> <p>2b2.Score/Analysis: H-1; M-4; L-2; I-0</p> <p>TAP Discussion: The TAP discussed the attribution of costs six months before the procedure as too long of a period for a physician based measure. With the current attribution method, it appears to be more appropriate at a plan or system-level rather than an individual provider. These attribution approaches were submitted as guidelines only.</p> <p>2b3. Exclusions: H-0; M-2; L-4; I-1</p>

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

<p>TAP Discussion: The TAP wanted more information on why low cost outliers were excluded and high cost outliers were winsorized; a sensitivity analysis of this decision was recommended by the TAP. The TAP also recommended that the measure should include a count of high cost outliers if they are going to be winsorized. Information about the high cost outliers might actually drive targeted interventions.</p> <p>2b4. Risk Adjustment: H-0; M-0; L-6; I-1</p> <p>TAP Discussion: The TAP wanted more information on severity levels on how they related to the risk adjustment model. The TAP agreed that not all of the comorbidities provided in the submission seem appropriate for the population in the measure.</p> <p>2b5. Identification of statistically significant/meaningful differences:</p> <p>TAP Discussion: There was general agreement that the complexities of the score may make it difficult to discern meaningful differences between providers.</p> <p>2b6. Multiple data sources: N/A</p> <p>Overall Validity: H-0; M-1; L-5; I-0</p> <p>2c. Stratification for disparities: H-1; M-0; L-4; I-2</p> <p>TAP Discussion: Administrative data is limited in its ability to stratify based on race.</p>
<p>Overall Scientifically Acceptable: Yes [Y-11; N-5 (Committee Vote)]</p> <p>Overall Reliability: H-2; M-14; L-0; I-0</p> <p>Overall Validity: H-1; M-9; L-6; I-0</p> <p>Committee Discussion: The Steering Committee was concerned with the lack of specification regarding the measure's use of MSDRG's in the risk-adjustment methodology. Ingenix explained that among the population of patients who undergo knee or hip replacements, there is minimal variation in the underlying co-morbidities. Therefore, the methodology required to adequately risk adjust is much less stringent than it would be if looking at a more complicated condition such as coronary artery disease.</p>
<p>3. Usability:</p> <p>3a. Measure performance results are publicly reported: H-0; M-5; L-2; I-0</p> <p>TAP Discussion: The TAP was concerned that this ETG was not currently being used as a stand-alone measure and it was unclear if it was currently being publicly reported.</p> <p>3b. Measure results are meaningful/useful for public reporting and quality improvement: H-0; M-4; L-3; I-0</p> <p>TAP Discussion: The TAP was concerned that this ETG was not currently being used as a stand-alone measure which may impact the need for public reporting.</p> <p>3c. Data and results can be decomposed for transparency and understanding: H-0; M-3; L-4; I-0</p> <p>TAP Discussion: The TAP expressed concern over the difficulty in understanding the clinical hierarchy and risk model. The lack of clarity in these aspects of the measure makes it difficult to deconstruct the measure for transparency and understanding.</p> <p>3d. Harmonized or justification for differences: N/A</p>
<p>Overall Usability: H-0; M-12; L-4; I-1</p> <p>Committee Discussion: The Steering Committee iterated their concern that, because the measure is used as part of a grouper, it is unclear if it is useful as a standalone measure. Additionally, based on the nature of the Ingenix product, hip and knee replacements had been combined into a single measure, which was not believed by some to be the most clinically relevant approach.</p>
<p>4. Feasibility:</p> <p>4a. Data elements routinely generated during care process: H-5; M-2; L-0; I-0</p> <p>TAP Discussion: The TAP believes this subcriterion has been met; data is a byproduct of care.</p> <p>4b. Data elements available electronically: H-6; M-1; L-0; I-0</p> <p>TAP Discussion: The TAP believes this subcriterion has been met; data elements that are available electronically.</p> <p>4c. Susceptibility to inaccuracies/ unintended consequences identified: H-0; M-3; L-4; I-0</p> <p>TAP Discussion: The TAP agrees that much of this surgery is dependent on patient preferences thus the measure should account for these preferences in inclusion and exclusion criteria of the measure. Additionally, providers who treat their patients conservatively can appear to be high cost users since the only patients who get surgery are those who are more severe.</p> <p>4d. Data collection strategy can be implemented: H-1; M-5; L-1; I-0</p> <p>TAP Discussion: No additional issues were raised by the TAP.</p>
<p>Overall Feasibility: H-1; M-8; L-7; I-0</p> <p>Committee Discussion: See Ingenix feasibility discussion above.</p>

870
871

NATIONAL QUALITY FORUM

872 **Candidate Consensus Standards Not Recommended for Endorsement**

873 Six candidate consensus standards were not recommended for endorsement because they did not
874 meet NQF criteria; two did not pass scientific acceptability, and the remaining had issues with
875 other criteria.

876 The evaluation summary tables follow the list of measures and summarize the results of the
877 TAP’s and Committee’s evaluation of and voting on the candidate consensus standards not
878 recommended for endorsement. Hyperlinks are provided from each summary table to the
879 detailed measure specifications. To access the meeting transcripts and recordings in which these
880 measures are discussed, refer to the [project web page](#).

881

882 ***Cardiovascular***

883 (1591) ETG-based congestive heart failure (CHF) cost of care measure (Ingenix).....42

884 (1594) ETG-based coronary artery disease (CAD) cost of care measure (Ingenix)45

885 ***Non-Condition Specific***

886 (1599) ETG-based non-condition specific cost of care measure (Ingenix)48

887 ***Bone/Joint***

888 (1603) ETG-based hip fracture cost of care measure (Ingenix)51

889 ***Pulmonary***

890 (1605) ETG-based asthma cost of care measure (Ingenix)53

891 (1608) ETG-based chronic obstructive pulmonary disease (COPD) cost of care measure
892 (Ingenix)56

893 **Evaluation Summary—Candidate Consensus Standards Not Recommended for** 894 **Endorsement**

[1591: ETG Based Congestive Heart Failure \(CHF\) cost of care measure \(Ingenix\)](#)

Description: The measure focuses on resources used to deliver episodes of care for patients with Congestive Heart Failure (CHF). CHF episodes are defined using the Episode Treatment Groups (ETG) methodology and describe the unique presence of the condition for a patient and the services involved in diagnosing, managing and treating CHF. A number of resource use measures are defined for CHF episodes, including overall cost of care, cost of care by type of service, and the utilization of specific types of services. Each resource use measure is expressed as a cost or a utilization count per episode and comparisons with internal and external benchmarks are made using risk adjustment to support valid comparisons. As requested by NQF, the focus of this submission is for CHF episodes and will cover both measures at the CHF base and severity level and also a CHF composite measure where CHF episode results are combined across CHF severity levels. At the most detailed level, the measure is defined as the base condition of CHF and an assigned level of

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

severity (e.g., resources per episode for CHF, severity level 1 episodes). Composite measures can then be created using these measurement units to meet a specific need. For example, a composite measure for CHF is derived by combining CHF episode results across CHF severity levels. Appropriate risk adjustment is applied to support comparisons (e.g., for physician measurement, adjusting for a physician's mix of CHF episodes by severity level when supporting a CHF composite comparison). The focus of this measure is on CHF. However, CHF episode results could also be included in a "cardiology", "chronic care", or other clinical composite for a physician, combining episodes in clinical areas similar to CHF. Further, an "overall" composite for a physician can be created, again by aggregating episode results across appropriate conditions and severity levels and applying proper risk adjustment when making comparisons.

Resource Use Type: Per episode

Data Type: Administrative claims, other

Resource Use Service Category: Inpatient services: Inpatient facility services, Inpatient services: Admissions/discharges, Ambulatory services: Outpatient facility services, Ambulatory services: Emergency Department, Ambulatory services: Pharmacy, Ambulatory services: Evaluation and management, Ambulatory services: Procedures and surgeries, Ambulatory services: Imaging and diagnostic
Ambulatory services: Lab services

Care Setting: Ambulatory Care : Ambulatory Surgery Center (ASC), Ambulatory Care : Clinic/Urgent Care, Ambulatory Care : Clinician Office, Emergency Medical Services, Ambulance, Home Health, Hospice, Hospital/Acute Care Facility, Imaging Facility, Laboratory

Level of Analysis: Clinician : Group/Practice, Clinician : Individual, Clinician : Team, Facility, Health Plan, Integrated Delivery System
Population : Community, Population : County or City, Population : National, Population : Regional, Population : states

Measure Developer: Ingenix, 950 Winter Street, Suite 3800, Waltham, Massachusetts, 02451

Committee Recommendation for Endorsement: Y-6; N-8; Abstain-0 (re-vote) [Y-10; N-8; Abstain-0 (initial vote)]

Conditions/Questions for Developer:

1. Why are some of the codes, typically seen in congestive heart failure measures, excluded?
2. How are hospitalizations that occur during the course of the measure handled?
3. Does the episode include events that occur before and/or after the episode?

Developer Response:

1. Ingenix excluded the codes that were specific to diastolic heart failure (as this is a systolic and diastolic/systolic mix measure); if those codes were included it would have created another episode. Ingenix includes codes that were both systolic and diastolic, and used them as a marker to increase the severity score for the episode.
2. Hospital admissions that occurred during the course of the measure that are coded for congestive heart failure are included in the measure; hospitalizations are not used for severity adjustment. If the hospital admission date occurs during the measurement year, then the admission is included in that measurement year.
3. No, this measure is insulated from events that occur before or after the episode.

1. Importance to Measure and Report

1a. High Impact: H -8; M-0; L-0; I-0

TAP Discussion: The TAP believes this is a high impact, high cost area that is important to measure and report.

1b. Resource use/cost problems: H -8; M-0; L-0; I-0

TAP Discussion: The TAP believes this subcriterion has been met.

1c. Purpose clearly described: H -5; M-3; L-0; I-0

TAP Discussion: The TAP believes the purpose of the measure is clearly described.

1d. Resource use service categories consistent and representative: H -7; M-1; L-0; I-0

TAP Discussion: The TAP believes this subcriterion has been met; the resource use service categories are consistent and representative of the measure.

Overall Importance: Yes [Y-17; N-1 (Committee Vote)]

Committee Discussion: The Steering Committee believes this is a high impact, high cost area and that the measure has been clearly described. This criterion has been met.

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

2. Scientific Acceptability of Measure Properties:

2a. Reliability:

2a1. Well defined/precise specifications: H -3; M-4; L-0; I-1

TAP Discussion: The TAP believed there was a bit of confusion around the term, “congestive heart failure”, it was brought up that not all “heart failure” is necessarily “congestive” and there needs to be more clarification around the use of this term. The TAP agrees that this measure is targeting systolic heart failure and then a mix of systolic/diastolic heart failure. Ingenix also has a diastolic heart failure measure, but it has not been submitted for NQF endorsement. When the ICD9 code exists for systolic and diastolic – it’s a marker for severity adjustment. Overall, the TAP believes that the clinical and construction logic of the measure was described in sufficient detail and users will be able to implement the measure as described.

2a2. Reliability testing: H -7; M-1;L-0;I-0

TAP Discussion: The TAP believes this measure has demonstrated extensive benchmarking and comparisons; however they would have liked to see more external comparisons. The testing data submitted was from nine health care organizations, all large commercial insurers that vary geographically. Ingenix demonstrated reliability by performing parallel development of the data by using two independent approaches. These two different approaches led to the same results as levels near 99.9%The data was tested primarily on commercial databases, however some Part C plan Medicare patients were also included. It is important to note that this measure was submitted for use in the commercial, less than 65 years old population.

2b. Validity:

2b1. Specifications consistent with resource use/cost problem: H -2;M-2; L-0; I-0

TAP Discussion: The TAP agrees that the specifications are consistent with the resource use.

2b2. Validity testing: H -4; M-4; L-0; I-0

TAP Discussion: The TAP believes Ingenix has sufficiently demonstrated face validity.

2b3. Exclusions: H -4; M-3; L-1; I-0

TAP Discussion: There are no exclusions within this measures, the TAP believes this subcriterion has been met.

2b4. Risk adjustment: H -4; M-2; L-0; I-1

TAP Discussion: The TAP believes that this risk adjustment appears to be somewhat circular – the measure is risk adjusted if the individual was hospitalized during the year – if the provider is using a large amount of resources, inevitably there will be more diagnoses in that measurement period, which would in turn also affect severity level category. Ingenix has made it clear that they are not using utilization to directly risk-adjust the cost of the episode. There is a lack of information in terms of the variables selected for inclusion in the calibration of the risk model, the risk groups selected in terms of a cutoff for the severity score, and there is no rationale presented for why this cutoff point has been chosen.

2b5. Identification of statistically significant/meaningful differences: H -2; M-1; L-3; I-1

TAP Discussion: The TAP believes there is little information to compare statistical versus practical significance for this measure. The measure allows the user to determine what is clinically significant based on confidence intervals. The sample size appears sufficient enough to obtain a confidence interval that it will be useful to establish differences that are clinically and statistically significant. Ingenix has created confidence intervals around the observed to expected ratio The minimum sample size to detect statistically significant differences depends upon the case mix of the providers and the variation in performance across providers..

2b6. Multiple data sources: N/A

2c. Stratification for disparities: H-0; M-0; L-0; I-0; N/A-8

TAP Discussion: Due to the limitations in the administrative claims data, at this time the measure does not stratify for disparities.

Overall Reliability: H-3, M-12, L-2, I-0

Overall Validity: H-1, M-13, L-4, I-0

Overall Scientific Acceptability: Yes [Y-14; N-4(Committee Vote)]

Committee Discussion: The Steering Committee discussion focused on how clearly specified the codes used with the measure are, and how well they capture systolic heart failure. This is a measure of systolic heart failure, a paired measure of diastolic heart failure from Ingenix exists but they did not submit it to the project. Because the Steering Committee could not take into account the existence of the diastolic measure, there was concern around the completeness and accuracy with which this measure would capture systolic heart failure. The diagnosis codes specified are limited to the 428 codes that used the word “systolic”, they do not use some of the 404s and 402s that the other measures have used to capture the larger heart failure population. The measure specifications have been in use for a significant amount of time; Ingenix has demonstrated that if this measure is used in the same population, at the same time, then the result will be the same roughly 99.9% of the time. The Steering Committee discussed how there are carve outs for mental health & pharmacy data and therefore comparisons within the health plan are the same or likely to be the same. However, when comparing

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

<p>across health plans or across physician groups validity may become an issue when there are differences in the completeness of the data submitted. The Steering Committee expressed concerns over the reliability, validity and risk adjustment method. Specifically, that the measure may be adjusting for comorbidities identified during the measurement period as opposed to comorbidities identified prior to the episode. There was also concern that the risk adjustment may be “over –adjusting”, or possibly “adjusting away” significant differences.</p>
<p>3. Usability:</p> <p>3a. Measure performance results are publicly reported: H-1; M-1; L-2;I-2 <i>TAP Discussion:</i> The TAP was concerned with the availability of this data to the public and requested clarification from NQF on what is required for “public reporting”. The measures are widely used by providers to compare to one another. The results of this measure also allow for provider profiling, provider report cards and there is a cost base analysis for the members to estimate what the cost of the service would be, including the out of pocket expense. Since this measure is reported within a suite of measures, it has not been broken out individually for reporting or use in quality improvement.</p> <p>3b. Measure results are meaningful/useful for public reporting and quality improvement: H-3; M-1; L-0; I-2 <i>TAP Discussion:</i> The TAP agrees that more information would be needed to explain the results of this measure to the public and to be used for internal quality improvement.</p> <p>3c. Data and results can be decomposed for transparency and understanding: H-0; M-2; L-3;I-1 <i>TAP Discussion:</i> The TAP agrees there are challenges for the use of this measure, which include its complexity and lack of clarity in the specifications. TAP also agrees it is difficult to assess the extent to which the measure can be decomposed as it is currently specified.</p> <p>3d. Harmonized or justification for differences: N/A</p>
<p>Overall Usability: H-0; M-10; L-7; I-0, N/A-0 Committee Discussion: The Steering Committee discussed the fact that more information would be needed to explain the results of this measure to the public and to be used for internal quality improvement. The Steering Committee believes there are challenges for the use of this measure, which include its complexity and lack of clarity in the specifications. The Steering Committee agrees it is difficult to assess the extent of which the measure can be decomposed as it is currently specified.</p>
<p>4. Feasibility:</p> <p>4a. Data elements routinely generated during care process: H-5; M-0; L-0; I-1 <i>TAP Discussion:</i> The TAP believes that this sub criterion has been met; all of the data elements are generated during the care process.</p> <p>4b. Data elements available electronically: H-5;M-0; L-0;I-1 <i>TAP Discussion:</i> The TAP believes that this sub criterion has been met; all of the data is available electronically.</p> <p>4c. Susceptibility to inaccuracies/ unintended consequences identified: H-0; M-4; L-1; I-1 <i>TAP Discussion:</i> The TAP noted that Ingenix does not have a formal audit system to ensure that all of the numbers are included & correct. In general, when dealing with any measure that uses administrative data there are various inaccuracies, pertaining particularly to coding inaccuracies and variation.</p> <p>4d. Data collection strategy can be implemented: H-3; M-0; L-1; I-2 <i>TAP Discussion:</i> The majority of the TAP agreed that barriers to use are minimal. (NQF Note: This is prior to the submission of product pricing information shared only with the Steering Committee)</p>
<p>Overall Feasibility: H-2; M-8; L-7; I-1 Committee Discussion: See Ingenix feasibility discussion above.</p>

895

<p>1594 ETG Based Coronary Artery Disease (CAD) cost of care measure (Ingenix)</p> <p>Description: The measure focuses on resources used to deliver episodes of care for patients with CAD. CAD episodes are defined using the Episode Treatment Groups (ETG) methodology and describe the unique presence of the condition for a patient and the services involved in diagnosing, managing and treating CAD. A number of resource use measures are defined for CAD episodes, including overall cost of care, cost of care by type of service, and the utilization of specific types of services. Each resource use measure is expressed as a cost or a utilization count per episode and comparisons with internal and external benchmarks are made using risk adjustment to support valid comparisons. As requested by NQF, the focus of this submission is for CAD episodes and will cover both measures at the CAD base and severity level and also a CAD composite measure where CAD episode results are combined across CAD severity levels. At the most detailed level, the measure is defined as the base condition of CAD and an assigned level of severity (e.g., resources per episode for CAD, severity level 1 episodes). Composite measures can then be created using these measurement units to meet a specific need. For example, a composite measure for CAD is derived by combining CAD episode results across CAD</p>

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

severity levels. Appropriate risk adjustment is applied to support comparisons (e.g., for physician measurement, adjusting for a physician's mix of CAD episodes by severity level when supporting a CAD composite comparison). The focus of this measure is on CAD. However, CAD episode results could also be included in a "cardiology", "chronic care", or other clinical composite for a physician, combining episodes in clinical areas similar to CAD. Further, an "overall" composite for a physician can be created, again by aggregating episode results across appropriate conditions and severity levels and applying proper risk adjustment when making comparisons.

Resource Use Type: Per episode

Data Type: Administrative claims, other

Resource Use Service Category: Inpatient services: Inpatient facility services, Inpatient services: Admissions/discharges, Ambulatory services: Outpatient facility services, Ambulatory services: Emergency Department, Ambulatory services: Pharmacy, Ambulatory services: Evaluation and management, Ambulatory services: Procedures and surgeries, Ambulatory services: Imaging and diagnostic
Ambulatory services: Lab services

Care Setting: Ambulatory Care : Ambulatory Surgery Center (ASC), Ambulatory Care : Clinic/Urgent Care, Ambulatory Care : Clinician Office, Emergency Medical Services, Ambulance, Home Health, Hospice, Hospital/Acute Care Facility, Imaging Facility
Laboratory

Level of Analysis: Clinician : Group/Practice, Clinician : Individual, Clinician : Team , Facility, Health Plan, Integrated Delivery System, Population : Community, Population : County or City, Population : National, Population : Regional, Population : states

Measure Developer: Ingenix, 950 Winter Street, Waltham, Massachusetts, 02154

Committee Recommendation for Endorsement: Y-5; N-9; Abstain – 0 (re-vote) [Y-8; N-10; Abstain-0 (initial vote)]

1. Importance to Measure and Report

1a. High Impact: H-5; M-0; L-0; I-0

TAP Discussion: The TAP believes this is a high impact, high cost area; this sub criterion has been met.

1b. Resource use/cost problems: H-5; M-0; L-0; I-0

TAP Discussion: The TAP believes this subcriterion has been met.

1c. Purpose clearly described: H-5; M-0; L-0; I-0

TAP Discussion: The TAP believes this subcriterion has been met; the measure purpose is clearly described.

1d. Resource use service categories consistent and representative: H-3; M-2; L-0; I-0

TAP Discussion: The TAP believes this subcriterion has been met; the resource use categories are consistent and representative.

Overall Importance: Y-16, N-1 (Committee Vote)

Committee Discussion: The Steering Committee believes this is a high impact, high cost area and that the measure has been clearly described. This criterion has been met.

NATIONAL QUALITY FORUM

2. Scientific Acceptability of Measure Properties:

2a. Reliability:

2a1. Well defined/precise specifications: H-3; M-1; L-0; I-0

TAP Discussion: The diagnoses codes for this measure are the 410s through 414s and then the 429s, all of which represent complications of myocardial infarction. These codes seem comprehensive for identifying patients with coronary artery disease; however, the Steering Committee raised the question if the populations are similar enough that the user can reasonably make inferences about the resource use needed for each type of cardiac episode. Overall, the measure is very well specified and is being used across different health plans.

2a2. Reliability testing: H-3; M-1; L-0; I-0

TAP Discussion: The measure is specified in a way that it has been used over a long period of time, Ingenix demonstrated that if the user uses the same measure in the same population then the result will be the same. The TAP believes this subcriterion has been met.

2b. Validity:

2b1. Specifications consistent with resource use/cost problem: H-3; M-1; L-0; I-0

TAP Discussion: The TAP believes this subcriterion has been met; a specific population is defined and measured.

2b2. Validity testing: H-3; M-0; L-0; I-0

TAP Discussion: The TAP believes Ingenix has sufficiently demonstrated face validity.

2b3. Exclusions: H-2; M-1; L-0; I-0

TAP Discussion: There are no exclusions within this measures, the TAP believes this subcriterion has been met.

2b4. Risk adjustment: H-2; M-1; L-0; I-0

TAP Discussion: The TAP requested that the developer demonstrate proof of the concept that this is accurately accounting for differences in the population – the risk adjustment method does not appear to be robust. Additional information the model's goodness of fit was requested. NQF staff is working with Ingenix to supply this information to the Steering Committee.

2b5. Identification of statistically significant/meaningful differences: H-1; M-0; L-1; I-1

TAP Discussion: The Steering Committee believes that this measure did not identify statistically significant or meaningful differences across groups. There was general concern that something may be classified as statistically significant, when it is not clinically significant.

2b6. Multiple data sources: N/A

TAP Discussion: N/A

2c. Stratification for disparities: H-0; M-0; L-0; I-0; N/A-8

TAP Discussion: Due to the limitations in the administrative claims data, at this time the measure does not stratify for disparities.

Overall Reliability : H-5; M-11; L-2; I-0

Overall Validity: H-2; M-10; L-6; I-0

Overall Scientifically Acceptable: Yes [Y-12; N- 5 (Committee vote)]

Committee Discussion: The Steering Committee agreed that the measure accurately identified the primary incurring diagnoses codes as 410s through 414s. Within those strata there is a range of conditions – ranging from chronic, stable coronary artery disease to patients with cardiogenic shock complicated by a flail mitral posterior leaflet. The Steering Committee discussed how there is a large spectrum of risk adverse outcomes within this population. Furthermore, this carries the risk of different resource use for each specific condition included in the measure. The measure was submitted for implementation across various levels of analysis, however for individual clinicians there is not a sample size guideline. Regarding specific reliability testing, the measure is specified in a way that it has been used over a long period of time. The Steering Committee discussed how there are carve outs for mental health & pharmacy data and therefore comparisons within the health plan are the same or likely to be the same. However, when comparing across health plans or across physician groups validity may become an issue. There were concerns around the risk adjustment method. Specifically, the Committee was concerned that the measure may be adjusting for comorbidities identified during the measurement episode as opposed to comorbidities identified prior to the episode. There was also concern that the risk adjustment may be “over –adjusting”, or possibly “adjusting away” significant differences.

3. Usability:

3a. Measure performance results are publicly reported: H-0; M-1; L-1; I-1

TAP Discussion: The TAP was concerned with the availability of this data to the public and requested clarification from NQF on what is required for “public reporting”. The measures are widely used by providers to compare to one another. The results of this measure also allow for provider profiling, provider report cards and there is a cost base analysis for the members to estimate what the cost of the service would be, including the out of pocket expense. Since this measure is reported within a suite of measures, it has not been broken out individually for reporting or use in quality improvement.

3b. Measure results are meaningful/useful for public reporting and quality improvement: H-0; M-2; L-1; I-0

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

<p>TAP Discussion: The TAP agrees that more information would be needed to explain the results of this measure to the public and to be used for internal quality improvement.</p> <p>3c. Data and results can be decomposed for transparency and understanding: H-0; M-3;L-0;I-0</p> <p>TAP Discussion: The TAP agreed there are challenges for the use of this measure, which include its complexity and lack of clarity in the specifications. TAP also agreed it is difficult to assess the extent of which the measure can be deconstructed for understanding as it is currently specified.</p> <p>3d. Harmonized or justification for differences: N/A</p>
<p>Overall Usability: H-1; M-11; L-4; I-1</p> <p>Committee Discussion: The Steering Committee agrees that more information would be needed to explain the results of this measure to the public and to be used for internal quality improvement. The Steering Committee discussed the challenges for the use of this measure, which include its complexity and lack of clarity in the specifications. The Steering Committee agrees it is difficult to assess the extent of which the measure can be decomposed as it is currently specified.</p>
<p>4. Feasibility:</p> <p>4a. Data elements routinely generated during care process: H-3; M-0; L-0; I-0</p> <p>TAP Discussion: The TAP believes that this sub criterion has been met; all of the data elements are generated during the care process.</p> <p>4b. Data elements available electronically: H-3; M-0; L-0;I-0</p> <p>TAP Discussion: The TAP believes that this sub criterion has been met; all of the data is available electronically.</p> <p>4c. Susceptibility to inaccuracies/ unintended consequences identified: H-2; M-1; L-0; I-0</p> <p>TAP Discussion: The TAP noted that Ingenix does not have a formal audit system to ensure that all of the numbers are included & correct. In general, when dealing with any measure that uses administrative data there are various inaccuracies, pertaining particularly to coding inaccuracies and variation.</p> <p>4d. Data collection strategy can be implemented: H-2;M-0; L-1; I-0</p> <p>TAP Discussion: The majority of the TAP agreed that barriers to use are minimal. (NQF Note: This is prior to the submission of product pricing information shared only with the Steering Committee)</p>
<p>Overall Feasibility: H-3; M-8; L-6; I-1</p> <p>Committee Discussion: See Ingenix feasibility discussion above.</p>

896

<p>1599: ETG Based Non-Condition Specific cost of care measure (Ingenix)</p> <p>Description: The measure focuses on resources used to diagnose, manage and treat a population of patients (non-condition specific) during a defined 12-month period of time. The population included in the measurement can be described generally. Examples include a population of individuals enrolled with a health plan, individuals assigned to a patient-centered medical home or accountable care organization (ACO), or a panel of individuals managed by a primary care physician (PCP). A number of resource use measures are defined for this measure set, including overall cost of care, cost of care by type of service, and the utilization of specific types of services. Each resource use measure is expressed as a cost or a utilization count per member per month and comparisons with internal and external benchmarks are made using risk adjustment to support valid comparisons. Risk adjustment is based on the measure of risk assigned to each individual using the Episode Risk Group (ERG) methodology.</p> <p>Resource Use Type: Per capita (population- or patient-based)</p> <p>Data Type: Administrative claims</p> <p>Resource Use Service Category: Inpatient services: Inpatient facility services, Inpatient services: Admissions/discharges, Ambulatory services: Outpatient facility services, Ambulatory services: Emergency Department, Ambulatory services: Pharmacy, Ambulatory services: Evaluation and management, Ambulatory services: Procedures and surgeries, Ambulatory services: Imaging and diagnostic Ambulatory services: Lab services</p> <p>Care Setting: Ambulatory Care : Ambulatory Surgery Center (ASC), Ambulatory Care : Clinic/Urgent Care, Ambulatory Care : Clinician Office, Emergency Medical Services, Ambulance, Home Health, Hospice, Hospital/Acute Care Facility, Imaging Facility, Laboratory</p> <p>Level of Analysis: Clinician : Group/Practice, Clinician : Individual, Clinician : Team, Facility, Health Plan, Integrated Delivery System</p> <p>Population : County or City, Population : National, Population : Regional, Population : states</p> <p>Measure Developer: Ingenix, 950 Winter Street, Suite 3800, Waltham, Massachusetts, 02451</p> <p>Committee Recommendation for Endorsement: Y-5; N-9; Abstain-0 (re-vote) [Y-12; N-6; Abstain-0 (initial vote)]</p> <p>Conditions/Questions for Developer:</p> <ol style="list-style-type: none"> How does the risk score correlate with the actual expenditures?

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

2. What is the distinction between ETGs and ERGs?
3. Can this measure be applied to the Medicare population?
4. Have there been any changes in the underlying risk model used in the ETGs since what has been published on the Ingenix web site a year ago?
5. How are the carve outs, pharmacy and mental health data handled? How was this data validated?

Developer Response:

1. Ingenix provides options for expenditure thresholds for a patient's annual member costs: \$25,000, \$100,000, and \$250,000. Ingenix explained that these thresholds would vary depending on the application.
2. ETGs are episode-based measures. For example, an episode of diabetes, congestive heart failure or COPD--the severity models are built separately for each of the conditions which allows for risk adjustment for each separate condition-based episode. The results are then tagged for each episode for a member not only by condition, but also by the level of severity. There are hundreds of ETGs that map into the ERGs. Ingenix maps to the ERG designed for the population-based risk adjustment; they weight each of the ERG markers to the final ERG score. The ERGs looks at age, in which case they may be applied to the Medicare population, however not all of the ETGs take age into account in the risk adjustment model. During the developer testing they didn't find that age had much explanatory power so they are not included in all of the ERGs. The ERG will point to a different weight depending on the age of the individual. However, since this measure has only been tested in a commercial database, per NQF policy, it can only be endorsed for use in commercial populations.
3. The ETG models and the risk models related to the ETGs have not been updated or recalibrated within the last year; therefore the information on the Ingenix website is still applicable.
4. Ingenix works with a population that has pharmacy and medical data. Mental health is excluded because the claims are not often available in addition to lack of coding for mental health services. Pharmacy data hasn't been an issue because it's up to the user whether they want to include and compare populations who have pharmacy data. The methodology can be adjusted, you are able to have a mixed population of both medical and pharmacy benefits, and the user is able to isolate the medical resource use data if they choose to.

1. Importance to Measure and Report :Y-16; N-0

Committee Discussion: This criterion was also discussed during the June 6 conference call. To access the summary of this call, [click here](#).

1a. High Impact: H-15; M-1; L-0; I-0

Committee Discussion: The Steering Committee has deemed the measure focus to be high impact.

1b. Resource use/cost problems: H-13; M-3; L-0; I-0

Committee Discussion: The Steering Committee agrees this criterion has been met.

1c. Purpose clearly described: H-12; M-4; L-0; I-0

Committee Discussion: The Steering Committee believes the measure has met this sub criterion, as the measure's purpose is clearly described.

1d. Resource use service categories consistent and representative: H-8; M-8; L-0; I-0

Committee Discussion: The resource use service categories are representative of the measure intent and focus.

2. Scientific Acceptability of Measure Properties: Yes [Y-9; N-6 (Committee Vote)]

2a. Overall Reliability: H-8; M-7; L-1; I-0

Committee Discussion: The Ingenix team has a robust system where they double code the data – the steps that lead to the production of the data has a 99.9% match between the two approaches.. The Committee agreed that tables present measure results it is unclear if they actually represent that the measure is reliable.

2a1. Measure well defined and precisely specified: H-10; M-5; L-1; I-0

Committee Discussion: This measure appears to be well defined and specified. This methodology is used in a number of organizations and appears to work well. This sub criterion has been met.

2a2. Reliability Testing: H-9; M-7; L-0; I-0

Committee Discussion: The Committee agreed that this sub criterion has been met; the results have shown to be repeatable. The Committee suggested more robust reliability testing methods should be explored.

2b. Overall Validity: H-2; M-10; L-3; I-0

Committee Discussion: In the submission, Ingenix states that they apply the methodology to data from several different organizations, but this is not detailed in any of the results. Face validity was tested however there is not any description of the results within the submission. The tables that were submitted to demonstrate validity are not clearly labeled or defined.

2b1. Specifications consistent with intent: H-7; M-8; L-1; I-0

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

Committee Discussion: The Committee agrees the specifications are consistent with the intent.

2b2. Validity Testing: H-0; M-8; L-6; I-0

Committee Discussion: This measure has been demonstrated to meet the requirement for face validity.

2b3. Exclusions: H-9; M-4; L-2; I-0

Committee Discussion: There are no exclusions based on cost or other criteria. The Committee reiterated concerns with comparability for plans that have pharmacy carve outs or do not have pharmacy data to those that do.

2b4. Risk Adjustment: H-6; M-8; L-1; I-0

Committee Discussion: When looking at the ETG codes, a severity score is assigned; the methodology then takes into account the ETG severity score and the number of comorbidities. A retrospective model contains the observed episodes that may occur during that year, but a user will not be able to observe any markers or costs for people who did not undergo services. The ERG risk level determines the individual's ERG risk score which drives the risk adjustment. The Committee acknowledged this methodology is very complex and not completely understood by all members.

2b5. Identification of statistically significant/meaningful differences: H-5; M-7; L-3; I-0

Committee Discussion: There is a way to stratify those with or without pharmacy data. The Committee expressed concern that valid comparisons cannot be made across organizations with different levels of data completeness and consistency.

2b6. Multiple data sources: N/A

2c. Stratification for disparities: H-0; M-4; L-2; I-9

Committee Discussion: This measure does not stratify by race and ethnicity. This may be possible in the future, but at the present time this information is not available.

3. Overall Usability: H-0; M-10; L-5; I-0

Committee Discussion: The Committee questioned on whether this measure has been featured in peer reviewed articles; the developer was unaware of any that could be shared with the Committee. The developers explained that this measure is currently being used to profile physicians. They are unaware of any efforts to publicly report the results, even within health plans to their covered lives.

3a. Measure performance results are publicly reported: H-0; M-4; L-6; I-4

Committee Discussion: Ingenix conducted a survey of their customers, some users are publicly reporting the data and others are sharing information with physicians for incentive based programs. Some users have decided to put the information on a website that goes to their providers, which allows them to access their risk scores and score card. Providers are then able to drill down on the scorecard to the claim base level, the patient level and then the overall claims level.

3b. Measure results are meaningful/useful for public reporting and performance improvement:

Committee Discussion: H-3; M-6; L-3; I-3

3c. Data and results can be decomposed for transparency and understanding: H-1; M-8; L-5; I-1

Committee Discussion: While Ingenix has a transparency website open to the public which explains the methodology and approach to measuring resources, the submission reviewed by the Committee was admittedly complex and at times difficult to identify the relevant information.

3d. Harmonized or justification for differences: N/A

4. Feasibility: H-3; M-8, L-6, I-0

4a. Data elements routinely generated during care process: H-13; M-2; L-2; I-0

Discussion: The Steering Committee believes that this sub criterion has been met; all of the data elements are generated during the care process.

4b. Data elements available electronically: H-14, M-4, L-0, I-0

Discussion: The Steering Committee believes that this sub criterion has been met; all of the data is available electronically.

4c. Susceptibility to inaccuracies/ unintended consequences identified: H-5, M-9, L-3; I-0

Discussion: Mental health is not available and pharmacy data rarely is, when pharmacy data is included it is stratified. Ingenix does not have a formal audit system to ensure that all of the numbers are included & correct. In general, when dealing with any measure that uses administrative data there are various inaccuracies, pertaining particularly to coding inaccuracies and variation. Ingenix provides guidelines how to use small volumes/ sample sizes, however there is not content available to demonstrate this approach. This measure appears less prone to "gaming", as there is not much a user can do to manipulate the start or end of an episode.

4d. Data collection strategy can be implemented: H-1, M-10, L-13, I-1

Discussion: See Ingenix feasibility discussion above.

NATIONAL QUALITY FORUM

898

<p>1603: ETG/ PEG Based Hip Fracture Cost of Care measure (Ingenix)</p> <p>Description: The measure focuses on resources used to deliver episodes of care for patients with Hip Fracture. Hip Fracture episodes are defined using the Episode Treatment Groups (ETG) methodology and describe the unique presence of the condition for a patient and the services involved in diagnosing, managing and treating Hip Fracture. A number of resource use measures are defined for Hip Fracture episodes, including overall cost of care, cost of care by type of service, and the utilization of specific types of services. Each resource use measure is expressed as a cost or a utilization count per episode and comparisons with internal and external benchmarks are made using risk adjustment to support valid comparisons. As requested by NQF, the focus of this submission is for Hip Fracture episodes and will cover both measures at the Hip Fracture base and severity level and also a Hip Fracture composite measure where Hip Fracture episode results are combined across Hip Fracture severity levels. At the most detailed level, the measure is defined as the base condition of Hip Fracture and an assigned level of severity (e.g., resources per episode for Hip Fracture, severity level 1 episodes). Composite measures can then be created using these measurement units to meet a specific need. For example, a composite measure for Hip Fracture is derived by combining Hip Fracture episode results across Hip Fracture severity levels. Appropriate risk adjustment is applied to support comparisons (e.g., for physician measurement, adjusting for a physician's mix of Hip Fracture episodes by severity level when supporting a Hip Fracture composite comparison). The focus of this measure is on Hip Fracture. However, Hip Fracture episode results could also be included in an "orthopedics", "acute care", or other clinical composite for a physician, combining episodes in clinical areas similar to Hip Fracture. Further, an "overall" composite for a physician can be created, again by aggregating episode results across appropriate conditions and severity levels and applying proper risk adjustment when making comparisons.</p> <p>Resource Use Type: Per episode</p> <p>Data Type: Administrative claims, Other</p> <p>Resource Use Service Categories: Inpatient services: Inpatient facility services; Admissions/discharged; Ambulatory services: Outpatient facility services; Emergency Department; Pharmacy; Evaluation and management; Procedures and surgeries; Imaging and diagnostic; Lab services</p> <p>Care Setting: Ambulatory Care : Ambulatory Surgery Center (ASC), Ambulatory Care : Clinic/Urgent Care, Ambulatory Care : Clinician Office, Emergency Medical Services/Ambulance, Home Health, Hospice, Hospital/Acute Care Facility, Imaging Facility, Laboratory, Post Acute/Long Term Care Facility: Nursing Home/Skilled Nursing Facility, Post Acute/Long Term Care Facility: Rehabilitation</p> <p>Level of Analysis: Clinician : Group/Practice, Clinician : Individual, Clinician : Team, Facility, Health Plan, Integrated Delivery System, Population: Community, Population: County or City, Population : National, Population : Regional, Population: State</p> <p>Measure Developer: Ingenix, 950 Winter Street, Suite 3800, Waltham, Massachusetts, 02451</p>
<p>Committee Recommendation for Endorsement: This measure did not pass the scientific acceptability criterion, and is not recommended for endorsement.</p>
<p>Conditions/Questions for Developer:</p> <ol style="list-style-type: none"> 1. Why are different age groups assigned the same risk coefficients, when they will have extremely different risk factors? 2. How does the episode grouper work in terms of low and high outliers? Are you able to provide information on exactly how many episodes have been excluded? 3. Why do you cut the low cost episodes from being included in the measure? <p>Developer Response:</p> <ol style="list-style-type: none"> 1. This represents a limitation of the data set. Due to the minimal number of people over 65 in commercial programs, we didn't have the numbers to further stratify. 2. We exclude cases that are low in cost. We have the data to talk about the number of cases that are excluded by varying a low outlier, yes. 3. The hypothesis that that these low cost episodes – ones under 2.5 percent – are either mistakes or miscodes. They are probably incomplete episodes, so we don't count them.
<p>1. Importance to Measure and Report</p> <p>1a.High Impact: H-2; M-1; L-2; I- 0</p> <p>TAP Discussion: There was general agreement that hip fracture is a major cause of morbidity, mortality and high resource use. The TAP did, however, question the importance of measuring hip fractures in a predominately under 65 group of patients. Ingenix acknowledged that this was a significant limitation of using administrative data.</p> <p>1b. Resource use/cost problems: H-2; M-2; L-1; I-0</p> <p>TAP Discussion: No issues were identified.</p> <p>1c. Purpose clearly described: H-1; M-4; L-0; I-0</p>

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

<p>TAP Discussion: No issues were identified.</p> <p>1d. Resource use service categories consistent and representative: H-2; M-2; L-1; I-0</p> <p>TAP Discussion: The TAP were concerned that resource use service categories omit nursing homes and inpatient or outpatient rehab services.</p>
<p>Overall Importance: Y-10, N-6</p> <p>Committee Discussion: The Committee agreed that hip fractures are a high impact area of healthcare. They were concerned, however, that the measure did not include populations of patients over 65, where the vast majority of hip fractures would occur, and where the nature of hip fractures is a significantly different than it is for younger populations. Ingenix reminded the Committee that the measure was tested in a commercial database, not a Medicare database, and would therefore be endorsed as such. The Committee ultimately questioned whether it was important to measure hip fractures in a younger population at all.</p>
<p>2. Scientific Acceptability of Measure Properties:</p> <p>2a. Overall Reliability: H-1; M-0; L-4; I-0</p> <p>2a1. Measure well defined and precisely specified: H-1; M-2; L-2; I-0</p> <p>TAP Discussion: The TAP was concerned that the measure didn't capture certain co-morbid conditions such as dementia which are critical to understanding resource use for this clinical condition. There was substantial unease that the data does not examine the Medicare population, where the majority of hip-fractures occur.</p> <p>2a2. The results are repeatable: H-1; M-2; L-2; I-0</p> <p>TAP Discussion: The panel questioned whether one could infer grouper reliability from the tables submitted by Ingenix. Ingenix explained that the tables illustrate expected variability in results and point to a relatively consistent cost across health care organizations.</p> <p>2b. Overall Validity: H-0; M-1; L-3; I-0</p> <p>2b1. Evidence is consistent with intent: H-0; M-0; L-5; I-0</p> <p>TAP Discussion: The TAP reiterated their concern that the measure hasn't captured the patient population most likely to be affected by hip fractures. Therefore, the measure may have limited applicability, due to the limitations of using only commercial data. The panel also felt that hip fractures in younger populations versus older populations represent two very different clinical situations.</p> <p>2b2. Score/Analysis: H-0; M-1; L-4; I-0</p> <p>TAP Discussion: The TAP was uncomfortable with the fact that all age groups were assigned the same risk coefficients. Ingenix explained that this also represents a limitation of the data set, where they did not have the numbers over 65 to further stratify. Members of the panel believed that certain clinically relevant co-morbidities and complications such as dementia and post-op delirium should be reported on in a hip-fracture measure.</p> <p>2b3. Exclusions: H-0; M-1; L-4; I-0</p> <p>TAP Discussion: The TAP felt that the reasoning behind the exclusion criteria was unclear and not based on clinical evidence.</p> <p>2b4. Risk Adjustment: H-0; M-0; L-4; I-1</p> <p>TAP Discussion: The developer described how the measure contains low dollar exclusions. The assumption is that these claims represent incomplete episodes.</p> <p>2b5. Identification of statistically significant/meaningful differences: H-0; M-0; L-4; I-1</p> <p>TAP Discussion: There was a discussion regarding the relative cost of care ratio and a question about what numbers represent statistically significant differences. Ingenix explained that the numbers would depend on the confidence interval, the underlying variance of episode cost and the number of total cases.</p> <p>2b6. Multiple data sources: N/A (using all administrative data)</p> <p>2c. Stratification for disparities: H-0; M-1; L-1; I-3</p> <p>TAP Discussion: Racial disparities were addressed in the submission, but the data limits a further examination into these disparities.</p>
<p>Overall Scientifically Acceptable: No [Y-7; N-10 (Committee Vote)]</p> <p>Overall Reliability: H-1; M-11; L-3; I-2</p> <p>Overall Validity: H-0; M-6; L-10; I-0</p> <p>Committee Discussion: The Committee believed the measure was limited in its clinical construction logic as a result of its reliance upon commercial data, where the population of patients with hip fractures was notably low. Thus, the testing completed by Ingenix for this measure represented a fairly <i>uncommon</i> condition – hip fractures in under 65's – when the majority of hip fractures are much more common and different clinically. The Committee agreed, therefore, that significant and meaningful differences could not be produced by this measure, particularly when reporting at an individual physician level. Furthermore, the Committee were concerned with the fact that the grouper function was not tested or reported on, and Ingenix provided no information comparing scoring of attribution over episodes of time</p>
<p>Usability:</p>

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

<p>3a. Measure performance results are publicly reported: H-0; M-2; L-3; I-0 <i>TAP Discussion:</i> The TAP believes this subcriterion has been met.</p> <p>3b. Measure results are meaningful/useful for public reporting and quality improvement: H-0; M-1; L-4; I-0 <i>TAP Discussion:</i> The TAP acknowledged the impressive amount of work Ingenix put into this measure, but again articulated concern that the measure would have limited meaningful use as it is not capturing the appropriate population. The panel was uneasy with the grouping of two clinically different age cohorts together into one measure; they felt that the clinical situation, treatment path and mortality for a younger population with hip fractures versus an older population were different enough to warrant two separate measures.</p> <p>3c. Data and results can be decomposed for transparency and understanding: H-0; M-2; L-3; I-0 <i>TAP Discussion:</i> The TAP agrees this subcriterion has been met.</p> <p>3d. Harmonized or justification for differences: N/A</p>
<p>Overall Usability: This measure did not pass the scientific acceptability criterion. As a result, the Committee did not discuss usability.</p>
<p>4. Feasibility:</p> <p>4a. Data elements routinely generated during care process: H-3; M-1; L-1; I-0 <i>TAP Discussion:</i> The TAP agrees that this subcriterion has been met; all data is routinely generated through the care process.</p> <p>4b. Data elements available electronically: H-4; M-0; L-1; I-0 <i>TAP Discussion:</i> The TAP agrees that this subcriterion has been met; all data is available electronically.</p> <p>4c. Susceptibility to inaccuracies/ unintended consequences identified: H-1; M-1; L-3; I-0 <i>TAP Discussion:</i> The TAP believe that this subcriterion has been met, however Ingenix does not have a formal audit system in order to monitor for inaccuracies.</p> <p>4d. Data collection strategy can be implemented: H-0; M-2; L-2; I-1 <i>TAP Discussion:</i> The TAP believe that this subcriterion has been met. (NQF Staff Note: this is prior to the submission of product pricing information reviewed by the Steering Committee only.)</p>
<p>Overall Feasibility: This measure did not pass the scientific acceptability criterion. As a result, the Committee did not vote on feasibility.</p>

899

<p>1605: ETG Based Asthma Cost of Care Measure(Ingenix)</p> <p>Description: The measure focuses on resources used to deliver episodes of care for patients with Asthma. Asthma episodes are defined using the Episode Treatment Groups (ETG) methodology and describe the unique presence of the condition for a patient and the services involved in diagnosing, managing and treating asthma. A number of resource use measures are defined for asthma episodes, including overall cost of care, cost of care by type of service, and the utilization of specific types of services. Each resource use measure is expressed as a cost or a utilization count per episode and comparisons with internal and external benchmarks are made using risk adjustment to support valid comparisons. As requested by NQF, the focus of this submission is for Asthma episodes and will cover both measures at the Asthma base and severity level and also an Asthma composite measure where Asthma episode results are combined across Asthma severity levels. At the most detailed level, the measure is defined as the base condition of Asthma and an assigned level of severity (e.g., resources per episode for Asthma, severity level 1 episodes). Composite measures can then be created using these measurement units to meet a specific need. For example, a composite measure for Asthma is derived by combining Asthma episode results across Asthma severity levels. Appropriate risk adjustment is applied to support comparisons (e.g., for physician measurement, adjusting for a physician's mix of Asthma episodes by severity level when supporting an Asthma composite comparison). The focus of this measure is on Asthma. However, Asthma episode results could also be included in a "pulmonologist", "chronic care", or other clinical composite for a physician, combining episodes in clinical areas similar to Asthma. Further, an "overall" composite for a physician can be created, again by aggregating episode results across appropriate conditions and severity levels and applying proper risk adjustment when making comparisons.</p> <p>Resource Use Type: Per episode Data Type: Administrative claims, Other Resource Use Service Categories: Inpatient services: Inpatient facility services, Inpatient services: Admissions/discharges, Ambulatory services: Outpatient facility services, Ambulatory services: Emergency Department, Ambulatory services: Pharmacy, Ambulatory services: Evaluation and management, Ambulatory services: Procedures and surgeries, Ambulatory services: Imaging and diagnostic, Ambulatory services: Lab services Care Setting: Ambulatory Care : Ambulatory Surgery Center (ASC), Ambulatory Care : Clinic/Urgent Care, Ambulatory Care : Clinician Office, Emergency Medical Services/Ambulance, Home Health, Hospice, Hospital/Acute Care Facility, Imaging Facility, Laboratory, Post</p>

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

<p>Acute/Long Term Care Facility: Nursing Home/Skilled Nursing Facility Level of Analysis: Clinician : Group/Practice, Clinician : Individual, Clinician : Team, Facility, Health Plan, Integrated Delivery System, Population: Community, Population: County or City, Population : National, Population : Regional, Population: State Measure Developer: Ingenix, 950 Winter Street, Suite 3800, Waltham, Massachusetts, 02451</p>
<p>Committee Recommendation for Endorsement: Y-7; N-9; Abstain-0</p>
<p>Conditions/Questions for Developer:</p> <ol style="list-style-type: none"> 1. Can you give us more information on how repeatability and "consistency" were determined? The results don't appear consistent. 2. Are patients with COPD excluded? 3. How are results reported and interpreted? 4. How would a smaller health plan implement this measure? It seems it might be too complex and burdensome. <p>Developer Response:</p> <ol style="list-style-type: none"> 1. Repeatability was demonstrated by programming the measure in SAS code and the Ingenix software and comparing results. Because there are differences in what geographies these health plans are pulling from, variation is expected. But while differences across HCO's are expected, whether the differences are too high or low is difficult to know. 2. Patients are excluded from the asthma episode if they have more costs attributable to COPD than asthma. 3. The main measurement is the O/E ratio metric - the numerator of which is the cost of all the episodes of asthma, and the denominator which is the expected costs. 4. The burden depends on the plan's familiarity with ETGs and similar products, and for those who are just starting out, there is unlimited training involved (i.e. help desk support, etc.). There is another option where Ingenix takes the data and runs it themselves - or uses their PCQ Connect product that prepared the data into report-ready formats.
<p>1.Importance to Measure and Report 1a.High Impact: H-9; M-0; L-0 <i>TAP Discussion:</i> The TAP agrees that asthma is a very important health care area to measure. 1b. Resource use/cost problems: H-8; M-1; L-0 ; I-0 <i>TAP Discussion:</i> The TAP agrees the Measure demonstrates cost problems and opportunity for improvement. 1c. Purpose clearly described: H-7; M-2; L-0; I-0 <i>TAP Discussion:</i> The TAP believes the purpose and objective of the measure are clear. 1d. Resource use service categories consistent and representative: H-7; M-2; L-0; I-0 <i>TAP Discussion:</i> The TAP feel this subcriterion has been met.</p>
<p>Overall Importance: Y-16, N-0 Committee Discussion: The Steering Committee agreed that asthma constitutes a high impact healthcare area.</p>
<p>2. Scientific Acceptability of Measure Properties: 2a. Reliability: H-0; M-8; L-1; I-0 2a1.Measure well defined and precisely specified: H-2; M-6; L-1; I-0 <i>TAP Discussion:</i> This measure is one that's part of a suite of episodes around diseases and conditions included in Ingenix's episode treatment grouper. This product identifies claims that should be part of an episode of asthma and divides them into year-long segments, looking at asthma as a chronic disease. The episodes are severity adjusted using clinical markers called condition status factors. Anchor episodes, or face-to-face encounters, are merged together into one episode (i.e. "asthma"). 2a2. The results are repeatable: H-3; M-5; L-1; I-0 <i>TAP Discussion:</i> The TAP didn't understand why Ingenix used three different population samples, rather than taking a portion of the larger population and testing it multiple times. They would like better communication on the approach as well as more detailed depiction of the data. Repeatability was generally determined to be demonstrated adequately, but for the above reasons, some did question the reliability of the measure score. 2b. Overall Validity: H-0; M-6; L-1; I-2 2b1. Evidence is consistent with intent: H-2; M-5; L-1; I-1 <i>TAP Discussion:</i> It was unclear to the panel whether Ingenix is actually measuring asthma costs as intended. The determination of what is an asthma cost and what is not isn't transparent. They also agreed that any results are going to be questioned when potentially over 50% of the costs (the pharmacy costs) are not represented. There were suggestions to stratify those health plans that have pharmacy carve-out arrangements. 2b2.Score/Analysis: H-1; M-4; L-2; I-2 <i>TAP Discussion:</i> Face validity was determined to be appropriate. The TAP continued to express concern about the exclusion of</p>

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

<p>pharmacy costs, which were agreed to be a significant component of asthma care. Pharmacy data is not a requirement to get into the episode (for all ETGs).</p> <p>2b3. Exclusions: H-1; M-7; L-1; I-0 <i>TAP Discussion:</i> The TAP was concerned about the lack of transparency regarding which costs were excluded, and why. Confusion existed around what the grouper identified as outliers or exclusions. Winsorizing very high cost episodes, the top 2%, effectively excludes those kinds of patients that would be important to know about. Addition information such as sensitivity analyses would have helped explain the impact of these high cost cases.</p> <p>2b4. Risk Adjustment: H-1; M-4; L-2; I-2 <i>TAP Discussion:</i> The TAP expressed the same concerns regarding the risk-adjustment methodology as they had for previous Ingenix measures. The TAP was apprehensive that because the measure doesn't require use of standardized costs, the playing field is not level and it can't be implemented consistently across organizations if one is using standard and another actual pricing. To examine how refined the risk-adjustment is, R-squares for different severity levels and how they predict resource utilization should be provided.</p> <p>2b5. Identification of statistically significant/meaningful differences: H-0; M-8; L-0; I-1 <i>TAP Discussion:</i> The TAP felt confident in Ingenix's methodology after it was explained.</p> <p>2b6. Multiple data sources: N/A (using all administrative data)</p> <p>2c. Stratification for disparities: H-2; M-6; L-0; I-1 <i>TAP Discussion:</i> Gender and age can be stratified, but race data is not available.</p>
<p>Overall Reliability: H-1; M-14; L-1; I-0 Overall Validity: H-0; M-8; L-8; I-0 Overall Scientifically Acceptable: Yes [Split vote [Y-8; N-8 (Committee Vote)]] Committee Discussion: The Committee struggled with the circuitous reasoning behind asthma with acute exacerbation being a condition status and then having that condition status factor into the assignment of severity levels. Ingenix defended this methodology by explaining that for all measures, everything related to severity is based on utilization, which, although circular, is the best possible option. The Committee reiterated the TAP's concern that over half of asthma resource use costs are not captured in this measure since pharmacy data is not collected. They expressed unease about the incomparability of entities that have pharmacy data to those that do not.</p>
<p>Usability:</p> <p>3a. Measure performance results are publicly reported: H-2; M-4; L-2; I-1 <i>TAP Discussion:</i> This product is generally used with a suite of ETG's, usually in combination with the pneumonia and COPD measures. There was uncertainty about the measure's usefulness on its own. Since Ingenix can't ascertain if this measure is being used individually the concern from the panel is how the individual measure could be used.</p> <p>3b. Measure results are meaningful/useful for public reporting and quality improvement: H-0; M-6; L-2; I-1 <i>TAP Discussion:</i> The TAP was concerned about the possibility of misinterpretation of results because of the transparency and usability of the results of this measure.</p> <p>3c. Data and results can be decomposed for transparency and understanding: H-3; M-5; L-1; I-0 <i>TAP Discussion:</i> The TAP reiterated their concern of the transparency of the score. Ingenix clarified that there are ways to drill into different aspects of care to see how they might be driving the score.</p> <p>3d. Harmonized or justification for differences: N/A</p>
<p>Overall Usability: H-0; M-9; L-6; I-1 Committee Discussion: Several Steering Committee members challenged the idea that asthma should be thought of in terms of "episodes," as it is a chronic condition.</p>
<p>4. Feasibility:</p> <p>4a. Data elements routinely generated during care process: H-7; M-2; L-0; I-0 <i>TAP Discussion:</i> The TAP believes this subcriterion has been met; data is a byproduct of care.</p> <p>4b. Data elements available electronically: H-7; M-2; L-0; I-0 <i>TAP Discussion:</i> The TAP agrees this subcriterion has been met; data is available electronically.</p> <p>4c. Susceptibility to inaccuracies/ unintended consequences identified: H-1; M-8; L-0; I-0 <i>TAP Discussion:</i> The TAP was generally comfortable with the error checks built into the product.</p> <p>4d. Data collection strategy can be implemented: H-4; M-4; L-0; I-1 <i>TAP Discussion:</i> The TAP expressed some concern about the burden this measure would place on a programmer to implement, particularly at smaller health plans.</p>
<p>Overall Feasibility: H-1; M-8; L-7; I-0</p>

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

900

Committee Discussion: See Ingenix feasibility discussion above.

1608: ETG Based Chronic Obstructive Pulmonary Disease Cost of Care Measure (COPD) (Ingenix)

Description: The measure focuses on resources used to deliver episodes of care for patients with COPD. COPD episodes are defined using the Episode Treatment Groups (ETG) methodology and describe the unique presence of the condition for a patient and the services involved in diagnosing, managing and treating COPD. A number of resource use measures are defined for COPD episodes, including overall cost of care, cost of care by type of service, and the utilization of specific types of services. Each resource use measure is expressed as a cost or a utilization count per episode and comparisons with internal and external benchmarks are made using risk adjustment to support valid comparisons.

As requested by NQF, the focus of this submission is for COPD episodes and will cover both measures at the COPD base and severity level and also a COPD composite measure where COPD episode results are combined across COPD severity levels. At the most detailed level, the measure is defined as the base condition of COPD and an assigned level of severity (e.g., resources per episode for COPD, severity level 1 episodes). Composite measures can then be created using these measurement units to meet a specific need. For example, a composite measure for COPD is derived by combining COPD episode results across COPD severity levels. Appropriate risk adjustment is applied to support comparisons (e.g., for physician measurement, adjusting for a physician's mix of COPD episodes by severity level when supporting a COPD composite comparison). The focus of this measure is on COPD. However, COPD episode results could also be included in a "pulmonary" "chronic care", or other clinical composite for a physician, combining episodes in clinical areas similar to COPD. Further, an "overall" composite for a physician can be created, again by aggregating episode results across appropriate conditions and severity levels and applying proper risk adjustment when making comparisons.

Resource Use Type: Per episode

Data Type: Administrative claims, Other

Resource Use Service Categories:

Inpatient services: Inpatient facility services, Inpatient services: Admissions/discharges, Ambulatory services: Outpatient facility services, Ambulatory services: Emergency Department, Ambulatory services: Pharmacy, Ambulatory services: Evaluation and management, Ambulatory services: Procedures and surgeries, Ambulatory services: Imaging and diagnostic, Ambulatory services: Lab services

Care Setting: Ambulatory Care : Ambulatory Surgery Center (ASC), Ambulatory Care : Clinic/Urgent Care, Ambulatory Care : Clinician Office, Emergency Medical Services/Ambulance, Home Health, Hospice, Hospital/Acute Care Facility, Imaging Facility, Laboratory, Post Acute/Long Term Care Facility: Nursing Home/Skilled Nursing Facility

Level of Analysis: Clinician : Group/Practice, Clinician : Individual, Clinician : Team, Facility, Health Plan, Integrated Delivery System, Population: Community, Population: County or City, Population : National, Population : Regional, Population: State

Measure Developer: Ingenix, 950 Winter Street, Suite 3800, Waltham, Massachusetts, 02451

Committee Recommendation for Endorsement: This measure did not pass the scientific acceptability criterion, and is not recommended for endorsement.

Conditions/Questions for Developer:

1. What was the clinical logic of using 180 days, particularly since your Asthma measure had used 365 days, and both are similar chronic conditions?

Developer Response:

1. We will have to examine that further.

1. Importance to Measure and Report

1a. High Impact: H-7; M-0; L-0; I-0

TAP Discussion: The TAP agreed Ingenix did well with articulating the high impact of COPD.

1b. Resource use/cost problems: H-7; M-0; L-0; I-0

TAP Discussion: The TAP believe that COPD represents a resource use issue that can be addressed.

1c. Purpose clearly described: H-7; M-0; L-0; I-0

TAP Discussion: The TAP feel the purpose and objective are clear.

1d. Resource use service categories consistent and representative: H-6; M-1; L-0; I-0

TAP Discussion: The TAP believes this subcriterion has been met.

Overall Importance: Y-16, N-0

Committee Discussion: There was unanimous agreement that asthma constitutes a high impact area of healthcare.

2. Scientific Acceptability of Measure Properties:

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

<p>2a. Overall Reliability: H-4; M-3; L-0; I-0 2a1. Measure well defined and precisely specified: H-4; M-3; L-0; I-0 <i>TAP Discussion:</i> The TAP discussion focused around the clinical logic around the timeframes chosen. 2a2. The results are repeatable: H-5; M-2; L-0; I-0 <i>TAP Discussion:</i> The TAP agrees that reliability for this measure is similar to the previously discussed Ingenix asthma measure. 2b. Overall Validity: H-0; M-7; L-0; I-0 2b1. Evidence is consistent with intent: H-2; M-5; L-0; I-0 <i>TAP Discussion:</i> The TAP believes this subcriterion has been met. 2b2. Score/Analysis: H-0; M-7; L-0; I-0 <i>TAP Discussion:</i> The TAP remained concerned about Ingenix's testing method for customization, the inability to compare actual versus standardized prices, and the high level of pharmacy exclusions. 2b3. Exclusions: H-1; M-6; L-0; I-0 <i>TAP Discussion:</i> There are no clinical exclusions, only administrative ones. The TAP felt it was unclear how tie-breaking logic works and noted that it was not specified in the submission how COPD and asthma ETG's interact. 2b4. Risk Adjustment: H-0; M-4; L-3; I-0 <i>TAP Discussion:</i> While Ingenix had a nice description of how they developed their risk-adjustment approach, the panel would have liked to see more description of the modeling presented in the submission. 2b5. Identification of statistically significant/meaningful differences: H-0; M-7; L-0; I-0 <i>TAP Discussion:</i> The TAP questioned whether the practical significance of the measure since it is a relative cost ratio. 2b6. Multiple data sources: N/A (using all administrative data) 2c. Stratification for disparities: H-2; M-5; L-0; I-0 <i>TAP Discussion:</i> Only gender and age are stratified for. Race data is not available.</p>
<p>Overall Reliability: H-3; M-10; L-2; I-0 Overall Validity: H-1; M-5; L-9; I-0 Overall Scientifically Acceptable: Yes [Y-5; N-10 (Committee Vote)] Committee Discussion: The Steering Committee appreciated the change Ingenix made to the measure's timeframe at the TAP's suggestion, from 180 to 365 days, to remain consistent with the asthma measure. It was felt the analysis of scientific acceptability for this measure would generally reflect the same analysis for measure 1560 Asthma.</p>
<p>Usability: 3a. Measure performance results are publicly reported: H-0; M-7; L-0; I-0 <i>TAP Discussion:</i> The TAP expressed doubts regarding whether the measure could be implemented in a user-friendly manner. 3b. Measure results are meaningful/useful for public reporting and quality improvement: H-0; M-7; L-0; I-0 <i>TAP Discussion:</i> The panel agreed that measure provides useful information for individual health plans. However, they expressed concern about how useful it would be to compare across health plans, due to the fact that standardized pricing is not required. 3c. Data and results can be decomposed for transparency and understanding: H-3; M-4; L-0; I-0 <i>TAP Discussion:</i> It was agreed that previous discussions regarding Ingenix transparency would also apply to this measure. 3d. Harmonized or justification for differences: N/A</p>
<p>Overall Usability: This measure did not pass the scientific acceptability criterion. As a result, the Committee did not discuss usability.</p>
<p>4. Feasibility: 4a. Data elements routinely generated during care process: H-5; M-2; L-0; I-0 <i>TAP Discussion:</i> The TAP believes this subcriterion has been met; data is a byproduct of care. 4b. Data elements available electronically: H-7; M-0; L-0; I-0 <i>TAP Discussion:</i> The TAP believes this subcriterion has been met; data available electronically. 4c. Susceptibility to inaccuracies/ unintended consequences identified: H-3; M-4; L-0; I-0 <i>TAP Discussion:</i> The TAP is comfortable that Ingenix can accurately identify inaccuracies and errors. 4d. Data collection strategy can be implemented: H-6; M-1; L-0; I-0 <i>TAP Discussion:</i> The TAP believes this subcriterion has been met.</p>
<p>Overall Feasibility: This measure did not pass the scientific acceptability criterion. As a result, the Committee did not vote on feasibility.</p>

901

902

NATIONAL QUALITY FORUM

903 **Candidate Consensus Standards with No Committee Consensus**

904 The Committee was unable to come to consensus on one candidate consensus standard.

905

906 The following evaluation summary table summarizes the results of the TAP’s and Committee’s
 907 evaluation of and voting on the candidate consensus standard that did not draw Committee
 908 consensus. A hyperlink is provided in summary table to the detailed measure specifications. To
 909 access the meeting transcripts and recordings in which this measure is discussed, refer to the
 910 [project web page](#).

911 (1595) ETG based diabetes cost of care measure (Ingenix).....58

912
 913 **Evaluation Summary—Candidate Consensus Standard with No Committee**
 914 **Consensus**

915

1595: ETG Based Diabetes Cost of Care Measure (Ingenix)
<p>Description: The measure focuses on resources used to deliver episodes of care for patients with Diabetes. Diabetes episodes are defined using the Episode Treatment Groups (ETG) methodology and describe the unique presence of the condition for a patient and the services involved in diagnosing, managing and treating diabetes. A number of resource use measures are defined for diabetes episodes, including overall cost of care, cost of care by type of service, and the utilization of specific types of services. Each resource use measure is expressed as a cost or a utilization count per episode and comparisons with internal and external benchmarks are made using risk adjustment to support valid comparisons. The focus of this submission is for Diabetes episodes and will cover both measures at the Diabetes base and severity level and also a Diabetes composite measure where Diabetes episode results are combined across Diabetes severity levels. At the most detailed level, the measure is defined as the base condition of diabetes and an assigned level of severity (e.g., resources per episode for diabetes, severity level 1 episodes).</p> <p>Resource Use Measure Type: Per episode</p> <p>Data Source: Administrative claims, Other</p> <p>Resource Use Service Category: Inpatient services: Inpatient facility services; Inpatient services: Admissions/discharges; Ambulatory services: Outpatient facility services; Ambulatory services: Emergency Department; Ambulatory services: Pharmacy; Ambulatory services: Evaluation and management; Ambulatory services: Procedures and surgeries; Ambulatory services: Imaging and diagnostic; Ambulatory services: Lab services</p> <p>Care Setting: Ambulatory Care: Ambulatory Surgery Center (ASC), Ambulatory Care: Clinic/Urgent Care, Ambulatory Care: Clinician Office, Home Health, Hospice, Hospital/Acute Care Facility, Imaging Facility, Laboratory</p> <p>Level of Analysis: Clinician: Group/Practice, Clinician: Individual, Clinician: Team, Facility, Health Plan, Integrated Delivery System, Population: Community, Population: County or City, Population: National, Population : Regional</p> <p>Measure Developer: Ingenix</p>
Committee Recommendation for Endorsement: Y-7; N-7; Abstain -0 (re-vote) [Y-11; N-7; Abstain-0 (initial vote)]
<p>1. Importance to Measure and Report:</p> <p>1a. High Impact: H-9 ; M-0 ; L-0 ; I-0 <i>TAP Discussion:</i> The TAP believes this is a high cost, impact aspect of healthcare; this subcriteria has been met.</p> <p>1b. Resource use/cost problems: H- 3 ; M-6 ; L-0 ; I-0 <i>TAP Discussion:</i> The TAP would have liked to see more evidence of provider variation and other types of variation in treating diabetes in addition to the regional variation.</p> <p>1c. Purpose clearly described: H- 4 ; M-5 ; L-0 ; I-0 <i>TAP Discussion:</i> The TAP believes that the intent provided not specific to this diabetes measure, it is a very general statement.</p> <p>1d. Resource use service categories consistent and representative: H- 9 ; M-0 ; L-0 ; I-0</p>

NATIONAL QUALITY FORUM

<p>1595: ETG Based Diabetes Cost of Care Measure (Ingenix)</p> <p><i>TAP Discussion:</i> The TAP believes this subcriterion has been met.</p>
<p>Overall Importance: Y-18, N-0</p> <p><i>Committee Discussion:</i> The Steering Committee believes this is a high impact area that should be measured; this subcriterion has been met.</p>
<p>2. Scientific Acceptability of Measure Properties:</p> <p>2a1. Well defined/precise specifications: H- 5 ; M-3 ; L-1 ; I-0</p> <p><i>TAP Discussion:</i> Specifications for co-morbidities, severity levels, etc. are not clear. It is unclear if severity ratings are weighted based on services of comparable cost. Only costs that are mapped back to the diabetes code are counted in the episode. The measure is stratified by severity level not clinical condition. Concerns about how patients with pharmacy benefit (or who run out of pharmacy benefit) are compared to those with full pharmacy benefit.</p> <p>2a2. Reliability testing: H- 7 ; M-1; L-0 ; I-0</p> <p><i>TAP Discussion:</i> Demonstration of internal consistency was presented to demonstrate reliability. The Committee requested additional reliability tests in during maintenance. Additional detail in terms of the r2 of the risk adjustment model and calibration results was requested.</p> <p>2b1. Specifications consistent with resource use/cost problem: H- 1 ; M-6 ; L-1 ; I-0</p> <p><i>TAP Discussion:</i> TAP was unclear on whether diabetes education codes were included in the specifications?</p> <p>2b2. Validity testing: H- 4 ; M-3; L-0 ; I-1</p> <p><i>TAP Discussion:</i> The TAP believes adequate validity testing information provided. More robust methods should be considered in future evaluations.</p> <p>2b3. Exclusions: H-0; M-7 ; L1 ; I-0</p> <p><i>TAP Discussion:</i> TAP was unclear on how exclusions were identified.</p> <p>2b4. Risk adjustment : H-0 ; M-4 ; L-4 ; I-0</p> <p><i>TAP Discussion:</i> The TAP was concerned about the inability to distinguish between complications and comorbidities.</p> <p>2b5. Identification of statistically significant/meaningful differences: H- 0 ; M-4 ; L-4 ; I-0</p> <p><i>TAP Discussion:</i> Insufficient evidence that the sample size threshold and analysis at the physician level is meaningful at that level. Unclear how the 30 sample size was selected.</p> <p>2b6. Multiple data sources: N/A</p> <p>2c. Stratification for disparities: H- 0 ; M-0 ; L-0 ; I-0; N/A-9</p> <p><i>TAP Discussion:</i> Due to the limitations in the administrative claims data, at this time the measure does not stratify for disparities.</p>
<p>Overall Scientifically Acceptable: Yes [Y-10; N-8 (Committee Vote)]</p> <p><i>Committee Discussion:</i> As an introduction to the measure, the developer summarized their responses to the TAP concerns including that the diabetes education codes have been confirmed and are included in the specifications. Similar to the TAP, the Committee expressed concern about the minimum sample size guideline suggesting 30 cases per physician; the Committee questioned how this number was identified and if any statistical analysis was performed to support this guideline. In response to this concern, the developer explained that this sample size was borrowed from previous work done by NCQA on resource utilization and stated that from their perspective, while sample size can be important, ensuring results are statistically significant is more important. The Committee also requested explanation of the attribution model, finding that it was very complex, and questioned of the total sample from their analysis, what percent of physicians have a minimum sample size of 30. The developer explained that the attribution model seeks to identify the highest number of contacts between the physician and the patient related to diabetes; in case of a tie, the provider with the highest actual cost gets attributed the episode. Another concern identified by the Committee relates to how the measure captures costs related to the sequela of diabetes (e.g., renal disease, eye disease, CHF); the measure as presented does not currently account for these costs as they trigger alternate episodes. There was also discussion on how this measure (or measures like it) might be paired with quality (process) measures, as it measures resource use and adjusts for conditions <i>before</i> care is provided. The Committee also spent some time discussing and trying to understand the episode trigger mechanisms, such as when a patient enters the episode in the middle of the 12-episode; in this case the episode is marked incomplete. There was a question to the developer about what percentage of the claims was higher or lower than expected. The developer was unable to answer the question off hand but will get back to the Committee with this information. The issue of mental health and pharmacy carve outs was a prevalent issue throughout the discussion of these measures. For this measure mental health is not stratified for when it is carved out.</p>
<p>3. Usability:</p> <p>3a. Measure performance results are publicly reported: H- 0; M-1 ; L-1; I-6</p> <p><i>TAP Discussion:</i> The usability information submitted is not specific to diabetes, but for all Ingenix measures. TAP expressed concerns</p>

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

<p>1595: ETG Based Diabetes Cost of Care Measure (Ingenix)</p> <p>with the availability of this data to the public and requested clarification from NQF on what is required for "public reporting". The NQF CSAC and BOD continue to discuss this issue; NQF staff will continue to filter any new information on the refining of this policy to the TAP to facilitate final ratings of this usability criterion.</p> <p>3b. Measure results are meaningful/useful for public reporting and quality improvement: H- 0 ; M-4 ; L-2 ; I-2 <i>TAP Discussion:</i> The usability information submitted is not specific to diabetes, but for all Ingenix measures.</p> <p>3c. Data and results can be decomposed for transparency and understanding: H- 1 ; M-2 ; L-5 ; I-0 <i>TAP Discussion:</i> The usability information submitted is not specific to diabetes, but for all Ingenix measures. Challenges for the use of this measure include, complexity, lack of specificity in specifications. The TAP agrees it is difficult to assess the extent of which the measure can be decomposed as currently specified.</p> <p>3d. Harmonized or justification for differences: H-0 ; M-0 ; L-0 ; I-0 ; N-9 <i>TAP Discussion:</i> The usability information submitted is not specific to diabetes, but for all Ingenix measures.</p> <p>Overall Usability: H-0; M-9; L-6; I-3 Committee Discussion: While there is a transparency website for physicians to go to in order determine what a score means, it may take a lot of time to do this. The Steering Committee questioned whether this is a reasonable expectation and adequately demonstrates transparency. Other concerns raised by the Steering Committee were related to the attribution model and how the complexity of the methodology might impact how understandable the measure construction and results are. Because this measure is part of an episode grouper and is not used in isolation as an individual measure, the information the developer was able to present on its current use is not specific to the diabetes episode, but the product as a whole.</p> <p>4. Feasibility: 4a. Data elements routinely generated during care process: H- 8 ; M-0 ; L-0 ; I-0 <i>TAP Discussion:</i> The TAP agrees this subcriterion has been met; measures rely on administrative data. 4b. Data elements available electronically: H-8 ; M-0 ; L-0; I-0 <i>TAP Discussion:</i> The TAP agrees this subcriterion has been met; administrative data are in electronic format. 4c. Susceptibility to inaccuracies/ unintended consequences identified: H-2 ; M-2 ; L-4 ; I-0 <i>TAP Discussion:</i> The TAP does not feel this subcriterion was adequately met; there are current issues identified with specifications could result in inaccuracies and errors. 4d. Data collection strategy can be implemented: H- 5 ; M-2 ; L-1 ; I-0 <i>TAP Discussion:</i> The TAP agrees that barriers to use are minimal. (NQF Note: This is prior to the submission of product pricing information reviewed only by the Steering Committee).</p> <p>4. Feasibility: H-2; M-8; L-8; I-0 Committee Discussion: See Ingenix feasibility discussion above.</p>

916

917 **WITHDRAWN BY DEVELOPER**

918 The 12 measures listed below were withdrawn from the Cycle 2 review process by the
 919 developers for further refinement and testing.

920

921 **Pulmonary**

- 922 • (1577) Episode of care for patients with asthma over a one year period (ABMS-REF)
- 923 • (1581) Episode of care for patients with stable chronic obstructive pulmonary disease
- 924 over a one year period (ABMS-REF)
- 925 • (1582) Episode of care for patients with unstable chronic obstructive pulmonary disease
- 926 over a one year period (ABMS-REF)
- 927 • (1587) Episode of care for ambulatory pneumonia (ABMS-REF)

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

- 928 • (1588) Episode of care for community acquired pneumonia hospitalization (ABMS-REF)
929

930 **Cancer**

- 931 • (1578) Episode of care for 60-day period preceding breast biopsy (ABMS-REF)
932 • (1579) Episode of care for cases of newly diagnosed breast cancer over a 15 month
933 period (ABMS-REF)
934 • (1583) Episode of care for 21-day period around a colonoscopy (ABMS-REF)
935 • (1584) Episode of care for treatment of localized colon cancer (ABMS-REF)
936

937 **Bone/Joint**

- 938 • (1585) Episode of care for simple, non-specific lower back pain (acute and subacute)
939 (ABMS-REF)
940 • (1586) Episode of care for acute/subacute lumbar radiculopathy with or without lower
941 back pain (ABMS-REF)
942 • (1610) ETG based low back pain resource use measure (Ingenix)
943
944

945 **ADDITIONAL CONSIDERATIONS**

946 As the first NQF resource use measure review and evaluation process concludes, there is a great
947 opportunity to reflect on and provide recommendations for future efforts in this area. While
948 resource use measurement has been used in the commercial sector for many years, the emerging
949 interest in using these measures for public reporting and payment initiatives further highlights
950 the need for efforts such as this to explore the complexities and potential challenges for multiple
951 applications. The Committee was asked to provide guidance to the field on how measure of
952 resource use can be improved and provide special considerations for developing measures for the
953 Medicare program. Additionally, members provided insight on how measures of resource use
954 may be linked with quality measures to provide a true assessment of efficiency. Through this
955 exercise, the Committee offered recommendations in several areas. In doing so, several new
956 principles emerged.

957

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

958 1) Submitting and Evaluating Resource Use Measures

959 *Emerging Principle 1: While guidelines in measure components may be acceptable for internal*
960 *quality improvement, to promote measurement for comparison across entities nationally, the*
961 *entire resource use measure construct should be standardized in the form of specifications.*

962 *Emerging Principle 2: The risk adjustment model applied to the measure should be specific to the*
963 *intended population.*

964 *Emerging Principle 3: The factors in the risk adjustment model and severity model should be*
965 *confirmed to be a contributor to the outcome of the measure.*

966 *Emerging Principle 4: In addition to statistical significance, justification of the variables used in*
967 *the risk-adjustment model should be provided based on either clinical relevance or evidence in*
968 *the literature.*

969 2) Reliability and Validity Testing

970 *Emerging Principle 5: To demonstrate reliability of a resource use measure, developers can*
971 *focus on precision of the measure score.*

972 *Emerging Principle 6: When there is such limited variability in a data set that it does not*
973 *adequately distinguish performance differences among providers, reliability cannot simply rely*
974 *on confidence intervals; sample size should also be included in the reliability assessment.*

975 *Emerging Principle 7: The gold standard approach to determining the validity of data elements*
976 *based on administrative claims data in resource use measures is to assess the agreement of*
977 *claims data with source of the data elements in the chart.*

978 3) Data Quality and Comprehensiveness

979 *Emerging Principle 8: Data sets used to measure resources should be as comprehensive as*
980 *possible. Efforts to obtain clinical and carved-out data (e.g., pharmacy, behavioral health) should*
981 *be made to ensure the data set used to calculate resource use is robust, complete, and*
982 *representative.*

983 *Emerging Principle 9: Measure scores calculated and reported using data with carve-outs*
984 *should be labeled as such.*

985 *Emerging Principle 10: Comparisons of entities with and without carved-out data is inappropriate.*

NATIONAL QUALITY FORUM

986 *Emerging Principle 11: If a measure is intending to measure a clinical condition that encompasses*
987 *a predominant portion of its costs in pharmacy claims, consider whether costs should be measured*
988 *at all in the absence of these data. It is the developers' responsibility to conduct an analysis to*
989 *determine whether the lack of these data invalidates the measure score or comparisons.*

990 **4) Measuring Cost and Resource Use in the Medicare Population**

991 *Emerging Principle 12: A patient-centered approach should be used to describe the interaction*
992 *of conditions (and episodes) in the development of resource use measures for the Medicare*
993 *population.*

994 **5) Linking Quality and Cost to Develop Measures of Efficiency and Value**

995 *Emerging Principle 13: Efficiency measurement approaches should be patient-centered, building*
996 *upon previous efforts such as the NQF Patient-Centered Episodes of Care (EOC) Efficiency*
997 *Framework.*

998

999 Recommendations and emerging principles are discussed below:

1000 **1) Submitting and Evaluating Resource Use Measures**

1001 In an effort to minimize the confusion in the submission and evaluation processes, the
1002 Committee identified areas within the resource use measure specification modules that should be
1003 clarified so that the developer understands the information that is required for the measure to be
1004 considered fully. The Committee recognized that in an effort to improve the clarity of the
1005 measure submissions, there should also be attention paid to how new submission requirements
1006 will affect the burden on the developer to submit measures for consideration. While there are
1007 some areas of the submission that will need additional information and more clarity, there may
1008 be other pieces of information that may not be required.

1009 *Emerging Principle 1: While guidelines in measure components may be acceptable for internal*
1010 *quality improvement, to promote measurement for comparison across entities nationally, the*
1011 *entire resource use measure construct should be standardized in the form of specifications.*

1012 The data protocol module components were framed as flexible user instructions for missing data,
1013 data inclusion and exclusions, and data cleaning that could be submitted as specifications or

NATIONAL QUALITY FORUM

1014 guidelines. All of the measure submissions and all of the data protocol components were
1015 submitted as guidelines. Allowing for flexibility in this module led to some discomfort for the
1016 experts specifically related to handling missing data. Ensuring that the data used to run the
1017 resource measures are complete and representative is a critical first step to generating valid
1018 measure results. Allowing flexibility in these steps could allow for errors and inconsistent
1019 implementation of the data cleaning and data preparation steps. As such, the Committee
1020 recommends that the steps within the data protocol module be submitted as specifications going
1021 forward. Specifically, it should be indicated explicitly in the submission whether it is acceptable to
1022 implement the measure using a data set with carve-outs. Data cleaning steps should be explicitly
1023 stated as specifications.

1024
1025 Likewise, in the reporting module, while the attribution approach could be submitted as
1026 specifications or guidelines, the Committee was very concerned with how the models reviewed
1027 might be applied, even as guidelines. While there was not widespread agreement on any of the
1028 attribution approaches reviewed, the Committee recognized that there must be some attribution
1029 approach employed with the use of these measures to facilitate actionable measurement since many
1030 states and healthcare systems may require varied approaches for their unique market. This
1031 highlights the need for more discussion on how, if at all, attribution approaches should be evaluated
1032 in this process where the goal is to endorse standardized approaches to measurement.

1033 *Emerging Principle 2: The risk-adjustment model applied to the measure should be specific to the*
1034 *intended population.*

1035 After reviewing various risk-adjustment approaches presented in the measures submitted to this
1036 project, the Committee agreed measure developers need to demonstrate that the specified risk
1037 models are appropriate for the target population. For instance, if the hierarchical condition category
1038 (HCC) model is used to measure a commercial population, developers need to demonstrate that it is
1039 appropriate for use outside of a Medicare population. The Committee agreed that risk models have
1040 unique weights for comorbidities and may not include all relevant conditions (for example,
1041 pregnancy) when the risk-adjustment model is used outside of the population in which it is
1042 calibrated. Measure developers have the burden of demonstrating appropriateness through R-

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

1043 squared values and through a detailed clinical and statistical explanation of how variables were
1044 added to the risk model. Additional research is needed in this area to explore how various risk-
1045 adjustment approaches change the relative ranking of providers in terms of resource use and how
1046 the use of clinically enhanced administrative data may impact measure scores and the selection of
1047 factors added to the risk-adjustment models.

1048 *Emerging Principle 3: The factors in the risk-adjustment model and severity model should be*
1049 *confirmed to be a contributor to the outcome of the measure.*

1050 *Emerging Principle 4: In addition to statistical significance, justification of the variables used in*
1051 *the risk-adjustment model should be provided based on either clinical relevance or evidence in*
1052 *the literature.*

1053 The Committee agreed that measure developers need to demonstrate that variables included in
1054 the risk-adjustment model are not simply selected based on their statistical explanatory power,
1055 but rather, risk factors are well documented in clinical evidence. When variables are chosen for
1056 inclusion in the risk-adjustment model, developers are responsible for demonstrating a
1057 relationship to the outcome of the measure (i.e., resource use). Additional detail through a
1058 sensitivity analysis including various risk-adjustment variables can be provided in future
1059 evaluations to demonstrate the effect of variables included in the final risk-adjustment approach.

1060

1061 **2) Reliability and Validity Testing**

1062 The cumulative experience of the multiple TAPs and the Resource Use Steering Committee
1063 demonstrated that resource use measures developers are at various levels of measure testing
1064 sophistication.

1065

1066 Measures submitted as resource use national consensus standards need to improve their level of
1067 sophistication of reliability and validity testing. To balance the developer burden of testing for
1068 the initial evaluation of resource use measures with providing the experts the information needed
1069 to make a valid conclusion about reliability and validity, the TAPs and Steering Committee
1070 agreed that the scope of testing may be on a relatively small scale for initial endorsement. The

NATIONAL QUALITY FORUM

1071 Committee agreed further analysis by all developers would be required to support continued
1072 endorsement at the time of review in order to maintain NQF endorsement.

1073
1074 Reliability and validity testing is included in the NQF evaluation criteria, and NQF allows
1075 flexibility in the specific methods used in testing to allow measure developer flexibility. The
1076 Committee evaluated: 1) the scope of testing, 2) what tests of reliability and validity could be
1077 performed, and 3) how to weigh the results of this testing. The Steering Committee interpreted
1078 testing results within the unique context of the specific measure under review.

1079
1080 ***Reliability testing***

1081 The NQF evaluation criteria states that reliability testing should demonstrate that the data
1082 elements are repeatable, producing the same results a high proportion of the time when assessed
1083 in the same population in the same time period, or that the measure score is precise. The
1084 Committee agreed that developers can demonstrate that the measure score is precise by
1085 demonstrating an adequate ratio of signal to noise, or how well one can confidently distinguish
1086 the performance of one physician from another.¹² The signal is ability of the measure to identify
1087 real differences in performance, whereas the noise attributed to measurement error.

1088 Demonstrating reliability in this context relies on three major drivers: sample size, differences
1089 among physicians, and random variation in the measure scores, or measurement error.¹³

1090 *Emerging Principle 5: To demonstrate reliability of a resource use measure, developers can*
1091 *focus on precision of the measure score.*

1092 Reliability at the data element level of resource use measures submitted to this project relied on
1093 administrative claims and by virtue of their design as coded programs were repeatable. However,
1094 the Committee clarified that while coded programs may be repeatable at the data element level,
1095 measure developers need to demonstrate adequate validity testing at the data element level.

1096 *Emerging Principle 6: When there is such limited variability in a data set that it does not*
1097 *adequately distinguish performance differences among providers, reliability cannot simply rely*
1098 *on confidence intervals; sample size should also be included in the reliability assessment.*

NATIONAL QUALITY FORUM

1099 Reliability of resource use measures at the measure score level needs to demonstrate that the
1100 measure score is precise. Providing confidence intervals in measure reporting does not
1101 sufficiently demonstrate reliability of the measure.

1102
1103 NQF does not prescribe what tests of reliability could be performed, specific thresholds for
1104 results, or how to weigh the results of this testing since an evaluation should account for the
1105 context of the test, measure, and the data source. The evaluation should incorporate both
1106 empirical evidence and expert judgment to evaluate whether the specific measure under
1107 evaluation by the Committee has sufficiently demonstrated reliability through the measure
1108 submission.

1109

1110 *Validity testing*

1111 *Emerging Principle 7: The gold standard approach to determining the validity of data elements*
1112 *based on administrative claims data in resource use measures is to assess the agreement of*
1113 *claims data with source of the data elements in the chart.*

1114 Since the entire dataset may not be available for such validation, applying the resource use
1115 measure to a simulated data set that should return known values of the data elements and scores
1116 may be used. With either approach, when the results obtained for the resource use measure do
1117 not match known values in the simulated data set or the abstracted data, an analysis should be
1118 conducted to determine the source of error.¹⁵ If the error is related to the measure specifications,
1119 including code lists, clinical or construction logic, and computer readable programming
1120 language, the measure specification should be corrected before submitting for endorsement.

1121

1122 The NQF criteria state that validity testing must demonstrate that the measure data elements are
1123 correct or that the measure score correctly reflects the cost of care or resources provided,
1124 adequately distinguishing high and low resource use. Developers must demonstrate measures
1125 have undergone sufficient validity testing demonstrating that the resource use measure actually
1126 measures what it claims to measure. If only face validity is addressed, it must be assessed
1127 systematically. The Committee recommended that validity testing be demonstrated by

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

1128 correlating measure scores with other valid indicators or by showing that the score produces
1129 different results when applied to subgroups known to have differences in resource use. Correct
1130 conclusions about resource use can be made when validity tests demonstrate that claims used in
1131 the measure accurately reflect information in the charts of a representative sample of patients.
1132 The Committee considered that most developers submitting to this project may not have direct
1133 access to chart abstracted data; however, additional efforts are strongly recommended to ensure
1134 data elements used to develop the resource use measures are valid.

1135

1136 **3) Data Quality and Comprehensiveness**

1137 In an effort to address some of the underlying global issues affecting the use of administrative
1138 claims data for the purposes of measuring resource use, the Committee identified several areas in
1139 which healthcare stakeholders might engage and support additional efforts to improve the ability
1140 of resource use measures to capture all resource use fully. In doing so, a few new principles for
1141 resource use measurement emerged.

1142 *Emerging Principle 8: Data sets used to measure resources should be as comprehensive as*
1143 *possible. Efforts to obtain clinical and carved-out data (e.g., pharmacy, behavioral health) should*
1144 *be made to ensure the data set used to calculate resource use is robust, complete, and*
1145 *representative.*

1146 *Emerging Principle 9: Measure scores calculated and reported using data with carve-outs*
1147 *should be labeled as such.*

1148 *Emerging Principle 10: Comparisons of entities with and without carved-out data is inappropriate.*

1149 A major concern of the Committee throughout the evaluation process was the impact of carve-
1150 out arrangements on accurately capturing resources used. While there are some systems that are
1151 able to recapture these data from the outsourced entities, others do not have this capability.

1152 Furthermore, measure results derived for entities with carve-out arrangements should be labeled
1153 as such to prevent comparison between entities with and without such carve-out arrangements.

1154 The measures received during this project were specified using administrative claims data only.

1155 The use of administrative claims data presents certain limitations for measuring resource use
1156 performance, limitations that are present in quality performance measurement as well. Primarily

NATIONAL QUALITY FORUM

1157 the reliance of resource use measures on administrative claims data to count resources, or dollars
1158 spent, captures only the output on behalf of the provider—not the costs to the patient, nor the
1159 costs or resources for which there are no administrative codes. Recognizing this as a limitation of
1160 the data available to measure these types of resources, the Committee recommended that future
1161 efforts in resource use measurement focus not only on the costs to the provider, but to the user as
1162 well, through identifying those resources that are important to measure and determining how to
1163 capture this data. The Committee recognized that while administrative data are the primary data
1164 source used for measuring resources at this time, there is opportunity to integrate the data
1165 gathered through EHRs and other clinical data to measure resource use.

1166
1167 Since resource use measurement is a priority, efforts should be made to ensure the necessary data
1168 are available for accurate measurement. However, there are significant challenges to determining
1169 where the responsibility lies to ensure data are complete and the ways in which important but
1170 sensitive information is shared. For a number of measures submitted for evaluation in this
1171 project, the instructions within the data protocol module suggest that the measure implementer is
1172 responsible for ensuring data are complete and representative. The Committee agreed however
1173 that measure implementers often do not have the resources or technical expertise to audit data
1174 before use. Future efforts should explore a potential role for large data aggregators to identify
1175 thresholds and set standards for data quality.

1176
1177 *Emerging Principle 11: If a measure is intending to measure a clinical condition that has a*
1178 *predominant portion of its costs in pharmacy claims, consider whether costs should be measured at*
1179 *all in the absence of these data. The developer is responsible for determining whether the lack of*
1180 *these data invalidates the measure score or comparisons.*

1181 When developing resource use measures, careful consideration should be given to whether the
1182 importance of measuring resources/costs in an area outweigh the limitations of the data. For
1183 some conditions, the lack of robust data could distort the measure output. For example, to
1184 measure the resources for asthma patients where greater than 40 percent of the resource use is
1185 pharmacy related, data sets without pharmacy data are inherently misleading in providing useful

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

1186 insight into the cost of asthma care.¹⁶ For acute or procedural episodes (e.g., hip replacement)
1187 where the care is more standardized (e.g., pre- and post-surgical antibiotics) pharmacy and
1188 mental health data do not account for a major portion of the resource use and thus administrative
1189 data, and carve-out issues may not have a tremendous impact on the measure results.

1190

1191 **4) Measuring Cost and Resource Use in the Medicare Population**

1192 Measures evaluated in this project were mainly specified for the commercial population;
1193 however, the Steering Committee identified areas of consideration for organizations developing
1194 resource use measure for Medicare beneficiaries. Resource use measures for the Medicare
1195 population will have to consider multiple co-occurring conditions, as well as multiple sites where
1196 beneficiaries seek care and resource use at the end of life. This guidance was provided with the
1197 understanding that there will be an urgent need for measures specified for the Medicare
1198 population for use in bundled payment demonstrations, physician feedback reporting programs,
1199 and value-based purchasing programs.

1200 An important consideration in resource use approaches developed for the Medicare population is
1201 the presence of multiple co-occurring conditions. The Committee considered that more than half
1202 of all beneficiaries were treated for five or more conditions, accounting for three-fourths of total
1203 Medicare spending.¹⁷ In 2002, more than 92.2 percent of all Medicare healthcare spending was
1204 incurred by beneficiaries with three or more conditions during the measurement year.¹⁸ In the
1205 Committee's discussion of the approach to measure resource use in patients with multiple co-
1206 occurring conditions, they concluded that cost estimates should be based on the time and
1207 attention a provider should be reasonably expected to deliver on a patient's multiple co-occurring
1208 conditions beyond the acute disease and its immediate complications for which the patient
1209 sought care.

1210 *Emerging Principle 12: A patient-centered approach should be used to describe the interaction*
1211 *of conditions (and episodes) in the development of resource use measures for the Medicare*
1212 *population.*

NATIONAL QUALITY FORUM

1213 Episode approaches attempting to assign claims to specific episodes should create a transparent
1214 hierarchy with rules to assess resource use in the Medicare population accurately. One approach
1215 the Committee suggested would allow flexibility in the assignment of individual claims to a
1216 single episode or to multiple open episodes. This patient-centered approach could allow an
1217 individual office visit for evaluation and management to be assigned to multiple episodes.

1218 Efforts to develop resource use measures for the Medicare population should consider the NQF
1219 consensus measure framework for assessing the efficiency of care for individuals with multiple
1220 chronic conditions (MCCs). MCC framework guiding principles include promoting shared
1221 accountability with members of the healthcare system, a multi-dimensional measure approach
1222 that incorporates various types of measures, a focus on shared decision making in concordance
1223 with a patient's preferences, and prioritization of measures across time that are most relevant to
1224 achieving desired outcomes as determined by the care plan.

1225 The Committee recognized the cost contribution of individual conditions to the total cost of
1226 managing a beneficiary may vary differently depending on other conditions present in each
1227 beneficiary. The following classification system of four types of overlapping episodes helps to
1228 illustrate the patterns of treatment discussed by the Committee: (1) linear additive episodes, (2)
1229 interactive episodes-cost increasing, (3) interactive episodes-cost savings, and (4) dominant
1230 episodes.¹⁹ Linear additive episodes occur when the patterns of illness are not overlapping, and
1231 episodes can be considered independent of one another. For example, a fracture of the radius
1232 and strep throat would be considered independent of one another.²⁰ Interactive episodes can be
1233 cost increasing when there are two or more conditions in which the presence of multiple
1234 conditions increases the level of resources required to treat all of the conditions. An example
1235 would include the treatment of diabetes in the presence of obesity.²¹ Under this condition, the
1236 cost of the combined condition is more than the sum of the individual parts. Interactive episodes
1237 can also be cost saving since the cost of treating overlapping conditions is not likely to require
1238 significantly different resources (e.g., the treatment of otitis media and bronchitis).²² Finally,
1239 dominant and mild disease combinations in which the presence of a dominant disease episode
1240 becomes the principle focus of care (e.g., the treatment of end-stage renal disease in the presence

NATIONAL QUALITY FORUM

1241 of mild asthma). These methods for overlapping episodes should be considered in developing
1242 approaches for assessing resource use in the Medicare population.

1243 The nature of the interaction between chronic and acute conditions should be considered when
1244 developing resource use measures. When developing measures to assess the resource use for a
1245 chronic condition, the resource use for an acute complication for that condition should be
1246 considered. The Committee considered the example of misinterpreting lower CAD resource use
1247 as better performance when, in fact, a per-capita assessment may demonstrate higher resource
1248 use. The higher resource use may be derived from higher rates of AMI in the measured
1249 population due to poor CAD management.

1250 Additional efforts are needed to propose alternative attribution approaches to encourage team-
1251 based care along the patient episode of care. Resource use measures developed for the Medicare
1252 population should also consider that beneficiaries often seek care from multiple sites. The typical
1253 Medicare beneficiary sees two primary care physicians and five specialists working in four
1254 different practices.²³ The Committee discussed how current attribution models assign treatment
1255 of the patient to an individual provider based on the number of visits or the highest proportion of
1256 costs. However, in a patient-centered model all providers who treat the patient should have
1257 responsibility for the care delivered.

1258 Episode-based approaches for the Medicare population should carefully consider their approach
1259 to dealing with end of life (EOL). Simply including EOL patients in estimates of episode-based
1260 resource use has the potential to introduce inappropriate incentives. Resource use measures that
1261 include EOL patients should be reported with balancing mortality measures to ensure that
1262 providers are not inadvertently reported as providing more efficient care when they have higher
1263 rates of mortality. On the other hand, with resource use during the last year of life accounting for
1264 more than a quarter of Medicare payments,²⁴ EOL patients should not be excluded from the
1265 analysis of resource use. Future evaluation of resource measures for the Medicare population
1266 should consider how measure developers handle EOL patients in profiling providers.

1267

1268 **5) Linking Quality and Cost to Develop Measures of Efficiency and Value**

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

1269 Developing measures of efficiency and value is critical to reducing the healthcare cost growth
1270 rate. In a first step toward developing efficiency measures, resource use measures must
1271 demonstrate they are important to measure, have scientifically acceptable properties, and are
1272 usable and feasible. Resource use measures that meet these criteria may be used in conjunction
1273 with quality measures to assess efficiency. The Steering Committee reflected on the mechanism
1274 and future work needed to achieve this goal.

1275 *Emerging Principle 13: Efficiency measurement approaches should be patient-centered, building*
1276 *upon previous efforts such as the NQF Patient-Centered Episodes of Care (EOC) Efficiency*
1277 *Framework.*

1278 Measures components may need to be aligned between quality and resource use measures.
1279 Components that may be aligned include the handling of exclusions, level of analysis, risk
1280 adjustment, and stratification approach. For example, the Committee recommended that quality
1281 and resource use measures be aligned in terms of their inclusion and exclusion criteria to ensure
1282 similar populations are being measured in both the resource use and quality performance.

1283 The Committee recommended future work to define the type of quality and resource use
1284 measures that can be used to assess quality. Considerations should include the measure type
1285 (e.g., outcome, process, patient experience), measurement period (e.g., single point in time,
1286 spanning the measurement year), and the number of quality measures that should be paired with
1287 a resource use measure. The Committee also considered that quality measures may be used to
1288 monitor for underuse on needed care. Assessments of efficiency will require careful
1289 consideration of the mechanism in which quality and resource use measures are linked.

1290 Future efforts should explore approaches to ensure that providers are benchmarked on cost
1291 performance against providers with similar or better quality performance. Benchmarking cohorts
1292 of providers based on quality performance allows for accurate interpretation of cost. Specifically
1293 this method ensures that the resource use performance is compared to only those providers with
1294 equal or higher quality performance.^{25,26} When available, the Committee agreed that outcome
1295 and patient experience of care measures with sufficient reliability (signal to noise) and validity
1296 should be selected to assess efficiency.²⁷

1297

NATIONAL QUALITY FORUM

1298 NEXT STEPS

1299 This project enabled first-hand experience in reviewing and understanding some of the various
1300 approaches for measuring resources and costs in healthcare, and while many lessons were
1301 learned, there is still abundant opportunity to apply the principles and recommendations that
1302 emerged from this work in future efforts. Ongoing work in the public sector to develop a public
1303 episode grouper for the Medicare population and explore ways to measure efficiency using a
1304 patient-centered approach will be the focus of future NQF efforts in this area. Additionally, using
1305 the recommendations from the Committee on improving the evaluation process, updates to the
1306 NQF resource use measure submission forms and evaluation criteria will be explored as we
1307 continue to enhance the endorsement process for measure submitters and evaluators.

1308 NOTES

- 1309 1. Catlin A, et al., National Health Spending in 2006: A Year of Change for Prescription Drugs,
1310 *Health Affairs*, 2008; 27(1):14–29.
1311
- 1312 2. Banks J, et al., Disease and disadvantage in the United States and in England, *JAMA*,
1313 2006;295(17):2037–2045.
1314
- 1315 3. Hoyert DL, et al., Annual summary of vital statistics: 2004, *Pediatrics*, 2006; 117(1):168–
1316 183.
1317
- 1318 4. Weiss JE, Mushinski M, International mortality rates and life expectancy: selected countries,
1319 *Statistical Bulletin—Metropolitan Life Insurance Company*, 1999;80(1):13–21.
1320
- 1321 5. Catlin A, et al.
1322
- 1323 6. CMS Grouper Episode for Medicare. Federal Business Opportunities. Obtained from:
1324 [https://www.fbo.gov/index?s=opportunity&mode=form&id=452b7af0f5752ba0c444853793e](https://www.fbo.gov/index?s=opportunity&mode=form&id=452b7af0f5752ba0c444853793e8024d&tab=core&_cview=1)
1325 [8024d&tab=core&_cview=1](https://www.fbo.gov/index?s=opportunity&mode=form&id=452b7af0f5752ba0c444853793e8024d&tab=core&_cview=1). Last Accessed August 2011.
- 1326 7. Institute of Medicine (IOM), *Crossing the Quality Chasm: A New Health Systems for the 21st*
1327 *Century*, Washington, DC: National Academies Press; 2001.
- 1328 8. AQA Principles of “Efficiency” Measures April 2006. Available from:
1329 <http://www.aqaalliance.org/performancewg.htm>. Last accessed August 2011.
1330

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

- 1331 9. Atlas RF, The role of PBMs in implementing the Medicare prescription drug benefit, *Health*
1332 *Affairs*, 2004; Jul-Dec;Suppl Web Exclusives:W4-504-15. Review.
- 1333 10. Grazier K L and Eselius LL, Mental health carve-outs: effects and implications, *Medical*
1334 *Care Research and Review*, 1999;56(suppl 2):37.
- 1335 11. Hasings C, Mosteller F, Tukey JW, Winsor CP, Low moments for small samples: a
1336 comparative study of order statistics, *Annals of Mathematical Statistics*, 1947;18:413–426.
- 1337 12. Scholle, S. H., J. Roski, et al. (2008). "Benchmarking physician performance: reliability of
1338 individual and composite measures." *The American Journal of Managed Care* 14(12): 833.
- 1339 13. Adams J, McGlynn E, et al., *Incorporating statistical uncertainty in the use of physician cost*
1340 *profiles*, *BMC Health Services Research*, 2010;10(1): 57.
- 1341
- 1342 14. Ibid.
- 1343
- 1344 15. [NQF Testing Task Force Report](#) (2011).
- 1345
- 1346 16. Bahadori K, Doyle-Waters, et al., Economic burden of asthma: a systematic review, *BMC*
1347 *Pulmonary Medicine*, 2009;9(1):24.
- 1348
- 1349 17. Thorpe KE and Howard DH, The rise in spending among Medicare beneficiaries: the role of
1350 chronic disease prevalence and changes in treatment intensity, *Health Affairs*,
1351 2006;25(5):w378.
- 1352
- 1353 18. Ibid..
- 1354
- 1355 19. Hornbrook M, Hurtado A, et al., Health care episodes: definition, measurement and use,
1356 *Medical Care Research and Review*, 1985;42(2):163.
- 1357
- 1358 20. Hornbrook, et al.
- 1359
- 1360 21. Ibid.
- 1361
- 1362 22. Ibid.
- 1363
- 1364 23. Pham HH, Schrag D, et al., Care patterns in Medicare and their implications for pay for
1365 performance, *New Eng J Med*, 2007;356(11): 1130-1139.
- 1366
- 1367 24. Hogan C, Lunney J, et al., Medicare beneficiaries' costs of care in the last year of life, *Health*
1368 *Affairs*, 2001;20(4):188.
- 1369

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

1370 25. Tompkins CP, Higgins AR, et al., Measuring outcomes and efficiency in Medicare value-
1371 based purchasing, *Health Affairs*, 2009;28(2):w251.

1372
1373 26. Chung J, Kaleba E, Wozniak G, A Framework for Measuring Healthcare Efficiency and
1374 Value. PCPI Work Group on Efficiency and Cost of Care. Aug 2008. White Paper.

1375 27. Hussey PS, De Vries H, et al., A systematic review of health care efficiency measures,
1376 *Health Services Research*, 2009;44(3):784-805.

1377

1378

1379

1380

1381

1382

1383

1384

1385

1386

1387

1388

1389

1390

1391

1392

1393

1394

1395

1396

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

1397 **APPENDIX A—SPECIFICATIONS FOR COST AND RESOURCE USE MEASURES**
 1398 **2011 (Cycle 2)**

1399
 1400 The following tables present the detailed measure specifications for the recommended consensus
 1401 standards. All information summarized here has been derived directly from the measure
 1402 developers without modification or alteration (except where measure developers agreed to such
 1403 modifications) and is current as of August 15, 2011. All proposed voluntary consensus standards
 1404 are open source, meaning they are fully accessible and disclosed.

1405 **Pulmonary**

1406 (1560) Relative Resource Use for People with Asthma (NCQA).....77
 1407 (1561) Relative Resource Use for People with COPD (NCQA).....79
 1408 (1611) ETG Based Pneumonia cost of care (Ingenix).....81

1409 **Bone/Joint**

1410 (1609) ETG/PEG Based Hip/Knee Replacement cost of care measure (Ingenix).....83

1411

1560: Relative Resource Use for People with Asthma	
Steward	NCQA
Description	The risk-adjusted relative resource use by health plan members with asthma during the measurement year.
Resource Use Measure Type	Per capita (population- or patient-based)
Data Source	Administrative claims Electronic Clinical Data : Electronic Health Record, Imaging/Diagnostic Study, Laboratory, Pharmacy Paper Records
Level of Analysis	Clinician : Group/Practice Health Plan Integrated Delivery System Population : National, Regional
Clinical Framework Description	2 Eligibility Criteria: An encounter with a diagnosis; or by multiple asthma medication events. An organization must use both methods to identify the eligible population, but a member only needs to be identified by one to be included in the measure. To identify the eligible population for measurement: Step 1: Health Plan members are identified as having persistent asthma by meeting at least one of the following criteria during both the measurement year and the year prior to the measurement year. Criteria need not be the same across years. <ul style="list-style-type: none"> • At least one ED visit (Tables ASM-A and ASM-B) with asthma as the principal diagnosis, or • At least one acute inpatient claim/encounter (Tables ASM-B) with asthma as the principal diagnosis (Table ASM-A), or • At least four outpatient asthma visits (Table ASM-B) with asthma as one of the listed diagnoses (Table

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

	1560: Relative Resource Use for People with Asthma
	<p>ASM-A) and at least two asthma medication dispensing events (Table ASM-C), or</p> <ul style="list-style-type: none"> • At least four asthma medication dispensing events (Table ASM-C) <p>Step 2: Since a member can be identified as having persistent asthma using only leukotriene modifiers as the sole asthma medication dispensed in that year, these members must also have at least one diagnosis of asthma (Table ASM-A), in any setting, in the same year as the leukotriene modifier prescription (e.g. measurement year or year prior to the measurement year).</p> <p>Exclusions:</p> <ol style="list-style-type: none"> 1) Active cancer. Exclude members who had at least one face-to-face encounter, in any setting, with any diagnosis of cancer in conjunction with any treatment code (Table RRU-A), during the measurement year. 2) ESRD. Exclude members who had at least one face-to-face encounter with any code to identify ESRD (Table RRU-B), during the measurement year. 3) Organ transplant. Exclude members who had at least one face-to-face encounter, in any setting, with any code to identify organ transplant (Table RRU-C), during the measurement year. 4) HIV/AIDS. Exclude members who had at least two face-to-face encounters in an outpatient or nonacute inpatient setting, or at least one face-to-face encounter in an acute inpatient or ED setting, with any diagnosis of HIV (Table RRU-D), with different dates of service during the measurement year. Refer to Table RRU-E for codes to identify visit type. 5) Members diagnosed with emphysema, COPD, cystic fibrosis or acute respiratory failure (Table ASM-E) on or prior to December 31 of the measurement year.
Costing Method	<p>RRU measures use NCOA's standardized prices. The organization does not report prices based on its contracts and fee schedules, rather it applies a standard price to each service, multiplies it by the number of units of service and reports the resulting standard cost. The standard pricing approach is based on the following sources of data:</p> <ul style="list-style-type: none"> • Relative values from the Medicare Fee Schedule (Resource-Based Relative Value Scale, or RBRVS) • Pharmacy prices published by First Bank Data • Inpatient prices based on a model that uses a broad set of averages, representing different local, regional and national health plans across the country. <p>A plan maps a standard price to each service, multiplies it by the number of units of service and reports the resulting standard cost. It then calculates total standard costs for eligible members across different areas of clinical care and aggregates standard costs across services and members to compute the overall relative resource use.</p> <p>All RRU measures report the standard cost for the following categories.</p> <ul style="list-style-type: none"> o Inpatient Facility o Surgery and Procedure o Inpatient Services o Outpatient Services o Evaluation and Management (E&M) o Inpatient Services o Outpatient Services o Diagnostic Laboratory Services o Diagnostic Imaging Services o Pharmacy, Ambulatory
Tested Population	<p>Commercial; Medicaid</p>
Resource Use Service Categories	<p>Inpatient services: Inpatient facility services; Evaluation and management; Procedures and surgeries; Imaging and diagnostic; Lab services; Admissions/discharges</p> <p>Ambulatory services: Outpatient facility services; Emergency Department; Pharmacy; Evaluation and management; Procedures and surgeries; Imaging and diagnostic; Lab services</p>

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

1560: Relative Resource Use for People with Asthma	
Attribution Approach	Specifications: Relative resource use is calculated at the plan-level and no attribution of resource use is made below this level. Attribution of resource use to a particular NCQA submission is based on the product line and reporting type of the plan that the member was enrolled in as of the end of the measure year.
Risk Adjustment	The current risk model utilized by NCQA is based on components of the CMS-HCC risk adjustment methodology and accounts for age, gender, and HCC-RRU risk classifications that predict cost variability. For each condition, members are assigned to a clinical cohort category that provides a more specific classification of the condition. A members age, gender, and HCC category determines their risk score (cohort). NCQA then calculates the average per-member per-month (PMPM) cost for each cohort then weights that cost by the total member months within each cohort. Each plan will have its own weight for each cohort since case-mix varies across plans. These weighted cohort PMPMs are then summed across all cohorts to estimate total resource use that would be expected if the “average” plan had the same case-mix as the plan in question. The ratio of the observed- to-expected PMPM utilization indicates the degree to which a plan deviates from expected performance. This is known as indirect standardization.
Stratification	NCQA collects resource measures at the plan level and summarizes across reporting cohorts along the following dimensions: Product line (3 levels): Commercial, Medicaid, and Medicare; Reporting type (2 levels): HMO and PPO; Area level (2 levels): national and region; Resource use or utilization (11 levels): inpatient facility, procedure and surgery (inpatient and outpatient), evaluation and management (inpatient and outpatient), laboratory services, imaging services, ambulatory pharmacy, inpatient discharges, emergency department discharges. Stratification of RRU results to control for individual confounding variables is not performed since age, gender and risk variables (comorbidity and disease interactions) that affect healthcare costs are accounted for in the RRU-HCC risk adjustment process. These include age and gender along with one of the 13 assigned HCC-RRU risk categories (e.g. male 18-44 HCC-RRU 1; male 18-44 HCC-RRU 2; male 18-44 HCC-RRU 3; etc...).

1412

1561: Relative Resource Use for People with COPD	
Steward	NCQA
Description	The risk-adjusted relative resource use by health plan members with COPD during the measurement year.
Resource Use Measure Type	Per capita (population- or patient-based)
Data Source	Administrative claims Electronic Clinical Data: Electronic Health Record, Imaging/Diagnostic Study, Laboratory, Pharmacy Paper Records
Level of Analysis	Clinician : Group/Practice Health Plan, Integrated Delivery System, Population : Community, National, Regional
Clinical Framework Description	Members are identified for the eligible population of the measure with a diagnosis of COPD (Table SPR-A) present anytime during the measurement year and who were continuously enrolled for a two year period (the measurement year and the year prior). Codes to Identify COPD: Chronic bronchitis-ICD-9 Diagnosis: 491 Emphysema -ICD-9 Diagnosis: 492 COPD -ICD-9 Diagnosis: 496 Exclusions: 1) Active cancer. Exclude members who had at least one face-to-face encounter, in any setting, with any diagnosis of cancer in conjunction with any treatment code (Table RRU-A), during the measurement year.

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

	1561: Relative Resource Use for People with COPD
	<p>2) ESRD. Exclude members who had at least one face-to-face encounter with any code to identify ESRD (Table RRU-B), during the measurement year.</p> <p>3) Organ transplant. Exclude members who had at least one face-to-face encounter, in any setting, with any code to identify organ transplant (Table RRU-C), during the measurement year.</p> <p>4) HIV/AIDS. Exclude members who had at least two face-to-face encounters in an outpatient or nonacute inpatient setting, or at least one face-to-face encounter in an acute inpatient or ED setting, with any diagnosis of HIV (Table RRU-D), with different dates of service during the measurement year. Refer to Table RRU-E for codes to identify visit type.</p> <p>5) Members diagnosed with emphysema, COPD, cystic fibrosis or acute respiratory failure (Table ASM-E) on or prior to December 31 of the measurement year.</p>
Costing Method	<p>RRU measures use NCOA's standardized prices. The organization does not report prices based on its contracts and fee schedules, rather it applies a standard price to each service, multiplies it by the number of units of service and reports the resulting standard cost. The standard pricing approach is based on the following sources of data:</p> <ul style="list-style-type: none"> • Relative values from the Medicare Fee Schedule (Resource-Based Relative Value Scale, or RBRVS) • Pharmacy prices published by First Bank Data • Inpatient prices based on a model that uses a broad set of averages, representing different local, regional and national health plans across the country. <p>A plan maps a standard price to each service, multiplies it by the number of units of service and reports the resulting standard cost. It then calculates total standard costs for eligible members across different areas of clinical care and aggregates standard costs across services and members to compute the overall relative resource use.</p> <p>All RRU measures report the standard cost for the following categories.</p> <ul style="list-style-type: none"> o Inpatient Facility o Surgery and Procedure o Inpatient Services o Outpatient Services o Evaluation and Management (E&M) o Inpatient Services o Outpatient Services o Diagnostic Laboratory Services o Diagnostic Imaging Services o Pharmacy, Ambulatory
Tested Population	Commercial; Medicaid; Medicare
Resource Use Service Categories	<p>Inpatient services: Inpatient facility services, Evaluation and management, Procedures and surgeries, Imaging and diagnostic, Lab services</p> <p>Admissions/discharges</p> <p>Ambulatory services: Outpatient facility services, Emergency Department; Pharmacy, Evaluation and management, Procedures and surgeries, Imaging and diagnostic, Lab services</p>
Attribution Approach	Specifications: Relative resource use is calculated at the plan-level and no attribution of resource use is made below this level. Attribution of resource use to a particular NCOA submission is based on the product line and reporting type of the plan that the member was enrolled in as of the end of the measure year.
Risk Adjustment	The current risk model utilized by NCOA is based on components of the CMS-HCC risk adjustment methodology and accounts for age, gender, and HHC-RRU risk classifications that predict cost variability. For each condition, members are assigned to a clinical cohort category that provides a more specific classification of the condition. A members age, gender, and HCC category determines their risk score (cohort). NCOA then calculates the average per-member per-month (PMPM) cost for each cohort then weights that cost by the total member months within each cohort. Each plan will have its own weight for each

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

	1561: Relative Resource Use for People with COPD
	cohort since case-mix varies across plans. These weighted cohort PMPMs are then summed across all cohorts to estimate total resource use that would be expected if the “average” plan had the same case-mix as the plan in question. The ratio of the observed- to-expected PMPM utilization indicates the degree to which a plan deviates from expected performance. This is known as indirect standardization.
Stratification	<p>NCQA collects resource measures at the plan level and summarizes across reporting cohorts along the following dimensions:</p> <p>Product line (3 levels): Commercial, Medicaid, and Medicare;</p> <p>Reporting type (2 levels): HMO and PPO;</p> <p>Area level (2 levels): national and region;</p> <p>Resource use or utilization (11 levels): inpatient facility, procedure and surgery (inpatient and outpatient), evaluation and management (inpatient and outpatient), laboratory services, imaging services, ambulatory pharmacy, inpatient discharges, emergency department discharges.</p> <p>Stratification of RRU results to control for individual confounding variables is not performed since age, gender and risk variables (comorbidity and disease interactions) that affect healthcare costs are accounted for in the RRU-HCC risk adjustment process. These include age and gender along with one of the 13 assigned HCC-RRU risk categories (e.g. male 18-44 HCC-RRU 1; male 18-44 HCC-RRU 2; male 18-44 HCC-RRU 3; etc...).</p>

1413

1414

	1609: ETG Based hip/knee replacement cost of care measure
Steward	Ingenix
Description	<p>This submission is for Hip/Knee Replacement procedure episodes and will cover both measures at the Hip Replacement and Knee Replacement PEGs. The measure focuses on resources used to deliver episodes of care for patients who have undergone a hip or knee replacement and assigns a level of severity (e.g., resources per episode for Knee Replacement, severity level 1 episodes). Hip Replacement and Knee Replacement episodes are initially defined using the Episode Treatment Groups (ETG) methodology and describe the unique presence of the condition for a patient and the services involved in diagnosing, managing and treating the condition. The Procedure Episode Group (PEG) methodology uses the ETG results and further logic to creating a procedure episode that focuses on the Hip Replacement and Knee Replacement component of the care. Procedure episodes identify a unique procedure event as well as the related services performed before and after the procedure including workup and therapy prior to the procedure as well as post-op activities such as repeated surgery and patient follow-up. Together, the ETG and PEG methodologies identify the services involved in diagnosing, managing and treating patients with Hip/Knee Replacements. A methodology to assign a severity level to each episode is employed to group Hip and Knee Replacement episodes by level of risk.</p> <p>Multiple types of resources can be measured for Hip/Knee Replacement episodes, including overall cost of care, cost of care by type of service, and the utilization of specific types of services. Each resource use measure is expressed as a cost or a utilization count per episode and comparisons with internal and external benchmarks are made using risk adjustment to support valid comparisons.</p>
Resource Use Measure Type	Per episode
Data Source	Administrative claims
Level of Analysis	<p>Clinician : Group/Practice, Individual, Team, Facility, Health Plan, Integrated Delivery System</p> <p>Population : Community, County or City, National, Regional, State</p>
Clinical Framework Description	This measure identifies patients with Hip/Knee Replacement and creates Hip/Knee Replacement episodes of care using the ETG and PEG methodologies described in the ETG_PEG Construction Logic attached in our response to S.2. Each procedure episode of Hip/Knee Replacement is characterized by a PEG Anchor Category ID that specifies the type of procedure; the PEG Anchor Category ID representing Hip Replacement

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

	1609: ETG Based hip/knee replacement cost of care measure
	<p>is 71518 and the PEG Anchor Category ID representing Knee Replacement is 71918.</p> <p>An ETG/PEG episode of Hip/Knee Replacement will contain all clinically relevant information related to the procedure. The Hip/Knee Replacement episode clinical framework is defined by the services, or claim lines, that can begin an episode, the primary and incidental diagnosis relationships involved and how records group to an episode, including relative strength of relationship.</p>
Costing Method	The financial amounts used should be complete and valid, reflecting the total payments related to the service. The financial amount used in resource measurement should reflect all payments for a service, including those made to the provider by payer, patient and other entities. Allowed payments will reflect both the quantity of different services provided as well as the actual unit price of those same services.
Tested Population	Commercial
Resource Use Service Categories	<p>Inpatient services: Inpatient facility services, Admissions/discharges</p> <p>Ambulatory services: Outpatient facility services, Emergency Department, Pharmacy, Evaluation and management, Procedures and surgeries, Imaging and diagnostic, Lab services</p>
Attribution Approach	<p>Guidelines: For physician measurement, the primary surgeon is typically attributed the episode, although applications of attribution could be developed to support an alternate approach. Both activity-based and population-based approaches should be supported. As a guideline, four different general options for physician episode attribution can be considered to attribute episodes to individual providers – three activity-based and one population-based approach.</p> <p>Approach 1 - Physician Episode Attribution using Professional Service Costs. This attribution approach identifies the responsible physician for an episode as that provider rendering the greatest amount of professional service costs during the episode.</p> <p>Approach 2 - Physician Episode Attribution using Episode Clusters. This attribution approach identifies the responsible physician for an episode as that provider in the peer group owning the greatest number of “clusters” within the episode.</p> <p>Approach 3 - Physician Episode Attribution using Non-Acute Evaluation and Management (E/M) Visits. This attribution approach identifies the responsible physician for an episode as that physician providing the greatest number of non-acute E/M visits within the episode.</p> <p>Approach 4 - Physician Episode Attribution using a Primary Care, Population-based Approach. This approach requires two important steps: 1) Identification of a PCP for each member. 2) Identify the patient's assigned PCP during the episode period.</p>
Risk Adjustment	<p>The level of severity assigned to an episode is used to support risk adjustment. The risk adjustment approach includes three important steps:</p> <ol style="list-style-type: none"> 1. Compute the observed experience for the provider being measured, across all episodes to be included in the comparison; 2. Compute the experience for peers or a best practice benchmark. Compute this experience at the level of the risk adjustment, in this case base procedure (hip or knee replacement) and severity level. For a peers benchmark, average cost per episode across all peers for the base procedure and severity level can be computed; 3. Compare the observed experience with the risk adjusted peers or benchmark experience – often called the “expected” result. This expected result is adjusted to reflect both the peers/benchmark levels of performance and also the provider's own case mix of episodes by condition and level of severity. The ratio of observed to expected results can be termed the relative cost ratio and is a risk adjusted measure.
Stratification	The severity level can then be used to stratify episodes by severity, measured as resource consumption.

1415

1416

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

	1611: ETG Based Pneumonia cost of care measure
Steward	Ingenix
Description	The measure focuses on resources used to deliver episodes of care for patients with pneumonia. Pneumonia episodes are defined using the Episode Treatment Groups (ETG) methodology and describe the unique presence of the condition for a patient and the services involved in diagnosing, managing and treating pneumonia. A number of resource use measures are defined for pneumonia episodes, including overall cost of care, cost of care by type of service, and the utilization of specific types of services. Each resource use measure is expressed as a cost or a utilization count per episode and comparisons with internal and external benchmarks are made using risk adjustment to support valid comparisons.
Resource Use Measure Type	Per episode
Data Source	Administrative claims, Other: Both medical and pharmacy administrative service records (claims or encounters) are used to support the measures. Member enrollment span, pharmacy benefit status and age and gender are also required. Provider characteristics, including specialty and unique provider identifier also have importance to support episode grouping, attribution and definition of peers.
Level of Analysis	Clinician : Group/Practice, Individual, Team Facility Health Plan Integrated Delivery System Population: Community, County or City, National, Regional, State
Clinical Framework Description	The pneumonia measure's episodes are defined using the Episode Treatment Group (ETG) methodology. The pneumonia ETG episode building process that supports pneumonia resource use measures has four important steps: Step 1: Identify Records; Assign Record Type and Anchor Records, Classify Diagnoses and Procedures Step 2: Build Episodes from Anchor Records Step 3: Group Non-Anchor Records to Episodes Step 4: Finalize the Episodes (identify co-morbidities and complicating factors, and assign episode severity)
Costing Method	The financial amounts used should be complete and valid, reflecting the total payments related to the service. The financial amount used in resource measurement should reflect all payments for a service, including those made to the provider by payer, patient and other entities. Allowed payments will reflect both the quantity of different services provided as well as the actual unit price of those same services.
Tested Population	Commercial
Resource Use Service Categories	Inpatient services: Inpatient facility services; Inpatient services: Admissions/discharges; Ambulatory services: Outpatient facility services; Ambulatory services: Emergency Department; Ambulatory services: Pharmacy; Ambulatory services: Evaluation and management; Ambulatory services: Procedures and surgeries; Ambulatory services: Imaging and diagnostic; Ambulatory services: Lab services
Attribution Approach	Guidelines: Both activity-based and population-based approaches should be supported. As a guideline, four different general options for physician episode attribution can be considered to attribute episodes to individual providers – three activity-based and one population-based approach. Approach 1 - Physician Episode Attribution using Professional Service Costs. This attribution approach identifies the responsible physician for an episode as that provider rendering the greatest amount of professional service costs during the episode. Approach 2 - Physician Episode Attribution using Episode Clusters. This attribution approach identifies the responsible physician for an episode as that provider in the peer group owning the greatest number of "clusters" within the episode. Approach 3 - Physician Episode Attribution using Non-Acute Evaluation and Management (E/M) Visits. This attribution approach identifies the responsible physician for an episode as that physician providing the greatest number of non-acute E/M visits within the episode.

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

	1611: ETG Based Pneumonia cost of care measure
	Approach 4 - Physician Episode Attribution using a Primary Care, Population-based Approach. This approach requires two important steps: 1) Identification of a PCP for each member. 2) Identify the patient's assigned PCP during the episode period.
Risk Adjustment	<p>ETG first assesses the observed co-morbidities and condition status factors for an episode and the patient's age and gender. ETG then assigns a weight to each factor found to influence the relative risk of an episode of pneumonia. These weights and factors are condition-specific and were estimated using pneumonia episode results for a large population. The overall severity score for an episode is the sum of these weights for all factors observed. Using the severity score, a severity level is created, with each pneumonia episode assigned to one of four severity levels. The level of severity assigned to an episode is used to support risk adjustment. The risk adjustment approach includes three important steps:</p> <ol style="list-style-type: none"> 1. Compute the observed experience for the provider being measured, across all episodes to be included in the comparison; 2. Compute the experience for peers or a best practice benchmark. Compute this experience at the level of the risk adjustment, in this case base procedure (hip or knee replacement) and severity level. For a peers benchmark, average cost per episode across all peers for the base procedure and severity level can be computed; 3. Compare the observed experience with the risk adjusted peers or benchmark experience – often called the “expected” result. This expected result is adjusted to reflect both the peers/benchmark levels of performance and also the provider’s own case mix of episodes by condition and level of severity. The ratio of observed to expected results can be termed the relative cost ratio and is a risk adjusted measure.
Stratification	ETG stratifies episodes by the intensity of service, or total cost. For a given episode, a severity score is assigned based on demographic factors (gender and age) and the presence of comorbidities and complications. Once a severity score is determined, a severity level, a number between 1 and 4 is assigned based on a table that relates severity levels to severity scores for each ETG. The severity level can then be used to stratify episodes by severity, measured as resource consumption.

1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429

NATIONAL QUALITY FORUM

1430 APPENDIX B—STEERING COMMITTEE

1431

1432 **Tom Rosenthal, MD (Co-Chair)**

1433 UCLA School of Medicine, Los Angeles, CA

1434

1435 **Bruce Steinwald, MBA (Co-Chair)**

1436 Independent Consultant, Washington, DC

1437

1438 **Paul G. Barnett, PhD**

1439 VA Palo Alto Health Care System, Menlo Park, CA

1440

1441 **Jack Bowhan**

1442 Wisconsin Collaborative for Healthcare Quality, Middleton, WI

1443

1444 **Jeptha P. Curtis, MD**

1445 Yale University School of Medicine, New Haven, CT

1446

1447 **Kurtis S. Elward, MD, MPH**

1448 Family Medicine of Albemarle, Charlottesville, VA

1449

1450 **William E. Golden, MD**

1451 Arkansas Medicaid, Little Rock, AR

1452

1453 **Lisa M. Grabert, MPH**

1454 American Hospital Association, Washington, DC

1455

1456 **Ethan A. Halm, MD, MPH**

1457 University of Texas Southwestern Medical Center, Dallas, TX

1458

1459 **Ann L. Hendrich, RN, MSN, PhD(c)**

1460 Ascension Health, St. Louis, MO

1461

1462 **Thomas H. Lee, MD**

1463 Partners HealthCare System, Inc., Boston, MA

1464

1465 **Jack Needleman, PhD**

1466 University of California, Los Angeles School of Public Health

1467

1468 **Mary Kay O'Neill, MD, MBA**

1469 CIGNA HealthCare, Seattle, WA

1470

1471 **David F. Penson, MD, MPH**

1472 Vanderbilt University Medical Center, Nashville, TN

1473

1474 **Doris Peter, PhD**

1475 Consumer Reports, Yonkers, NY

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

- 1476
1477 **Steve Phillips, MPA**
1478 Johnson & Johnson Health Care Systems Inc., Washington, DC
1479
1480 **David Redfearn, PhD**
1481 WellPoint, Las Vegas, NV Woodland Hills, CA
1482
1483 **Jeffrey B. Rich, MD**
1484 Mid-Atlantic Cardiothoracic Surgeons Ltd., Norfolk, VA
1485
1486 **William L. Rich, III, MD**
1487 Northern Virginia Ophthalmology Associates, Falls Church, VA
1488
1489 **Barbara A. Rudolph, PhD, MSSW**
1490 The Leapfrog Group, Fitchburg, WI
1491
1492 **Joseph Stephansky, PhD**
1493 Michigan Health & Hospital Association, Lansing, MI
1494
1495 **James N. Weinstein, DO, MS**
1496 The Dartmouth Institute for Health Policy and Clinical Practice & The Dartmouth-Hitchcock Clinic,
1497 Lebanon, NH
1498
1499 **Dolores Yanagihara, MPH**
1500 Integrated Healthcare Association, Oakland, CA
1501
1502
1503 **NQF Staff**
1504
1505 **Helen Burstin, MD, MPH**
1506 Senior Vice President, of Performance Measures
1507
1508 **Heidi Bossley, MBA, MSN**
1509 Vice President, of Performance Measures
1510
1511 **Taroon Amin, MA, MPH**
1512 Senior Director
1513
1514 **Ashlie Wilbon, RN, MPH**
1515 Senior Project Manager
1516
1517 **Lauralei Dorian**
1518 Project Manager
1519
1520 **Sarah Fanta**
1521 Project Analyst

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

1522 APPENDIX C—TECHNICAL ADVISORY PANELS

1523

1524 Cardiovascular/Diabetes Technical Advisory Panel

1525 Jeptha Curtis, MD, FACC (Co-Chair)

1526 Yale University School of Medicine, New Haven, CT

1527

1528 James Rosenzweig, MD (Co-Chair)

1529 Boston Medical Center and Boston University School of Medicine, Boston, MA

1530

1531 Mary Ann Clark, MHA

1532 Neocure Group, Washington, DC

1533

1534 Constance Hwang, MD, MPH

1535 Resolution Health, Inc., Columbia, MD

1536

1537 Thomas Marwick, MBBS, PhD

1538 Cleveland Clinic, Cleveland, OH

1539

1540 Michael O'Toole, MD

1541 Midwest Heart Specialists, Ltd., Downers Grove, IL

1542

1543 David Palestrant, MD

1544 Cedars-Sinai Medical Center, Los Angeles, CA

1545

1546 Brenda Parker, PharmD

1547 GlaxoSmithKline, Marietta, GA

1548

1549 Katherine Reeder, PhD, RN

1550 University of Kansas School of Nursing, Kansas City, KS

1551

1552 William Weintraub, MD

1553 Christiana Care Health System, Newark, DE

1554

1555

1556

1557

1558

1559

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

1560 **Pulmonary Technical Advisory Panel**

1561 **Kurtis Elward, MD, MPH (Co-Chair)**

1562 Family Medicine of Albermarle, Charlottesville, VA

1563

1564 **Janet Maurer, MD, MBA (Co-Chair)**

1565 American College of Chest Physicians, Northbrook, IL

1566

1567 **Gerene Bauldoff, PhD, RN**

1568 The Ohio State University, School of Nursing, Columbus, OH

1569

1570 **Kathryn Blake, PharmD**

1571 Nemours Children's Clinic, Jacksonville, FL

1572

1573 **Dale Bratzler, DO, MPH**

1574 University of Oklahoma, Health Sciences Center, Oklahoma City, OK

1575 **Zab Mosenifar, MD**

1576 Cedars Sinai Medical Center, Los Angeles, CA

1577

1578 **Linus Santo Tomas, MD, MS**

1579 Pulmonary & Critical Care, Medical College of Wisconsin, Milwaukee, WI

1580

1581 **Michael Schatz, MD, MS**

1582 Kaiser Permanente, Oakland, CA

1583

1584 **Richard Stanford, PharmD, MS**

1585 GlaxoSmithKline, Research Triangle Park, NC

1586

1587

1588

1589

1590

1591

1592

1593

1594

1595

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

1596 **Bone/Joint Technical Advisory Panel**

1597 **James Weinstein, DO, MS(Chair)**

1598 The Dartmouth Institute for Health Policy; Dartmouth-Hitch Clinic, Lebanon, NH

1599

1600 **Mary Kay O'Neill, MD, MBA**

1601 CIGNA HealthCare, Seattle, WA

1602

1603 **Elizabeth Paxton, MA**

1604 Kaiser Permanente, Oakland, CA

1605

1606 **John Ratliff, MD, FACS**

1607 Thomas Jefferson University, Philadelphia, PA

1608

1609 **Catherine Roberts, MD**

1610 Mayo Clinic, Phoenix, AZ

1611

1612 **Craig Rubin, MD**

1613 University of Texas Southwestern Medical School, Dallas, TX

1614

1615 **Patricia Sinnott, PT, PhD, MPH**

1616 VA health Economics Resource Center, Menlo Park, CA

1617

1618

1619

1620

1621

1622

1623

1624

1625

1626

1627

1628

1629

1630

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

- 1631 **Cancer Technical Advisory Panel**
- 1632 **David Penson, MD, MPH (Chair)**
1633 Vanderbilt University Medical Center, Nashville, TN
1634
- 1635 **Rohit Borker, PhD**
1636 GlaxoSmithKline, Philadelphia, PA
1637
- 1638 **Steven Chen, MD, MBA**
1639 California Medical Association, Camerillo, CA
1640
- 1641 **Timothy Gilligan, MD**
1642 Cleveland Clinic Taussig Cancer Institute, Cleveland, OH
1643
- 1644 **Stephen Grossbart, PHD**
1645 Catholic Healthcare Partners, Cincinnati, OH
1646
- 1647 **Dwight Kloth, PharmD**
1648 Fox Chase Cancer Center, Philadelphia, PA
1649
- 1650 **Louis Potters, MD, FACR**
1651 North Shore-Long Island Jewish Health System, New Hyde Park, NY
1652
- 1653 **Jay Schukman, MD**
1654 Anthem Blue Cross and Blue Shield, Richmond, VA
1655
- 1656 **John Skibber, MD**
1657 University of Texas-MD Anderson Cancer Center, Houston, TX
1658
- 1659 **Louise Walter, MD**
1660 University of California - San Francisco, San Francisco, CA
1661
- 1662
- 1663
- 1664
- 1665
- 1666
- 1667
- 1668

NQF REVIEW DRAFT—DO NOT CITE OR QUOTE

NQF MEMBER comments due November 21, 2011, 6:00 PM ET; PUBLIC comments due November 14, 2011 by 6:00 PM ET

NATIONAL QUALITY FORUM

1669 APPENDIX D—RESOURCE USE MEASUREMENT TERMS

1670 The following resource use measurement terms have been defined based on their use in the
1671 context of this project and are important to understanding the concepts in this report.

1672 **Attribution**—identifying and assigning of a responsible provider or entity (e.g., health plan) for
1673 the care delivered for an episode or population.

1674
1675 **Benchmarking**—the process of comparing the performance of accountable entities with that of
1676 their peers or with external best practice results. In developing comparative estimates, results
1677 should be risk adjusted for patient-level attributes to support the valid comparisons of these
1678 accountable entities.

1679
1680 **Carve-outs**—the outsourcing of services, such as behavioral health or pharmacy claims, to
1681 specialty health plans or claims processing entities or organizations.

1682
1683 **Clinical hierarchy**—an arrangement of clinical conditions that are ranked according to severity,
1684 as “high,” “below,” or “at the same level.” For example, if a patient has COPD and develops
1685 bronchitis, COPD would be assigned a greater weight than bronchitis.

1686

1687 **Exclusion criteria**—criteria applied before a measure is tested in order to remove any
1688 individuals with conditions that may skew the final measure score.

1689

1690 **Peer groups**—the ways in which resource use measures ensure providers and health plans are
1691 compared to similar providers and health plans.

1692

1693 **Per capita measure**—counts all services provided to a person within a specific population,
1694 regardless of condition or encounters with system.

1695

1696 **Per episode measure**—counts resources based on bundles of services that are part of a
1697 distinctive event provided by one or multiple entities (e.g., health services provided associated
1698 with an event or series of events for acute myocardial infarction).

1699

1700 **Resource use service categories**—categories of resource units or services provided care for a
1701 patient or population. Resource units are generally identified through claims data and
1702 grouped into categories with similar types of claims (e.g., x-rays grouped into imaging category).
1703 Categories are generally measured in terms of dollars, but also can also include resources
1704 not captured on a claim (e.g., nursing hours).

1705

1706 **Risk adjustment**—a corrective approach designed to reduce any negative or positive
1707 consequences associated with caring for patients of higher or lower health risk or propensity to
1708 require health services.

1709

NATIONAL QUALITY FORUM

1710 **Severity levels**—pre-determined levels of acuity used to rank and assign patients based on an
1711 assessment of the aggregate of their conditions/diagnosis codes.

1712
1713 **Standardized pricing**—pre-established uniform price for a service, typically based on historical
1714 price, replacement cost, or an analysis of completion in the market; removes variation in resource
1715 costs due to differences in negotiated prices or geographic differences based on labor or other
1716 input costs.

1717
1718 **Stratification**—division of a population or resource services into distinct, independent strata, or
1719 groups of similar data, enabling analysis of the specific subgroups. This type of adjustment can
1720 be used to show where disparities exist or where there is a need to expose differences in results.

1721