






ORIGINAL CONTRIBUTION

Leveling the field: Development of reliable scoring rubrics for quantitative and qualitative medical education research abstracts

Jaime Jordan MD, MAEd^{1,2}  | Laura R. Hopson MD³  | Caroline Molins MD, MSMEd⁴ | Suzanne K. Bentley MD, MPH⁵  | Nicole M. Deiorio MD⁶ | Sally A. Santen MD, PhD^{6,7}  | Lalena M. Yarris MD, MCR⁸  | Wendy C. Coates MD¹  | Michael A. Gisondi MD⁹ 

¹Department of Emergency Medicine, David Geffen School of Medicine at UCLA, Los Angeles, California, USA

²Department of Emergency Medicine, Ronald Reagan UCLA Medical Center, Los Angeles, California, USA

³Department of Emergency Medicine, University of Michigan, Ann Arbor, Michigan, USA

⁴AdventHealth Emergency Medicine Residency, Orlando, Florida, USA

⁵Icahn School of Medicine at Mount Sinai, New York, New York, USA

⁶Virginia Commonwealth University School of Medicine, Richmond, Virginia, USA

⁷University of Cincinnati College of Medicine, Cincinnati, Ohio, USA

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

⁹Department of Emergency Medicine, Stanford University, Palo Alto, California, USA

Correspondence

Jaime Jordan, UCLA Emergency Medicine, 924 Westwood Boulevard, Suite 300, Los Angeles, CA 90024, USA.
Email: jaimejordanmd@gmail.com

Abstract

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

Methods: We used a modified Delphi method to achieve expert consensus for scoring rubric items to optimize content validity. Eight education research experts participated in two separate modified Delphi processes, one to generate quantitative research items and one for qualitative. Modifications were made between rounds based on item scores and expert feedback. Homogeneity of ratings in the Delphi process was calculated using Cronbach's alpha, with increasing homogeneity considered an indication of consensus. Rubrics were piloted by scoring abstracts from 22 quantitative publications from *AEM Education and Training* "Critical Appraisal of Emergency Medicine Education Research" (11 highlighted for excellent methodology and 11 that were not) and 10 qualitative publications (five highlighted for excellent methodology and five that were not). Intraclass correlation coefficient (ICC) estimates of reliability were calculated.

Results: Each rubric required three rounds of a modified Delphi process. The resulting quantitative rubric contained nine items: quality of objectives, appropriateness of methods, outcomes, data analysis, generalizability, importance to medical education, innovation, quality of writing, and strength of conclusions (Cronbach's α for the third round = 0.922, ICC for total scores during piloting = 0.893). The resulting qualitative rubric contained seven items: quality of study aims, general methods, data collection, sampling, data analysis, writing quality, and strength of conclusions (Cronbach's α for the third round = 0.913, ICC for the total scores during piloting = 0.788).

Presented at Society for Academic Emergency Medicine Virtual Meeting, May 13, 2021.

Supervising Editor: Esther H. Chen, MD.

© 2021 by the Society for Academic Emergency Medicine

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

INTRODUCTION

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

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

Prior literature has described tools to assess quality in medical education research manuscripts, such as the Medical Education Research Study Quality Instrument (MERSQI) and the Newcastle-Ottawa Scale-Education (NOS-E).⁸ A limited attempt to utilize the MERSQI tool to retrospectively assess internal medicine medical education abstracts achieving manuscript publication showed increased scores for the journal abstract relative to the conference abstract.⁴ However, the MERSQI and similar tools were not developed specifically for judging abstracts, and there is a lack of

published validity evidence to support score interpretation based on these tools. To equitably assess the quality of education research abstracts to scholarly venues, which may have downstream effects on researcher scholarship, advancement, and reputation, there is a need for a rigorously developed abstract scoring rubric that is based on a validity evidence framework.^{9,10}

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

METHODS

Study design

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

Study setting and population

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

Study protocol

We followed Messick's framework for validity that includes five types of validity evidence; content, response process, internal

structure, relation to other variables, and consequential.¹⁵ Our study team drafted initial items for the scoring rubrics after a review of the literature and existing research abstract scoring rubrics to optimize content validity. We created separate items for research abstracts reporting quantitative and qualitative data. We sent the draft items to the Society for Academic Emergency Medicine (SAEM) education committee for review and comment to gather stakeholder feedback and for further content and response process validity evidence.¹⁶ One author (JJ) who was not a member of the Delphi panel then revised the initial lists of items based on committee feedback to create the initial Delphi surveys. We used an electronic survey platform (SurveyMonkey) to administer and collect data from the Delphi surveys.¹⁷ Experts on the Delphi panel rated the importance of including each item in a scoring rubric on a 1 to 9 Likert scale with 1 labeled as “not at all important” and 9 labeled as “extremely important.” The experts were invited to provide additional written comments, edits, and suggestions for each item. They were also encouraged to suggest additional items that they felt were important but not currently listed. We determined a priori that items with a mean score of 7 or greater advanced to the next round and items with a mean score of three or below were eliminated. The Delphi panel moderator (JJ) applied discretion for items scoring between 4 and 6, with the aim of both adhering to the opinions of the experts and creating a comprehensive scoring rubric. For example, if an item received a middle score but had comments supporting inclusion in a revised form, the moderator would make the suggested revisions and include the item in the next round.

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

We considered homogeneity of ratings in the Delphi process to be an indication of consensus. After consensus was achieved, we created final scoring rubrics for quantitative and qualitative medical education research abstracts. We then piloted the scoring rubrics to gather internal structure and further response process validity evidence. Five raters from the study group (JJ, LH, MG, CM, SB) participated in piloting. We piloted the final quantitative research rubric by scoring abstracts from publications identified in the most recent critical appraisal of EM education research by *Academic Emergency Medicine/AEM Education and Training*, “Critical Appraisal of Emergency Medicine Education Research: The Best Publications of 2016”.¹⁸ All 11 papers highlighted for excellent methodology in this issue were included in the pilot.¹⁸ Additionally, we included an equal number of randomly selected citations that were included in the issue but not selected as top papers, for a total of 22 quantitative publications.¹⁸ Given the limited number of qualitative studies cited in this issue of the critical appraisal series, we chose to pilot

the qualitative rubric on publications from this series from the last 5 years available (2012–2016).^{18–22} We randomly selected one qualitative publication that was highlighted for excellent methodology and one that was not from each year for a total of 10 qualitative publications.^{18–22} The same five raters who performed the quantitative pilot also conducted the qualitative pilot.

Data analysis

We calculated and reported descriptive statistics for item scoring during Delphi rounds. We used Cronbach's alpha to assess homogeneity of ratings in the Delphi process. Increasing homogeneity was considered to be an indication of consensus among the expert panelists. We used intraclass correlation coefficient (ICC) estimates to assess reliability among raters during piloting based on a mean rating ($\kappa = 5$), absolute agreement, two-way random-effects model. We performed all analyses in SPSS (IBM SPSS Statistics for Windows, Version 27.0).

RESULTS

Quantitative rubric

Three Delphi rounds were completed, each with 100% response rate. Mean item scores for each round are depicted in Table 1. After the first round, three items were deleted, one item was added, and five items underwent wording changes. After the second round, one item was deleted and eight items underwent wording changes. After the third round items were reordered for flow and ease of use but no further changes were made to content or wording. Cronbach's alpha for the third round was 0.922, indicating high internal consistency. The final rubric contained nine items: quality of objectives, appropriateness of methods, outcomes, data analysis, generalizability, importance to medical education, innovation, quality of writing, and strength of conclusions (Data Supplement S1, Appendix S1, available as supporting information in the online version of this paper, which is available at <http://onlinelibrary.wiley.com/doi/10.1002/aet2.10654/full>). The ICC for the total scores during piloting was 0.893, indicating excellent agreement. ICCs for individual rubric items ranged from 0.406 to 0.878 (Table 3).

Qualitative rubric

Three Delphi rounds were completed, each with 100% response rate. Mean item scores for each round are depicted in Table 2. After the first round 2 items were deleted, one item was added and nine items underwent wording changes. After the second round, three items were deleted and four underwent wording changes. After the third round no further changes were made. The resulting tool contained seven items reflecting the domains of quality of study aims, general

TABLE 1 Items and mean scores of expert review during Delphi process for quantitative scoring rubric

Item	Mean score (\pm SD), N = 8
<i>Round 1</i>	
Clarity of objectives	8.88 (\pm 0.35)
0 = No clear objective or hypothesis	
1 = Objective(s) are stated but unclear	
2 = Clearly stated objective(s)	
Quality of objectives	7.71 (\pm 1.70)
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	
Study design	6.5 (\pm 2.07)
0 = Inappropriate study design for objective(s)	
0.5 = Single-group cross-sectional or single-group posttest only	
1 = Single-group pretest and posttest	
1.5 = Two or more nonrandomized groups (quasi-experimental study)	
2 = Two or more randomized groups (experimental study)	
Sampling: institutions	5.38 (\pm 1.92)
0 = Single institution	
2 = Multi-institutional	
Sampling: response rate	5.29 (\pm 2.43)
0 = Less than 50% or not reported	
1 = 50%–74%	
2 = Greater than or equal to 75%	
Type of data	6.50 (\pm 2.67)
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	
Power/sample size	5.63 (\pm 2.83)
0 = No power/sample size calculation was performed	
2 = A power/sample size calculation was calculated and satisfied	
Data analysis	7.88 (\pm 0.99)
0 = No analysis described or inappropriate data analysis for study design	
1 = Descriptive analysis only (i.e., frequency, mean, median)	
2 = Beyond descriptive analysis (i.e., any comparative statistics or test of statistical inference)	
Generalizability	8.13 (\pm 0.64)

(Continues)

TABLE 1 (Continued)

Item	Mean score (\pm SD), N = 8
0 = Not at all generalizable, results are only applicable to very specific population/setting	
0.5 = Minimally generalizable	
1 = Moderately generalizable	
1.5 = Very generalizable, results apply to most EM educational populations/settings	
2 = Extremely generalizable, results apply to educational populations/settings beyond EM	
Relevance and importance of topic to medical education	7.5 (\pm 2.73)
0 = This topic is only of interest to a very small group of people and is unlikely to result in important knowledge	
0.5 = This is an important topic to EM medical education that will lead to information of interest to many EM educators and learners	
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 other specialties beyond EM and is likely to be important for every medical educator and learner to know	
Publication readiness/quality of writing	7.38 (\pm 1.85)
0 = Poorly written, unclear, difficult to understand	
1 = Generally well-written, but leaves room for confusion on some concepts or has one or two errors	
2 = Exceptionally well-written, clear, logical organization, and presentation of ideas.	
Outcome(s)	7.63 (\pm 2.13)
0.5 = Kirkpatrick level 1—satisfaction, attitudes, perceptions, opinions, general facts (i.e., demographics)	
1 = Kirkpatrick level 2—knowledge, skills (includes behaviors in a test setting such as simulation)	
1.5 = Kirkpatrick level 3—behaviors in real context or clinical setting	
2 = Kirkpatrick level 4 = patient or health care outcome (actual effects on real patients, programs, or society)	
Innovation of study	7.25 (\pm 1.39)
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	
Global rating	8.00 (\pm 1.31)

(Continues)

TABLE 1 (Continued)

Item	Mean score (±SD), N = 8
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	
Round 2	
Quality of objectives	9.00 (±0)
0 = No stated objective	
1 = Poorly chosen or ambiguous objective(s)	
2 = Clear, well-thought-out objective(s)	
Appropriateness of methods	8.38 (±1.06)
0 = Inappropriate methods for objective(s)	
1 = Chosen methods were suboptimal, 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)	
Study design	5.25 (±2.66)
0 = Study design not described	
0.5 = Single group cross-sectional or single-group postassessment only	
1 = Single-group pre- and postassessment	
1.5 = Two or more nonrandomized groups (quasi-experimental study)	
2 = Two or more randomized groups (experimental study)	
Data analysis	7.50 (±1.31)
0 = No analysis described or inappropriate data analysis for study design	
1 = Descriptive analysis only (i.e., frequency, mean, median)	
2 = Beyond descriptive analysis (i.e., any comparative statistics or test of statistical inference)	
Generalizability	7.00 (±1.51)
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	
Relevance and importance of topic to medical education	7.00 (±1.31)
0 = This topic is only of interest to a very small group of people and is unlikely to result in important knowledge	
0.5 = This is an important topic to EM medical education that will lead to information of interest to many EM educators and learners	

(Continues)

TABLE 1 (Continued)

Item	Mean score (±SD), N = 8
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	
Publication readiness/quality of writing	7.25 (±2.05)
0 = Poorly written, unclear, difficult to understand	
1 = Generally well written, but leaves room for confusion on some concepts or has one or two errors	
2 = Exceptionally well-written, clear, logical organization and presentation of ideas	
Outcome(s)	6.25 (±2.25)
0 = Chosen outcomes are inappropriate for study objective	
0.5 = Kirkpatrick level 1—satisfaction, attitudes, perceptions, opinions, general facts (i.e., demographics)	
1 = Kirkpatrick level 2—knowledge, skills (includes behaviors in a test setting such as simulation)	
2 = Kirkpatrick level 3—behaviors in real context or clinical setting	
3 = Kirkpatrick level 4—patient or health care outcome (actual effects on real patients, programs, or society)	
Innovation of study	7.75 (±1.04)
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	
Strength of conclusion(s)	7.00 (±1.51)
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	
Round 3	
Quality of objectives	8.63 (±0.52)
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	

(Continues)

TABLE 1 (Continued)

Item	Mean score (\pm SD), N = 8
Appropriateness of methods 0 = Inappropriate methods for objective(s) 1 = Chosen methods were suboptimal, but did address the objective(s) 2 = Chosen methods were the best feasible for the objective(s) (i.e., rigorous methods)	8.75 (\pm 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 (\pm 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 (\pm 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 (\pm 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 (\pm 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 (\pm 0.93)
Innovation of study 0 = Not innovative or novel	7.63 (\pm 1.19)

(Continues)

TABLE 1 (Continued)

Item	Mean score (\pm SD), N = 8
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)	
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 (\pm 0.89)

methods, data collection, sampling, data analysis, writing quality, and strength of conclusions (Appendix S2). Cronbach's alpha for the third round was 0.913, indicating high internal consistency. ICC for the total scores during piloting was 0.788, indicating good agreement. The item on writing quality had an ICC of -0.301 , likely due to the small scale of the item and sample size leading to limited variance. ICCs for the remainder of the items ranged from 0.176 to 0.897 (Table 3).

DISCUSSION

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

We deliberately created separate rubrics for abstracts reporting quantitative and qualitative data because each has unique methods. When designing a study, education researchers must decide the best method to address their questions. Often, in the exploratory phase of inquiry, a qualitative study is the most appropriate choice to identify key topics that merit further study. These often may be narrow in scope and may employ one or more qualitative methods (e.g., ethnography, focus groups, personal interviews). The careful and rigorous analysis may reveal points that can be studied

TABLE 2 Items and mean scores of expert review during Delphi process for qualitative scoring rubric

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

(Continues)

TABLE 2 (Continued)

Item	Mean score (±SD), N = 8
0 = No analysis described or inappropriate data analysis for study objectives/design	
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	7.50 (±2.07)
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	
Quality of writing	7.50 (±2.00)
0 = Poorly written, unclear, difficult to understand	
0.5 = Generally well written	
1 = Exceptionally well-written, clear, logical organization and presentation of ideas	
Innovation of study	6.00 (±2.51)
0 = Not innovative or novel	
1 = Moderately innovative	
2 = Innovative or novel	
Strength of conclusion(s)	7.63 (±1.69)
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	
<i>Round 2</i>	
Quality of study aims/objectives	8.75 (±0.46)
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	
General methods	8.13 (±0.83)
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) or paradigm (e.g., postpositivist, constructivist/interpretivist) not stated or not ideal	

(Continues)

TABLE 2 (Continued)

Item	Mean score (\pm SD), N = 8
2 = Specific qualitative approach and paradigm stated and aligned with study aim/ objective(s)	
Data collection	7.63 (\pm 1.06)
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 are ideal for study aim/objective(s)	
Sampling (sampling is defined as the process of selecting participants)	7.50 (\pm 0.76)
0 = Sampling not described	
1 = Sampling described, but flawed (e.g., unclear, inappropriate, not theoretically justified)	
2 = Sampling clearly described and theoretically justified	
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.)	6.88 (\pm 2.59)
0 = No clear description of methods to enhance trustworthiness	
1 = Provides some evidence of trustworthiness, but not comprehensive	
2 = Provides significant evidence of trustworthiness such as clear description of researcher role, member checking, audit trail, study context, or triangulation, with supported rationale	
Data analysis	7.75 (\pm 1.39)
0 = No analysis described or inappropriate data analysis for study objectives/design	
1 = Some description of data analyses, but unclear or not justified	
2 = In-depth description of systematic data analyses appropriate to study objective with clear description of how themes and concepts were derived	
Importance of topic to medical education	6.38 (\pm 2.50)
0 = This topic is unlikely to result in important knowledge	
1 = This topic is essential to EM medical education and is likely to be important for EM educators and learners to know	
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	

(Continues)

TABLE 2 (Continued)

Item	Mean score (\pm SD), N = 8
Quality of writing	7.13 (\pm 2.36)
0 = Poorly written, unclear, difficult to understand	
0.5 = Generally well written	
1 = Consistently well-written, clear, logical organization and presentation of ideas	
Strength of conclusion(s)	8.25 (\pm 0.89)
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	
Study implications	6.00 (\pm 1.93)
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	
<i>Round 3</i>	
Quality of study aims/objectives	8.88 (\pm 0.35)
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	
General methods	8.38 (\pm 0.52)
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) 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)	
Data collection	8.00 (\pm 1.07)
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)	
Sampling (sampling is defined as the process of selecting participants)	7.63 (\pm 0.74)

(Continues)

TABLE 2 (Continued)

Item	Mean score (±SD), N = 8
0 = Sampling not described	
1 = Sampling described, but flawed (e.g., unclear, inappropriate, not theoretically justified)	
2 = Sampling clearly described and theoretically justified	
Data analysis	8.50 (±0.76)
0 = No analysis described or inappropriate data analysis for study objectives/design	
1 = Some description of data analyses, but unclear or not justified	
2 = In-depth description of systematic data analyses appropriate to study objective with clear description of how themes and concepts were derived	
Quality of writing	8.00 (±1.20)
0 = Poorly written, unclear, difficult to understand	
1 = Consistently well-written, clear, logical organization and presentation of ideas	
Strength of conclusion(s)	8.38 (±0.74)
0 = No clear conclusions can be drawn or conclusions do not follow directly from results	
1 = Conclusions require reader inference to draw conclusions	
2 = Conclusions are well supported by results	

later via quantitative methods to test a hypothesis gleaned during the qualitative phase.²⁷ Specific standards for reporting on qualitative research have been widely disseminated and are distinct from standards for reporting quantitative research.²⁸ Even an impeccably designed and executed qualitative study would fail to meet major criteria for excellent quantitative studies. For example, points may be subtracted for lack of generalizability or conduct of the qualitative study in multiple institutions as well as for the absence of common quantitative statistical analytics. The qualitative abstract itself may necessarily lack the common structure of a quantitative report and lead to a lower score. The obvious problem is that a well-conducted study might not be shared with the relevant research community if it is judged according to quantitative standards. A similar outcome would occur if quantitative work were judged by qualitative standards; therefore, we advocate for using scoring rubrics specific to the type of research being assessed.

Our work has several possible applications. The rubrics we developed may be adopted as scoring tools for medical education research studies that are submitted for presentation to scientific conferences. The presence of specific scoring rubrics for medical education research may address disparities in acceptance rates and ensure presentation of rigorously conducted medical education research at scientific conferences. Further, publication of abstract scoring rubrics such as ours sets expectations for certain elements

TABLE 3 Inter-rater reliability results during piloting

Item	ICC [95% CI]
Quantitative rubric	
1. Quality of objectives	0.406 [-0.006 to 0.705]
2. Appropriateness of methods	0.821 [0.671 to 0.916]
3. Outcome(s)	0.661 [0.365 to 0.843]
4. Data analysis	0.753 [0.548 to 0.883]
5. Generalizability	0.878 [0.767 to 0.944]
6. Relevance and importance of topic to medical education	0.747 [0.530 to 0.882]
7. Innovation of study	0.786 [0.607 to 0.900]
8. Quality of writing	0.726 [0.500 to 0.870]
9. Strength of conclusions	0.739 [0.512 to 0.878]
Total score	0.893 [0.802 to 0.950]
Qualitative rubric	
1. Quality of objectives	0.176 [-0.466 to 0.711]
2. General methods	0.897 [0.749 to 0.971]
3. Data collection	0.635 [0.158 to 0.892]
4. Sampling	0.531 [-0.106 to 0.863]
5. Data analysis	0.874 [0.574 to 0.950]
6. Quality of writing	-0.301 [-1.083 to 0.489]
7. Strength of conclusions	0.753 [0.415 to 0.927]
Total score	0.788 [0.469 to 0.939]

to be included and defines an acceptable level of submission quality. Dissemination and usage of the rubrics may therefore help improve research excellence. The rubrics themselves can serve as educational tools in resident and faculty training. For example, the rubrics could serve as illustrations or practice material in teaching how to prepare a strong abstract for submission. The inclusive wording of the items allows the rubrics to be adapted to medical education work in any medical specialty. Medical educators may also benefit from using the methods described here to create their own scoring rubrics or provide evidence-based best practice approaches for other venues. Finally, this study provides a tool that could lay the groundwork for future scholarship on assessing the quality of educational research.

LIMITATIONS

Our study has several limitations. First, the modified Delphi technique is a consensus technique that can force agreement of respondents, and the existence of consensus does not denote a correct response.¹¹ Since the method is implemented electronically, there is limited discussion and elaboration. Second, the team of experts were all researchers in EM; therefore, the rubrics may not generalize to other specialties. The rubrics were intended for quantitative and qualitative education research abstract submission, so it may not perform well for abstracts that include *both* quantitative and qualitative data or those focused on early work, innovations, instrument development, validity evidence, or

program evaluation. Finally, there are two limitations to the pilot testing. An a priori power calculation to determine sample size was not possible since the rubrics were novel. The ICCs of individual items on the scoring rubrics were variable and we chose not to eliminate items with low ICCs given the small sample size during piloting and a desire to create a tool comprehensive of key domains. Future studies of use of these tools incorporating larger samples may provide data for additional refinement. Faculty who piloted the rubrics were familiar with the constructs and rubrics, and it is not known how the rubrics would have performed with general abstract reviewers nor what training might be required. The success of separate rubrics may rely on the expertise of the reviewers in the methodology being assessed.

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

CONCLUSIONS

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

ACKNOWLEDGMENTS

The authors acknowledge that this project originated to meet an SAEM Education Committee Objective and thank all the committee members for their support of this work.

CONFLICTS OF INTEREST

The authors have no potential conflicts to disclose.

AUTHOR CONTRIBUTIONS

Jaime Jordan and Michael A. Gisondi conceived the study. Jaime Jordan, Michael A. Gisondi, Laura R. Hopson, Caroline Molins, and Suzanne K. Bentley contributed to the design of the study. Jaime Jordan, Laura R. Hopson, Caroline Molins, Suzanne K. Bentley, Nicole M. Deiorio, Sally A. Santen, Lalena M. Yarris, Wendy C. Coates, and Michael A. Gisondi contributed to data collection. Jaime Jordan analyzed the data. Jaime Jordan, Laura R. Hopson, Caroline Molins, Suzanne K. Bentley, Nicole M. Deiorio, Sally A. Santen, Lalena M. Yarris, Wendy C. Coates, and Michael A. Gisondi contributed to drafting of the manuscript and critical revision.

ORCID

Jaime Jordan  <https://orcid.org/0000-0002-6573-7041>

Laura R. Hopson  <https://orcid.org/0000-0002-1183-4751>

Suzanne K. Bentley  <https://orcid.org/0000-0003-0192-3133>

Sally A. Santen  <https://orcid.org/0000-0002-8327-8002>

Lalena M. Yarris  <https://orcid.org/0000-0003-1277-2852>

Wendy C. Coates  <https://orcid.org/0000-0002-3305-8802>

Michael A. Gisondi  <https://orcid.org/0000-0002-6800-3932>

REFERENCES

1. Padayachy A, Rodrigues G, Tahar A. Comment rédiger un abstract scientifique ? [How to write a scientific abstract]. *Rev Med Suisse*. 2019;15(664):1703-1706.
2. Timmer A, Sutherland LR, Hilsden RJ. Development and evaluation of a quality score for abstracts. *BMC Med Res Methodol*. 2003;3:2.
3. Mitchell NS, Stolzmann K, Benning LV, Wormwood JB, Linsky AM. Effect of a scoring rubric on the review of scientific meeting abstracts. *J Gen Intern Med*. 2020. doi: <https://doi.org/10.1007/s11606-020-05960-6>
4. Poolman RW, Keijser LC, De Waal Malefijt MC, et al. Reviewer agreement in scoring 419 abstracts for scientific orthopedics meetings. *Acta Orthop*. 2007;78(2):278-284.
5. van der Steen LP, Hage JJ, Kon M, Mazzola R. Reliability of a structured method of selecting abstracts for a plastic surgical scientific meeting. *Plast Reconstr Surg*. 2003;111(7):2215-2222.
6. Kuczarski TM, Raja AS, Pallin DJ. How do medical societies select science for conference presentation? How should they? *West J Emerg Med*. 2015;16(4):543-550.
7. Stephenson CR, Vaa BE, Wang AT, et al. Conference presentation to publication: a retrospective study evaluating quality of abstracts and journal articles in medical education research. *BMC Med Educ*. 2017;17(1):193.
8. Cook DA, Reed DA. Appraising the quality of medical education research methods: the medical education research study quality instrument and the Newcastle-Ottawa scale-education. *Acad Med*. 2015;90(8):1067-1076.
9. Sullivan GM. A primer on the validity of assessment instruments. *J Grad Med Educ*. 2011;3:119-120.
10. Cook DA, Beckman TJ. Current concepts in validity and reliability for psychometric instruments: theory and application. *Am J Med*. 2006;119:166.e10-199.e16.
11. Hasson F, Keeney S, McKenna HP. Research guidelines for the Delphi survey technique research guidelines for the Delphi survey. *J Adv Nurs*. 2000;32(4):1008-1015.
12. Thangaratnam SK, Redman CW. The Delphi technique. *Obstet Gynaecol*. 2005;7:120-125.
13. Pines JM, Alfaraj S, Batra S, et al. Factors important to top clinical performance in emergency medicine residency: results of an ideation survey and Delphi panel. *AEM Educ Train*. 2018;2(4):269-276.
14. Eubank BH, Mohtadi NG, Lafave MR, et al. Using the modified Delphi method to establish clinical consensus for the diagnosis and treatment of patients with rotator cuff pathology. *BMC Med Res Methodol*. 2016;16:56.
15. Messick S. Validity. In: Linn R, ed. *Educational Measurement*. 3rd ed. American Council on Education and Macmillan; 1989.
16. Join a Committee. Society for Academic Emergency Medicine website. c2021. Accessed July 5, 2021. <https://www.saem.org/about-saem/saem-membership/committees>
17. SurveyMonkey. c2021. Accessed April 7, 2021. <https://www.surveymonkey.com/>
18. Dubosh NM, Jordan J, Yarris LM, et al. Critical appraisal of emergency medicine educational research: the best publications of 2016. *AEM Educ Train*. 2018;3(1):58-73.
19. Lin M, Fisher J, Coates WC, et al. Critical appraisal of emergency medicine education research: the best publications of 2012. *Acad Emerg Med*. 2014;21:322-333.

20. Farrell SE, Kuhn GJ, Coates WC, et al. Critical appraisal of emergency medicine education research: the best publications of 2013. *Acad Emerg Med*. 2014;21(11):1274-1283.
21. Yarris LM, Juve AM, Coates WC, et al. Critical appraisal of emergency medicine education research: the best publications of 2014. *Acad Emerg Med*. 2015;22(11):1327-1336.
22. Heitz CR, Coates WC, Farrell SE, Fisher J, Juve AM, Yarris LM. Critical appraisal of emergency medicine education research: the best publications of 2015. *AEM Educ Train*. 2017;1(4):255-268.
23. Moralejo D, Ogunremi T, Dunn K. Critical appraisal toolkit (CAT) for assessing multiple types of evidence. *Can Commun Dis Rep*. 2017;43(9):176-181.
24. Reed DA, Beckman TJ, Wright SM, Levine RB, Kern DE, Cook DA. Predictive validity evidence for medical education research study quality instrument scores: quality of submissions to JGIM's medical education special issue. *J Gen Intern Med*. 2008; 23(7):903-907.
25. Johnson JL, Adkins D, Chauvin S. A review of the quality indicators of rigor in qualitative research. *Am J Pharm Educ*. 2020;84(1):7120.
26. Yang LJ, Chang KW, Chung KC. Methodologically rigorous clinical research. *Plast Reconstr Surg*. 2012;129(6):979e-988e.
27. Schneider NC, Coates WC, Yarris LM. Taking your qualitative research to the next level: a guide for the medical educator. *AEM Educ Train*. 2017;1(4):368-378.
28. O'Brien BC, Harris IB, Beckman TJ, Reed DA, Cook DA. Standards for reporting qualitative research: a synthesis of recommendations. *Acad Med*. 2014;89(9):1245-1251.

SUPPORTING INFORMATION

Additional supporting information may be found online in the Supporting Information section.

How to cite this article: Jordan J, Hopson LR, Molins C, et al. Leveling the field: Development of reliable scoring rubrics for quantitative and qualitative medical education research abstracts. *AEM Educ Train*. 2021;5:e10654. <https://doi.org/10.1002/aet2.10654>