

1

2 DR. JAIME JORDAN (Orcid ID : 0000-0002-6573-7041)

3 DR. LAURA R HOPSON (Orcid ID : 0000-0002-1183-4751)

4 DR. SUZANNE BENTLEY (Orcid ID : 0000-0003-0192-3133)

5 DR. NICOLE M. DEIORIO (Orcid ID : 0000-0002-8123-1112)

6 DR. SALLY SANTEN (Orcid ID : 0000-0002-8327-8002)

7 DR. LALENA M YARRIS (Orcid ID : 0000-0003-1277-2852)

8 DR. WENDY C. COATES (Orcid ID : 0000-0002-3305-8802)

9 DR. MICHAEL A. GISONDI (Orcid ID : 0000-0002-6800-3932)

10

11

12 Article type : Original Contribution

13

14

15 **Title:** Leveling the field: Development of Reliable Scoring Rubrics for Quantitative and
16 Qualitative Medical Education Research Abstracts

17

18 **Running Title:** Abstract Scoring Rubrics...

19 **Authors:** Jaime Jordan, MD, MAEd,^{1,2} Laura R. Hopson, MD,³ Caroline Molins, MD,
20 MSMEd,⁴ Suzanne K. Bentley, MD, MPH,⁵ Nicole M. Deiorio, MD,⁶ Sally A. Santen, MD,
21 PhD,^{6,7} Lalena M. Yarris, MD, MCR,⁸ Wendy C. Coates, MD,¹ Michael A. Gisondi, MD⁹

22

23 **Affiliations:**

This is the author manuscript accepted for publication and has undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the [Version of Record](#). Please cite this article as [doi: 10.1002/AET2.10654](https://doi.org/10.1002/AET2.10654)

This article is protected by copyright. All rights reserved

24 ¹ Department of Emergency Medicine, David Geffen School of Medicine at UCLA, Los
25 Angeles, CA

26 ² Department of Emergency Medicine, Ronald Reagan UCLA Medical Center, Los Angeles, CA

27 ³Department of Emergency Medicine, University of Michigan, Ann Arbor, MI

28 ⁴AdventHealth Emergency Medicine Residency, Orlando, FL

29 ⁵Icahn School of Medicine at Mount Sinai, New York, NY

30 ⁶Virginia Commonwealth University School of Medicine, Richmond, VA

31 ⁷University of Cincinnati College of Medicine

32 ⁸Department of Emergency Medicine, Oregon Health & Science University, Portland, OR

33 ⁹Department of Emergency Medicine, Stanford University, Palo Alto, CA

34

35

36 **Please address correspondence to:**

37 Jaime Jordan

38 UCLA Emergency Medicine

39 924 Westwood Boulevard, Suite 300

40 Los Angeles, CA 90024

41 Tel: 310-794-0585

42 Fax: 310-794-0599

43 Email: jaimejordanmd@gmail.com

44

45 **Presentations:** This work was presented at SAEM Annual Meeting, May 13th, 2021.

46 **Financial support:** none

47 **Acknowledgements:** We would like to acknowledge that this project originated to meet an
48 SAEM Education Committee Objective and would like to thank all the committee members for
49 their support of this work.

50 **Conflicts of interest:** JJ, LRH, CM, SKB, NMD, SAS, LMY, WCC, MAG report no conflict of
51 interest.

52 **Author contributions:** JJ and MAG conceived the study. JJ, MAG, LRH, CM, and SKB
53 contributed to the design of the study. JJ, LRH, CM, SKB, NMD, SAS, LMY, WCC, MAG

54 contributed to data collection. JJ analyzed the data. JJ, LRH, CM, SKB, NMD, SAS, LMY,
55 WCC, MAG contributed to drafting of the manuscript and critical revision.

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75 **Abstract:**

76 **Background:** Research abstracts are submitted for presentation at scientific conferences,
77 however, criteria for judging abstracts are variable. We sought to develop two rigorous abstract
78 scoring rubrics for education research submissions reporting (1) quantitative data and (2)
79 qualitative data; and then to collect validity evidence to support score interpretation.

80

81 **Methods:** We used a modified Delphi method to achieve expert consensus for scoring rubric
82 items to optimize content validity. Eight education research experts participated in two separate
83 modified Delphi processes, one to generate quantitative research items and one for
84 qualitative. Modifications were made between rounds based on item scores and expert feedback.

85 Homogeneity of ratings in the Delphi process was calculated using Cronbach's alpha, with
86 increasing homogeneity considered an indication of consensus. Rubrics were piloted by scoring
87 abstracts from 22 quantitative publications from *Academic Emergency Medicine Education and*
88 *Training* "Critical Appraisal of Emergency Medicine Education Research" (11 highlighted for
89 excellent methodology and 11 that were not) and 10 qualitative publications (5 highlighted for
90 excellent methodology and 5 that were not). Intraclass correlation coefficient (ICC) estimates of
91 reliability were calculated.

92

93 **Results:** Each rubric required three rounds of a modified Delphi process. The resulting
94 quantitative rubric contained nine items: quality of objectives, appropriateness of methods,
95 outcomes, data analysis, generalizability, importance to medical education, innovation, quality of
96 writing, and strength of conclusions. Cronbach's alpha for the 3rd round=0.922; ICC for total
97 scores during piloting = 0.893. The resulting qualitative rubric contained 7 items: quality of
98 study aims, general methods, data collection, sampling, data analysis, writing quality, and
99 strength of conclusions. Cronbach's alpha for the 3rd round = 0.913; ICC for the total scores
100 during piloting =0.788.

101

102 **Conclusion:** We developed scoring rubrics to assess quality in quantitative and qualitative
103 medical education research abstracts to aid in selection for presentation at scientific meetings.
104 Our tools demonstrated high reliability.

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121 Introduction

122 The scientific abstract is the standard method for researchers to communicate brief
123 written summaries of their findings. The written abstract is the gatekeeper for selection for
124 presentation at professional society meetings.¹ A research presentation serves many purposes
125 including dissemination of new knowledge, an opportunity for feedback, and the prospect of
126 fostering an investigator's academic reputation. Beyond the presentation, abstracts, as written
127 evidence of scientific conference proceedings, often endure through publication in peer-reviewed
128 journals. Because of the above, abstracts may be assessed in a number of potentially high-stakes
129 situations.

130 Abstracts are selected for presentation at conferences through a competitive process
131 based on factors such as study rigor, importance of research findings, and relevance to the
132 sponsoring professional society. Prior literature has shown poor observer-agreement in the
133 abstract selection process.² Scoring rubrics are often used to guide abstract reviewers in an
134 attempt to standardize the process, reduce bias, support equity, and promote quality.³ There are
135 limited data describing the development and validity evidence of such scoring rubrics but the
136 data available suggest that rubrics may be based on quality scoring tools for full research reports
137 and published guidelines for abstracts.^{2, 4-5} Medical conferences often apply rubrics designed for
138 judging clinical or basic science submissions which reflect standard hypothesis-testing methods
139 and often use a single subjective Gestalt rating for quality decisions.⁶ This may result in the
140 systematic exclusion of studies which employ alternate, but equally rigorous methods, such as
141 research in medical education. Existing scoring systems, commonly designed for biomedical
142 research, may not accurately assess the scope, methods, and types of results commonly reported
143 in medical education research abstracts, which may lead to a disproportionately high rate of
144 rejection of these abstracts. There are additional challenges in reviewing qualitative research
145 abstracts using a standard hypothesis-testing rubric. In these qualitative studies, word-count
146 constraints may limit the author's ability to convey the study's outcome appropriately.⁷ It is

147 problematic for qualitative studies to be constrained to a standard quantitative abstract template,
148 which may lead to low scores by those applying the rubric and a potential systematic bias against
149 qualitative research.

150 Prior literature has described tools to assess quality in medical education research
151 manuscripts, such as the Medical Education Research Study Quality Instrument (MERSQI) and
152 the Newcastle-Ottawa Scale-Education (NOS-E).⁸ A limited attempt to utilize the MERSQI tool
153 to retrospectively assess internal medicine medical education abstracts achieving manuscript
154 publication showed increased scores for the journal abstract relative to the conference abstract.⁴
155 However, the MERSQI and similar tools were not developed specifically for judging abstracts,
156 and there is a lack of published validity evidence to support score interpretation based on these
157 tools. In order to equitably assess the quality of education research abstracts to scholarly venues,
158 which may have downstream effects on researcher scholarship, advancement, and reputation,
159 there is a need for a rigorously developed abstract scoring rubric that is based on a validity
160 evidence framework.⁹⁻¹⁰

161 The aim of this paper is to describe the development and pilot testing of a dedicated
162 rubric to assess the quality of both quantitative and qualitative medical education research
163 studies. We describe the development process, which aimed to optimize content and response
164 process validity, and initial internal structure and relation to other variables validity evidence to
165 support score interpretation using these instruments. The rubrics may be of use to researchers
166 developing studies, abstract and paper reviewers, and may be applied to medical education
167 research assessment in other specialties.

168

169 **Methods**

170

171 **Study Design**

172 We utilized a modified Delphi technique to achieve consensus on items for a scoring
173 rubric to assess quality of emergency medicine (EM) education research abstracts. The modified
174 Delphi technique is a systematic group consensus strategy designed to increase content
175 validity.¹¹ Through this method we developed individual rubrics to assess quantitative and
176 qualitative EM medical education research abstracts. This study was approved by the
177 Institutional Review Board of the David Geffen School of Medicine at UCLA.

178

179 Study setting and participants

180 The first author identified eight EM education researchers with successful publication
181 records from diverse regions across the United States and invited them to participate in the
182 Delphi panel. Previous work has suggested that 6-10 experts is an appropriate number for
183 obtaining stable results in the modified Delphi method.¹²⁻¹⁴ All invited panelists agreed to
184 participate. The panel included one assistant professor, two associate professors, and 5
185 professors. All panelists serve as reviewers for medical education journals and four hold
186 editorial positions. We collected data in September and October, 2020.

187

188 Study protocol

189 We followed Messick's framework for validity that includes five types of validity
190 evidence; content, response process, internal structure, relation to other variables, and
191 consequential.¹⁵ Our study team drafted initial items for the scoring rubrics after a review of the
192 literature and existing research abstract scoring rubrics to optimize content validity. We created
193 separate items for research abstracts reporting quantitative and qualitative data. We sent the
194 draft items to the Society of Academic Emergency Medicine (SAEM) education committee for
195 review and comment in order to gather stakeholder feedback and for further content and response
196 process validity evidence.¹⁶ One author (JJ) who was not a member of the Delphi panel then
197 revised the initial lists of items based on committee feedback to create the initial Delphi
198 surveys. We used an electronic survey platform (SurveyMonkey) to administer and collect data
199 from the Delphi surveys.¹⁷ Experts on the Delphi panel rated the importance of including each
200 item in a scoring rubric on a 1-9 Likert scale with 1 labeled as "not at all important" and 9
201 labeled as "extremely important". The experts were invited to provide additional written
202 comments, edits and suggestions for each item. They were also encouraged to suggest additional
203 items that they felt were important but not currently listed. We determined *a priori* that items
204 with a mean score of 7 or greater advanced to the next round and items with a mean score of 3 or
205 below were eliminated. The Delphi panel moderator (JJ) applied discretion for items scoring
206 between 4 and 6, with the aim of both adhering to the opinions of the experts and creating a
207 comprehensive scoring rubric. For example, if an item received a middle score but had

208 comments supporting inclusion in a revised form, the moderator would make the suggested
209 revisions and include the item in the next round.

210 Each item consisted of a stem and anchored choices with associated point-value
211 assignments. Panelists commented on the stems, content and assigned point-value of choices and
212 provided narrative unstructured feedback. The moderator made modifications between rounds
213 based on item scores and expert feedback. After each round, we provided panelists with
214 aggregate mean item scores, written comments, and an edited version of the item list derived
215 from the responses in the previous round. The panelists were then asked to rate the revised items
216 and provide additional edits or suggestions.

217 We considered homogeneity of ratings in the Delphi process to be an indication of
218 consensus. After consensus was achieved, we created final scoring rubrics for quantitative and
219 qualitative medical education research abstracts. We then piloted the scoring rubrics to gather
220 internal structure and further response process validity evidence. Five raters from the study
221 group (JJ, LH, MG, CM, SB) participated in piloting. We piloted the final quantitative research
222 rubric by scoring abstracts from publications identified in the most recent critical appraisal of
223 EM education research by Academic Emergency Medicine/Academic Emergency Medicine
224 Education and Training, “Critical Appraisal of Emergency Medicine Education Research: The
225 Best Publications of 2016”.¹⁸ All 11 papers highlighted for excellent methodology in this issue
226 were included in the pilot.¹⁸ Additionally, we included an equal number of randomly selected
227 citations that were included in the issue but not selected as top papers, for a total of 22
228 quantitative publications.¹⁸ Given the limited number of qualitative studies cited in this issue of
229 the critical appraisal series, we chose to pilot the qualitative rubric on publications from this
230 series from the last 5 years available (2012-2016).¹⁸⁻²² We randomly selected one qualitative
231 publication that was highlighted for excellent methodology and one that was not from each year
232 for a total of 10 qualitative publications.¹⁸⁻²² The same five raters who performed the quantitative
233 pilot also conducted the qualitative pilot.

234

235 **Statistical Analysis**

236 We calculated and reported descriptive statistics for item scoring during Delphi
237 rounds. We used Cronbach’s alpha to assess homogeneity of ratings in the Delphi process.
238 Increasing homogeneity was considered to be an indication of consensus among the expert

239 panelists. We used intraclass correlation coefficient (ICC) estimates to assess reliability among
240 raters during piloting based on a mean rating ($k=5$), absolute agreement, 2-way random-effects
241 model. We performed all analyses in SPSS (IBM SPSS Statistics for Windows, Version 27.0.
242 Armonk, NY: IBM Corp).

243

244 **Results**

245

246 **Quantitative rubric:**

247 Three Delphi rounds were completed, each with 100% response rate. Mean item scores
248 for each round are depicted in Table 1. After the first round, three items were deleted, one item
249 was added, and five items underwent wording changes. After the second round, one item was
250 deleted and eight items underwent wording changes. After the third round items were re-ordered
251 for flow and ease of use but no further changes were made to content or wording. Cronbach's
252 alpha for the third round was 0.922 indicating high internal consistency. The final rubric
253 contained nine items: quality of objectives, appropriateness of methods, outcomes, data analysis,
254 generalizability, importance to medical education, innovation, quality of writing, and strength of
255 conclusions (Appendix A). The ICC for the total scores during piloting was 0.893, indicating
256 excellent agreement. ICCs for individual rubric items ranged from 0.406 to 0.878 (Table 3).

257

258 **Qualitative rubric:**

259 Three Delphi rounds were completed, each with 100% response rate. Mean item scores
260 for each round are depicted in Table 2. After the first round two items were deleted, one item
261 was added and nine items underwent wording changes. After the second round, three items were
262 deleted and four underwent wording changes. After the third round no further changes were
263 made. The resulting tool contained 7 items reflecting the domains of quality of study aims,
264 general methods, data collection, sampling, data analysis, writing quality, and strength of
265 conclusions (Appendix B). Cronbach's alpha for the third round was 0.913, indicating high
266 internal consistency. ICC for the total scores during piloting was 0.788 indicating good
267 agreement. The item on writing quality had an ICC of -0.301, likely due to the small scale of the
268 item and sample size leading to limited variance. ICCs for the remainder of the items ranged
269 from 0.176 to 0.897 (Table 3).

270

271 **Discussion:**

272 We developed novel and distinct abstract scoring rubrics for assessing quantitative and
273 qualitative medical education abstract quality through a Delphi process. It is important to
274 evaluate medical education research abstracts that utilize accepted education methods as a
275 distinctly different class than basic, clinical, and translational research. Through our Delphi and
276 piloting processes we have provided multiple types of validity evidence in support of these
277 rubrics aligned with Messick's framework including content, response process and internal
278 structure.¹⁵ Similar to other tools assessing quality in medical education research, our rubrics
279 assess aspects such as study design, sampling, data analysis, and outcomes that represent the
280 underpinnings of rigorous research.^{8, 23-26} Unlike many medical education research assessments
281 published in the literature, our tool was designed specifically for the assessment of abstracts
282 rather than full text manuscripts and therefore the specific item domains and characteristics
283 reflect this unique purpose.

284 We deliberately created separate rubrics for abstracts reporting quantitative and
285 qualitative data as each has unique methods. When designing a study, education researchers must
286 decide the best method to address their questions. Often, in the exploratory phase of inquiry, a
287 qualitative study is the most appropriate choice to identify key topics that merit further study.
288 These often may be narrow in scope and may employ one or more qualitative methods (e.g.,
289 ethnography, focus groups, personal interviews). The careful and rigorous analysis may reveal
290 points that can be studied later via quantitative methods to test a hypothesis gleaned during the
291 qualitative phase.²⁷ Specific standards for reporting on qualitative research have been widely
292 disseminated and are distinct from standards for reporting quantitative research.²⁸ Even an
293 impeccably designed and executed qualitative study would fail to meet major criteria for
294 excellent quantitative studies. For example, points may be subtracted for lack of generalizability
295 or conduct of the qualitative study in multiple institutions, as well as for the absence of common
296 quantitative statistical analytics. The qualitative abstract itself may necessarily lack the common
297 structure of a quantitative report and lead to a lower score. The obvious problem is that a well-
298 conducted study might not be shared with the relevant research community if it is judged
299 according to quantitative standards. A similar outcome would occur if quantitative work were

300 judged by qualitative standards, therefore we advocate for using scoring rubrics specific to the
301 type of research being assessed.

302 Our work has several possible applications. The rubrics we developed may be adopted as
303 scoring tools for medical education research studies that are submitted for presentation to
304 scientific conferences. The presence of specific scoring rubrics for medical education research
305 may address disparities in acceptance rates and ensure presentation of rigorously conducted
306 medical education research at scientific conferences. Further, publication of abstract scoring
307 rubrics such as ours sets expectations for certain elements to be included and defines an
308 acceptable level of submission quality. Dissemination and usage of the rubrics may therefore
309 help improve research excellence. The rubrics themselves can serve as educational tools in
310 resident and faculty training. For example, the rubrics could serve as illustrations or practice
311 material in teaching how to prepare a strong abstract for submission. The inclusive wording of
312 the items allows the rubrics to be adapted to medical education work in any medical specialty.
313 Medical educators may also benefit from using the methods described here to create their own
314 scoring rubrics or provide evidence-based best practice approaches for other venues. Finally, this
315 study provides a tool that could lay the groundwork for future scholarship on assessing the
316 quality of educational research.

317

318 **Limitations**

319 Our study has several limitations. First, the modified Delphi technique is a consensus
320 technique which can force agreement of respondents and the existence of consensus does not
321 denote a correct response.¹¹ Since the method is implemented electronically, there is limited
322 discussion and elaboration. Second, the team of experts were all researchers in EM, therefore the
323 rubrics may not generalize to other specialties. The rubrics were intended for quantitative and
324 qualitative education research abstract submission, so it may not perform well for abstracts that
325 include *both* quantitative and qualitative data or those focused on early work, innovations,
326 instrument development, validity evidence, or program evaluation. Finally, there are two
327 limitations to the pilot testing. An *a priori* power calculation to determine sample size was not
328 possible since the rubrics were novel. The ICCs of individual items on the scoring rubrics were
329 variable and we chose not to eliminate items with low ICCs given the small sample size during
330 piloting and a desire to create a tool comprehensive of key domains. Future studies of use of

331 these tools incorporating larger samples may provide data for additional refinement. Faculty who
332 piloted the rubrics were familiar with the constructs and rubrics, and it is not known how the
333 rubrics would have performed with general abstract reviewers nor what training might be
334 required. The success of separate rubrics may rely on the expertise of the reviewers in the
335 methodology being assessed.

336 We offer two medical education abstract scoring rubrics with supporting preliminary
337 reliability and validity evidence. Future studies could add additional validity evidence including
338 use with trained and untrained reviewers and relationship to other variables, e.g. a comparison
339 between rubric scores and expert judgement. Additional studies could be done to provide
340 consequential validity evidence by comparing the number and quality of accepted medical
341 education abstracts before and after the rubric's implementation or whether the number of
342 abstracts that eventually lead to publication increases.

343

344 **Conclusions**

345 Using the modified Delphi technique for consensus building, we developed two scoring
346 rubrics to assess quality in quantitative and qualitative medical education research abstracts with
347 supporting validity evidence. Application of these rubrics demonstrated high reliability.

348

349

350 **References**

- 351 1. Padayachy A, Rodrigues G, Tahar A. Comment rédiger un abstract scientifique ? [How to
352 write a scientific abstract]. Rev Med Suisse 2019;15(664):1703-1706.
- 353 2. Timmer A, Sutherland LR, Hilsden RJ. Development and evaluation of a quality score for
354 abstracts. BMC Med Res Methodol 2003;3:2.
- 355 3. Mitchell NS, Stolzmann K, Benning LV, Wormwood JB, Linsky AM. Effect of a Scoring
356 Rubric on the Review of Scientific Meeting Abstracts. J Gen Intern Med 2020. Epub ahead of
357 print.
- 358 4. Poolman RW, Keijser LC, de Waal Malefijt MC, et al. Reviewer agreement in scoring
359 abstracts for scientific orthopedics meetings. Acta orthop 2007;78(2):278–284.
- 360 5. van der Steen LP, Hage JJ, Kon M, et al. Reliability of a structured method of selecting
361 abstracts for a plastic surgical scientific meeting. Plast Reconstr Surg 2003;111(7):2215–2222.

- 362 6. Kuczmariski TM, Raja AS, Pallin DJ. How do Medical Societies Select Science for
363 Conference Presentation? How Should They? West J Emerg Med 2015;16(4):543-50.
- 364 7. Stephenson CR, Vaa BE, Wang AT, et al. Conference presentation to publication: a
365 retrospective study evaluating quality of abstracts and journal articles in medical education
366 research. BMC Med Educ 2017;17(1):193.
- 367 8. Cook DA, Reed DA. Appraising the quality of medical education research methods: the
368 Medical Education Research Study Quality Instrument and the Newcastle-Ottawa Scale-
369 Education. Acad Med 2015;90(8):1067-76.
- 370 9. Sullivan GM. A primer on the validity of assessment instruments. J Grad Med Educ 2011;3:
371 119-120.
- 372 10. Cook DA, Beckman TJ. Current Concepts in Validity and Reliability for Psychometric
373 Instruments: Theory and Application. Am J of Medicine 2006;119:166e10-199.e16.
- 374 11. Hasson F, Keeney S, Mckenna HP. Research guidelines for the Delphi Survey Technique
375 Research guidelines for the Delphi survey. J Adv Nurs 2000;32(4):1008-15.
- 376 12. Thangaratinam SK, Redman CWE. The Delphi Technique. The Obstetrician and
377 Gynaecologist 2005;7:120-125.13. Pines JM, Alfaraj S, Batra S, et al. Factors important to top
378 clinical performance in emergency medicine residency: results of an ideation survey and Delphi
379 panel. AEM Educ Train 2018;2(4):269-276.
- 380 14. Eubank BH, Mohtadi NG, Lafave MR, et al. Using the modified Delphi method to establish
381 clinical consensus for the diagnosis and treatment of patients with rotator cuff pathology. BMC
382 Medical Research Methodology 2016;16:56.
- 383 15. Messick S. Validity. In Linn R, ed. Educational Measurement, 3rd Edition. New York, NY:
384 American Council on Education and Macmillan, 1989.
- 385 16. Society of Academic Emergency Medicine. Join a committee. (*Accessed on July 5th at*
386 *<https://www.saem.org/about-saem/saem-membership/committees>*)
- 387 17. SurveyMonkey. (*Accessed on April 7th at <https://www.surveymonkey.com/>*)
- 388 Scherer RW, Langenberg P, von Elm E. Full publication of results initially presented in
389 abstracts. Cochrane Database Syst Rev. 2007 Apr;18(2):MR000005.
- 390 18. Dubosh NM, Jordan J, Yarris LM, et al. Critical Appraisal of Emergency Medicine
391 Educational Research: The Best Publications of 2016. AEM Educ Train 2018;3(1):58-73.

- 392 19. Lin M, Fisher J, Coates WC, et al. Critical Appraisal of Emergency Medicine Education
393 Research: The Best Publications of 2012. *Acad Emerg Med* 2014;21:322-333.
- 394 20. Farrell SE, Kuhn GJ, Coates WC, et al. Critical Appraisal of Emergency Medicine
395 Education Research: The Best Publications of 2013. *Acad Emerg Med* 2014;21(11):1274-83.
- 396 21. Yarris LM, Juve AM, Coates WC, et al. Critical Appraisal of Emergency Medicine
397 Education Research: The Best Publications of 2014. *Acad Emerg Med* 2015; 22(11):1327-36.
- 398 22. Heitz CR, Coates WC, Farrell SE, Fisher J, Juve AM, Yarris LM. Critical Appraisal of
399 Emergency Medicine Education Research: The Best Publications of 2015. *AEM Educ Train*
400 2017;1(4):255-268.
- 401 23. Moralejo D, Ogunremi T, Dunn K. Critical appraisal toolkit (CAT) for assessing multiple
402 types of evidence. *Can Commun Dis Rep* 2017;43(9):176-181.
- 403 24. Reed DA, Beckam TJ, Wright SM, Levine RB, Kern DE, Cook DA. Predictive validity
404 evidence for medical education research study quality instrument scores: quality of submissions
405 to JGIM's medical education special issue. *J Gen Intern Med* 2008;23(7):903-907.
- 406 25. Johnson JL, Adkins D, Chauvin S. A review of the quality indicators of rigor in qualitative
407 research. *American Journal of Pharmaceutical Education* 2020;84(1):7120.
- 408 26. Yang LJ, Chang KW, Chung KC. Methodologically rigorous clinical research. *Plast Reconstr*
409 *Surg* 2012;129(6):979e-988e.
- 410 27. Schneider NC, Coates WC, Yarris LM. Taking Your Qualitative Research to the Next Level:
411 A Guide for the Medical Educator. *AEM Educ Train* 2017;1(4):368-378.
- 412 28. O'Brien BC, Harris IB, Beckman TJ, Reed DA, Cook DA. Standards for reporting qualitative
413 research: a synthesis of recommendations. *Acad Med.* 2014;89(9):1245-51.

414

415

416

417

418 Table 1. Items and mean scores of expert review during Delphi process for quantitative scoring
419 rubric

420

Item	Mean score
------	------------

	(SD) N= 8
Round 1	
Clarity of objectives 0 = No clear objective or hypothesis 1 = Objective(s) are stated but unclear 2 = Clearly stated objective(s)	8.88 (0.35)
Quality of objectives 0 = No stated objective or hypothesis 1 = Poorly chosen objective(s) or stated hypothesis is difficult to test 2 = Well thought out study objective(s) or testable hypothesis	7.71 (1.70)
Study design 0 = Inappropriate study design for objective(s) 0.5 = Single group cross sectional or single group post-test only 1 = Single group pre-test and post-test 1.5 = Two or more non-randomized groups (quasi-experimental study) 2 = Two or more randomized groups (experimental study)	6.5 (2.07)
Sampling: institutions 0 = Single institution 2 = Multi-institutional	5.38 (1.92)
Sampling: response rate 0 = Less than 50% or not reported 1 = 50-74% 2 = Greater than or equal to 75%	5.29 (2.43)
Type of data 0 = Not described 1 = Assessment by study participant 1.5 = Subjective assessment by someone other than the study participant (i.e. an observer) 2 = Objective assessment	6.50 (2.67)
Power/sample size 0 = No power/sample size calculation was performed 2 = A power/sample size calculation was calculated and satisfied	5.63 (2.83)

<p>Data Analysis</p> <p>0 = No analysis described or inappropriate data analysis for study design</p> <p>1 = Descriptive analysis only (i.e. frequency, mean, median)</p> <p>2 = Beyond descriptive analysis (i.e. any comparative statistics or test of statistical inference)</p>	7.88 (0.99)
<p>Generalizability</p> <p>0 = not at all generalizable, results are only applicable to very specific population/setting</p> <p>0.5 = Minimally generalizable</p> <p>1 = Moderately generalizable</p> <p>1.5 = Very generalizable, results apply to most EM educational populations/settings</p> <p>2 = Extremely generalizable, results apply to educational populations/settings beyond EM</p>	8.13 (0.64)
<p>Relevance and importance of topic to medical education</p> <p>0 = This topic is only of interest to a very small group of people and is unlikely to result in important knowledge</p> <p>0.5 = This is an important topic to EM medical education that will lead to information of interest to many EM educators and learners</p> <p>1 = This topic is essential to EM medical education and is likely to be important and relevant for every EM educator and learner to know</p> <p>2 = This topic is essential to medical education other specialties beyond EM and is likely to be important for every medical educator and learner to know</p>	7.5 (2.73)
<p>Publication readiness/quality of writing</p> <p>0 = Poorly written, unclear, difficult to understand</p> <p>1 = Generally well-written, but leaves room for confusion on some concepts or has one or two errors</p> <p>2 = Exceptionally well-written, clear, logical organization and presentation of ideas.</p>	7.38 (1.85)
<p>Outcome(s)</p> <p>0.5 = Kirkpatrick level 1 – satisfaction, attitudes, perceptions, opinions, general facts (i.e. demographics)</p> <p>1 = Kirkpatrick level 2 – knowledge, skills (includes behaviors in a test setting such as simulation)</p> <p>1.5 = Kirkpatrick level 3 – behaviors in real context or clinical setting</p>	7.63 (2.13)

2 = Kirkpatrick level 4 = patient or health care outcome (actual effects on real patients, programs, or society)	
Innovation of study 0 = Not innovative or novel 1 = Moderately innovative (i.e. new method of instructing in a standard environment or standard instructional method in a novel area/environment) 2 = Completely novel idea	7.25 (1.39)
Global Rating 0 = No clear conclusions can be drawn 0.5 = Results ambiguous but appears to show a trend 1 = Conclusions can probably be based on results 1.5 = Results are clear and likely to be true 2 = Results are unequivocal	8.00 (1.31)
Round 2	
Quality of objectives 0 = No stated objective 1 = Poorly chosen or ambiguous objective(s) 2 = Clear, well thought out objective(s)	9.00 (0)
Appropriateness of methods 0 = Inappropriate methods for objective(s) 1 = Chosen methods were sub-optimal, but did address the objective(s) (i.e. acceptable methodology) 2 = Chosen methods were the best feasible for the objective(s) (i.e. rigorous methodology)	8.38 (1.06)
Study design 0 = Study design not described 0.5 = Single group cross sectional or single group post-assessment only 1 = Single group pre- and post-assessment 1.5 = Two or more non-randomized groups (quasi-experimental study) 2 = Two or more randomized groups (experimental study)	5.25 (2.66)
Data Analysis 0 = No analysis described or inappropriate data analysis for study design 1 = Descriptive analysis only (i.e. frequency, mean, median)	7.50 (1.31)

2 = Beyond descriptive analysis (i.e. any comparative statistics or test of statistical inference)	
<p>Generalizability</p> <p>0 = Results are only applicable to a very specific population/setting</p> <p>1 = Results are applicable to most EM educational populations/settings</p> <p>2 = Results are applicable to educational populations/settings beyond EM.</p>	7.00 (1.51)
<p>Relevance and importance of topic to medical education</p> <p>0 = This topic is only of interest to a very small group of people and is unlikely to result in important knowledge</p> <p>0.5 = This is an important topic to EM medical education that will lead to information of interest to many EM educators and learners</p> <p>1 = This topic is essential to EM medical education and is likely to be important and relevant for every EM educator and learner to know</p> <p>2 = This topic is essential to medical education in other specialties beyond EM and is likely to be important for every medical educator and learner to know</p>	7.00 (1.31)
<p>Publication readiness/quality of writing</p> <p>0 = Poorly written, unclear, difficult to understand</p> <p>1 = Generally well-written, but leaves room for confusion on some concepts or has one or two errors</p> <p>2= Exceptionally well-written, clear, logical organization and presentation of ideas.</p>	7.25 (2.05)
<p>Outcome(s)</p> <p>0 = Chosen outcomes are inappropriate for study objective</p> <p>0.5 = Kirkpatrick level 1 – satisfaction, attitudes, perceptions, opinions, general facts (i.e. demographics)</p> <p>1 = Kirkpatrick level 2 – knowledge, skills (includes behaviors in a test setting such as simulation)</p> <p>2 = Kirkpatrick level 3 – behaviors in real context or clinical setting</p> <p>3 = Kirkpatrick level 4 = patient or health care outcome (actual effects on real patients, programs, or society)</p>	6.25 (2.25)
<p>Innovation of study</p> <p>0 = Not innovative or novel</p> <p>1 = Moderately innovative (i.e. new method of instructing in a standard environment or standard instructional method in a novel area/environment)</p>	7.75 (1.04)

2 = Completely novel idea	
Strength of conclusion(s) 0 = No clear conclusions can be drawn 0.5 = Results ambiguous but appears to show a trend 1 = Conclusions can probably be based on results 1.5 = Conclusions are clear and likely to be true 2 = Conclusions are unequivocal	7.00 (1.51)
Round 3	
Quality of objectives 0 = No stated objective 1 = Poorly chosen or ambiguous objective(s) 2 = Clear, well thought out objective(s) that logically follow from the background information	8.63 (0.52)
Appropriateness of methods 0 = Inappropriate methods for objective(s) 1 = Chosen methods were sub-optimal, but did address the objective(s) 2 = Chosen methods were the best feasible for the objective(s) (i.e. rigorous methods)	8.75 (0.46)
Data analysis 0 = No analysis described or inappropriate data analysis for study design 1 = Descriptive analysis only (e.g frequency, mean, median) 2 = Beyond descriptive analysis (e.g. any comparative statistics or test of statistical inference)	8.38 (0.74)
Generalizability 0 = Results are only applicable to a very specific population/setting 1 = Results are applicable to most EM educational populations/settings 2 = Results are applicable to educational populations/settings beyond EM.	7.25 (1.58)
Relevance and importance of topic to medical education 0 = This topic is only of interest to a very small group of people and is unlikely to result in important knowledge 1 = This topic is essential to EM medical education and is likely to be important and relevant for every EM educator and learner to know 2 = This topic is essential to medical education in other specialties beyond EM and is likely to be important for every medical educator and learner to know	6.88 (1.46)

Quality of writing 0 = Poorly written, unclear, difficult to understand 0.5 = Generally well-written 1 = Exceptionally well-written, clear, logical organization and presentation of ideas.	7.50 (1.93)
Outcome(s) 0 = Chosen outcomes are inappropriate for study objective 1 = Chosen outcomes are reasonable for study objective, but not the best measure 2 = Chosen outcomes are ideal for study objective	8.50 (0.93)
Innovation of study 0 = Not innovative or novel 1 = Moderately innovative (e.g. new method of instructing in a standard environment or standard instructional method in a novel area/environment) 2 = Completely novel idea (e.g. new method of instructing in a novel area/environment)	7.63 (1.19)
Strength of conclusion(s) 0 = No clear conclusions can be drawn or conclusions do not follow directly from results 1 = Conclusions can probably be based on results 2 = Conclusions are unequivocal	8.25 (0.89)

421

422

423

424

425

426

427

428

429

430 Table 2. Items and mean scores of expert review during Delphi process for qualitative scoring

431 rubric

432

Item	Mean score (SD)
------	--------------------

	N = 8
Round 1	
<p>Quality of objectives</p> <p>0 = No stated objective</p> <p>1 = Poorly chosen or ambiguous objective(s)</p> <p>2 = Clear, well thought out objective(s) that logically follow from the background information</p>	8.13 (1.36)
<p>Study design</p> <p>0 = Qualitative design is not appropriate for study objective(s)</p> <p>1 = Qualitative approach is appropriate for study objective, but specific design not identified (i.e. phenomenology, ethnography, grounded theory, etc.)</p> <p>2 = Specific qualitative design identified and appropriate for study objective</p>	8.25 (0.89)
<p>Data collection methods</p> <p>0 = Data collection methods (participant observation, interviews, document review, etc.) not identified</p> <p>1 = Data collection methods identified but inappropriate for study objective</p> <p>2 = Data collection methods identified and appropriate for study objective</p>	7.88 (1.64)
<p>Sampling: method (Sampling is defined as the process of selecting participants)</p> <p>0 = sampling method not described</p> <p>1 = sampling method described, but not clear or not theoretically justified</p> <p>2 = Clear description of sampling method that is theoretically justified</p>	7.25 (1.49)
<p>Sampling: saturation (Saturation is defined as the point at which no new information is being learned from continued data collection)</p> <p>0 = Saturation of data not achieved or not described</p> <p>2 = Saturation of data achieved</p>	4.75 (2.92)
<p>Trustworthiness (Trustworthiness is a marker of quality and can be supported with evidence of credibility, transferability, dependability, and confirmability)</p> <p>0 = No clear description of researcher role, study context, or triangulation</p> <p>1 = Provides some evidence of trustworthiness, but not comprehensive</p> <p>2 = Provides significant evidence of trustworthiness such as clear description of researcher role, study context and triangulation</p>	6.75 (1.49)
<p>Data Analysis</p> <p>0 = No analysis described or inappropriate data analysis for study objectives/design</p>	7.50 (2.00)

1 = Some description of data analyses, but not entirely clear 2= In depth description of systematic data analyses appropriate to study objective with clear description of how themes and concepts were derived	
Relevance and importance of topic to medical education 0 = This topic is only of interest to a very small group of people and is unlikely to result in important knowledge 1 = This topic is essential to EM medical education and is likely to be important and relevant for every EM educator and learner to know 2 = This topic is essential to medical education in other specialties beyond EM and is likely to be important for every medical educator and learner to know	7.50 (2.07)
Quality of writing 0 = Poorly written, unclear, difficult to understand 0.5 = Generally well-written 1= Exceptionally well-written, clear, logical organization and presentation of ideas.	7.50 (2.00)
Innovation of study 0 = Not innovative or novel 1 = Moderately innovative 2 = Innovative or novel	6.00 (2.51)
Strength of conclusion(s) 0 = No clear conclusions can be drawn or conclusions do not follow directly from results 1 = Conclusions can probably be based on results 2 = Conclusions are unequivocal	7.63 (1.69)
Round 2	
Quality of study aims/objectives 0 = No stated aim or objective 1 = Poorly chosen or ambiguous aim/objective(s) 2 = Clear, well thought out aim/objective(s) that logically follow from the background information	8.75 (0.46)
General methods 0 = Qualitative methods are not appropriate for study aim/objective(s) 1 = Qualitative methods are appropriate for study aim/objective(s), but specific approach (e.g. phenomenology, ethnography, grounded theory, etc.) or paradigm (e.g. postpositivist, constructivist/interpretivist) not stated or not ideal	8.13 (0.83)

2 = Specific qualitative approach and paradigm stated and aligned with study aim/objective(s)	
<p>Data collection</p> <p>0 = Data collection methods (observation, interviews, document review, etc.) not identified or inappropriate for study aim/objective(s)</p> <p>1 = Data collection methods appropriate for study aim/objective(s), but not ideal</p> <p>2 = Data collection methods are ideal for study aim/objective(s)</p>	7.63 (1.06)
<p>Sampling (Sampling is defined as the process of selecting participants)</p> <p>0 = Sampling not described</p> <p>1 = Sampling described, but flawed (e.g. unclear, inappropriate, not theoretically justified)</p> <p>2 = Sampling clearly described and theoretically justified</p>	7.50 (0.76)
<p>Trustworthiness (Trustworthiness is a marker of quality and can be supported with evidence of credibility, transferability, dependability, confirmability, and reflexivity. Examples of specific techniques used to enhance trustworthiness include member checking, audit trail, triangulation, etc.)</p> <p>0 = No clear description of methods to enhance trustworthiness.</p> <p>1 = Provides some evidence of trustworthiness, but not comprehensive</p> <p>2 = Provides significant evidence of trustworthiness such as clear description of researcher role, member checking, audit trail, study context or triangulation, with supported rationale</p>	6.88 (2.59)
<p>Data Analysis</p> <p>0 = No analysis described or inappropriate data analysis for study objectives/design</p> <p>1 = Some description of data analyses, but unclear or not justified</p> <p>2 = In depth description of systematic data analyses appropriate to study objective with clear description of how themes and concepts were derived</p>	7.75 (1.39)
<p>Importance of topic to medical education</p> <p>0 = This topic is unlikely to result in important knowledge</p> <p>1 = This topic is essential to EM medical education and is likely to be important for EM educators and learners to know</p> <p>2 = This topic is essential to medical education in other specialties beyond EM and is likely to be important for medical educators and learners to know</p> <p>2 = This topic is essential to medical education in other specialties beyond EM</p>	6.38 (2.50)

and is likely to be important for medical educators and learners to know	
Quality of writing 0 = Poorly written, unclear, difficult to understand 0.5 = Generally well-written 1 = Consistently well-written, clear, logical organization and presentation of ideas.	7.13 (2.36)
Strength of conclusion(s) 0 = No clear conclusions can be drawn or conclusions do not follow directly from results 1 = Conclusions can probably be based on results, but inference is necessary to draw conclusions 2 = Conclusions are well supported by results	8.25 (0.89)
Study implications 0 = Does not provide valuable information for future research 1 = Provides information that contributes to the field, but has limited implications for future research 2 = Provides a foundation for future hypothesis testing research	6.00 (1.93)
Round 3	
Quality of study aims/objectives 0 = No stated aim or objective 1 = Poorly chosen or ambiguous aim/objective(s) 2 = Clear, well thought out aim/objective(s) that logically follow from the background information	8.88 (0.35)
General methods 0 = Qualitative methods not appropriate for study aim/objective(s) 1 = Qualitative methods appropriate for study aim/objective(s), but specific approach (e.g. phenomenology, ethnography, grounded theory, etc.) or paradigm (e.g. postpositivist, constructivist/interpretivist) not stated or not ideal 2 = Specific qualitative approach and paradigm stated and aligned with study aim/objective(s)	8.38 (0.52)
Data collection 0 = Data collection methods (observation, interviews, document review, etc.) not identified or inappropriate for study aim/objective(s) 1 = Data collection methods appropriate for study aim/objective(s), but not ideal 2 = Data collection methods ideal for study aim/objective(s)	8.00 (1.07)

<p>Sampling (Sampling is defined as the process of selecting participants)</p> <p>0 = Sampling not described</p> <p>1 = Sampling described, but flawed (e.g. unclear, inappropriate, not theoretically justified)</p> <p>2 = Sampling clearly described and theoretically justified</p>	7.63 (0.74)
<p>Data Analysis</p> <p>0 = No analysis described or inappropriate data analysis for study objectives/design</p> <p>1 = Some description of data analyses, but unclear or not justified</p> <p>2= In depth description of systematic data analyses appropriate to study objective with clear description of how themes and concepts were derived</p>	8.50 (0.76)
<p>Quality of writing</p> <p>0 = Poorly written, unclear, difficult to understand</p> <p>1= Consistently well-written, clear, logical organization and presentation of ideas.</p>	8.00 (1.20)
<p>Strength of conclusion(s)</p> <p>0 = No clear conclusions can be drawn or conclusions do not follow directly from results</p> <p>1 = Conclusions require reader inference to draw conclusions</p> <p>2 = Conclusions are well supported by results</p>	8.38 (0.74)

433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449

450
 451
 452
 453
 454
 455
 456
 457
 458
 459
 460
 461
 462
 463
 464
 465
 466
 467
 468
 469
 470
 471
 472
 473
 474
 475
 476
 477
 478
 479
 480

Table 3. Interrater Reliability Results During Piloting

Item	Intraclass correlation coefficient [95% CI]
Quantitative rubric	

1. Quality of objectives	0.406 [-0.006, 0.705]
2. Appropriateness of methods	0.821 [0.671, 0.916]
3. Outcome(s)	0.661 [0.365, 0.843]
4. Data analysis	0.753 [0.548, 0.883]
5. Generalizability	0.878 [0.767, 0.944]
6. Relevance and importance of topic to medical education	0.747 [0.530, 0.882]
7. Innovation of study	0.786 [0.607, 0.900]
8. Quality of writing	0.726 [0.500, 0.870]
9. Strength of conclusions	0.739 [0.512, 0.878]
Total score	0.893 [0.802, 0.950]
Qualitative rubric	
1. Quality of objectives	0.176 [-0.466, 0.711]
2. General methods	0.897 [0.749, 0.971]
3. Data collection	0.635 [0.158, 0.892]
4. Sampling	0.531 [-0.106, 0.863]
5. Data analysis	0.874 [0.574, 0.950]
6. Quality of writing	-0.301 [-1.083, 0.489]
7. Strength of conclusions	0.753 [0.415, 0.927]
Total score	0.788 [0.469, 0.939]

481

482

483

484

485

486

487

488

489

490

491

492

493

494 Appendix A. Quantitative Education Research Abstract Scoring Rubric

495

496 1. Quality of objectives

497 0 = No stated objective

498 1 = Poorly chosen or ambiguous objective(s)

499 2 = Clear, well thought out objective(s) that logically follow from the background information

500

501 2. Appropriateness of methods

502 0 = Inappropriate methods for objective(s)

503 1 = Chosen methods were sub-optimal, but did address the objective(s)

504 2 = Chosen methods were the best feasible for the objective(s) (i.e. rigorous methods)

505

506 3. Outcome(s)

507 0 = Chosen outcomes are inappropriate for study objective

508 1 = Chosen outcomes are reasonable for study objective, but not the best measure

509 2 = Chosen outcomes are ideal for study objective

510

511 4. Data analysis

512 0 = No analysis described or inappropriate data analysis for study design

513 1 = Descriptive analysis only (e.g frequency, mean, median)

514 2 = Beyond descriptive analysis (e.g. any comparative statistics or test of statistical inference)

515

516 5. Generalizability

517 0 = Results are only applicable to a very specific population/setting

518 1 = Results are applicable to most EM educational populations/settings

519 2 = Results are applicable to educational populations/settings beyond EM.

520

521 6. Relevance and importance of topic to medical education

522 0 = This topic is only of interest to a very small group of people and is unlikely to result in

523 important knowledge

524 1 = This topic is essential to EM medical education and is likely to be important and relevant for
525 every EM educator and learner to know

526 2 = This topic is essential to medical education in other specialties beyond EM and is likely to be
527 important for every medical educator and learner to know

528

529 7. Innovation of study

530 0 = Not innovative or novel

531 1 = Moderately innovative (e.g. new method of instructing in a standard environment or standard
532 instructional method in a novel area/environment)

533 2 = Completely novel idea (e.g. new method of instructing in a novel area/environment)

534

535 8. Quality of writing

536 0 = Poorly written, unclear, difficult to understand

537 0.5 = Generally well-written

538 1 = Exceptionally well-written, clear, logical organization and presentation of ideas.

539

540 9. Strength of conclusion(s)

541 0 = No clear conclusions can be drawn or conclusions do not follow directly from results

542 1 = Conclusions can probably be based on results

543 2 = Conclusions are unequivocal

544

545

546

547

548

549

550

551

552

553

554

555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585

Appendix B. Qualitative Education Research Abstract Scoring Rubric

Are you familiar with qualitative research study design?

Yes: Proceed with scoring

No: Decline

1. Quality of study aims/objectives

0 = No stated aim or objective

1 = Poorly chosen or ambiguous aim/objective(s)

2 = Clear, well thought out aim/objective(s) that logically follow from the background information

2. General methods

0 = Qualitative methods not appropriate for study aim/objective(s)

1 = Qualitative methods appropriate for study aim/objective(s), but specific approach (e.g. phenomenology, ethnography, grounded theory, etc.) or paradigm (e.g. postpositivist, constructivist/interpretivist) not stated or not ideal

2 = Specific qualitative approach and paradigm stated and aligned with study aim/ objective(s)

3. Data collection

0 = Data collection methods (observation, interviews, document review, etc.) not identified or inappropriate for study aim/objective(s)

1 = Data collection methods appropriate for study aim/objective(s), but not ideal

- 586 2 = Data collection methods ideal for study aim/objective(s)
587
588 4. Sampling (Sampling is defined as the process of selecting participants)
589 0 = Sampling not described
590 1 = Sampling described, but flawed (e.g. unclear, inappropriate, not theoretically justified)
591 2 = Sampling clearly described and theoretically justified
592
593 5. Data Analysis
594 0 = No analysis described or inappropriate data analysis for study objectives/design
595 1 = Some description of data analyses, but unclear or not justified
596 2= In depth description of systematic data analyses appropriate to study objective with clear
597 description of how themes and concepts were derived
598
599 6. Quality of writing
600 0 = Poorly written, unclear, difficult to understand
601 1= Consistently well-written, clear, logical organization and presentation of ideas.
602
603 7. Strength of conclusion(s)
604 0 = No clear conclusions can be drawn or conclusions do not follow directly from results
605 1 = Conclusions require reader inference to draw conclusions
606 2 = Conclusions are well supported by results
607
608