CrossMark

# Assessment of surgical performance of laparoscopic benign hiatal surgery: a systematic review

Elif Bilgic[1] · Mohammed Al Mahroos[1] · Tara Landry[2] · Gerald M. Fried[1] · Melina C. Vassiliou[1] · Liane S. Feldman[1]

## Abstract

**Background** Operative skills correlate with patient outcomes, yet at the completion of training or after learning a new procedure, these skills are rarely formally evaluated. There is interest in the use of summative video assessment of laparoscopic benign foregut and hiatal surgery (LFS). If this is to be used to determine competency, it must meet the robust criteria established for high-stakes assessments. The purpose of this review is to identify tools that have been used to assess performance of LFS and evaluate the available validity evidence for each instrument.

**Methods** A systematic search was conducted up to July 2017. Eligible studies reported data on tools used to assess performance in the operating room during LFS. Two independent reviewers considered 1084 citations for eligibility. The characteristics and testing conditions of each assessment tool were recorded. Validity evidence was evaluated using five sources of validity (content, response process, internal structure, relationship to other variables, and consequences).

**Results** There were six separate tools identified. Two tools were generic to laparoscopy, and four were specific to LFS [two specific to Nissen fundoplication (NF), one heller myotomy (HM), and one paraesophageal hernia repair (PEH)]. Overall, only one assessment was supported by moderate evidence while the others had limited or unknown evidence. Validity evidence was based mainly on internal structure (all tools reporting reliability and item analysis) and content (two studies referencing previous papers for tool development in the context of clinical assessment, and four listing items without specifying the development procedures). There was little or no evidence supporting test response process (one study reporting rater training), relationship to other variables (two comparing scores in subjects with different clinical experience), and consequences (no studies). Two tools were identified to have evidence for video assessment, specific to NF.

**Conclusion** There is limited evidence supporting the validity of assessment tools for laparoscopic foregut surgery. This precludes their use for summative video-based assessment to verify competency. Further research is needed to develop an assessment tool designed for this purpose.
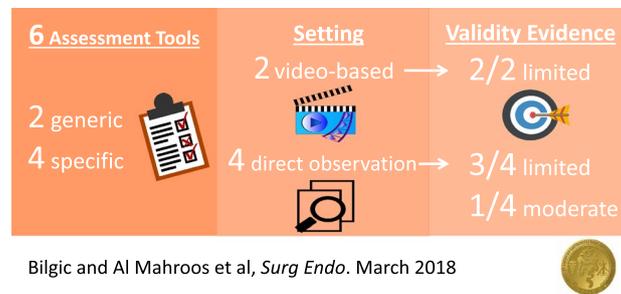
---

---

Mohammed Al Mahroos is Co-primary author.

---

Extended author information available on the last page of the article

## Graphical abstract



Assessing Surgical Performance – Benign Hiatel Surgery

6 Assessment Tools

Setting

Validity Evidence

2 video-based → 2/2 limited

2 generic
4 specific

4 direct observation → 3/4 limited
1/4 moderate

Bilgic and Al Mahroos et al, *Surg Endo*. March 2018

Surgical education is now focused on competency [1, 2]. Within this paradigm, assessment allows programs to track progress and make decisions about trainee skill level and readiness for independent practice [3]. Operative performance is a key measure of surgical practice that correlates with patient outcomes [4], yet no formal evaluation is required at the completion of training or after learning a new procedure. Traditionally, case numbers and end of rotation reports have been used for operative assessment. However, these measures are subjective and do not accurately assess trainee skill, which is why programs are shifting towards more objective assessments [5, 6]. To accurately measure skill and draw conclusions from the scores, assessment tools with appropriate supporting validity evidence need to be developed [3].

In a move towards objective assessments, the American Board of Surgery (ABS) has made the implementation of procedure-specific performance assessments mandatory, to keep track of trainee skill and ensure that all graduating residents have a certain skill level [7]. These assessments are completed through direct observation by the attending surgeon of each case (rater) using the Operative Performance Rating System (OPRS). In this setting, raters are not blinded to the training level or the operative experience of the trainees, which might cause them to have bias regarding the skill level of the person and affect the scores they assign [8]. An alternative to direct observation that would minimize rater bias is assessment using video-recorded performances, which would allow blinded raters to assess each trainee. In the United Kingdom National Training Program for laparoscopic colectomy (Lapco), a high-stakes, video-based assessment tool was developed in a transparent and rigorous fashion. After surgeons completed hands-on training with a supervising surgeon, they submitted videos scored by two independent experts using

the assessment tool, to be "signed-off" as competent to perform the procedure independently [9, 10].

The main focus of operative assessments have been on technical skills, which incorporate psychomotor skills such as bimanual coordination in addition to some cognitive skills such as knowledge of the procedure, anatomy, etc. However, in addition to technical skills and basic knowledge, advanced cognitive skills such as surgical planning, error prevention, and error recognition are central themes of expert performance in surgery, and should be included in operative assessments, to have a comprehensive understanding of surgeon competence [11].

Foregut and hiatal surgery are index operations for Advanced Gastrointestinal/Minimally Invasive Surgery Fellowships [12, 13]. For these complex procedures, patient outcomes have been shown to correlate with case experience of the surgeons [14, 15]. The Fellowship Council currently mandates program directors to use the Global Operative Assessment of Laparoscopic Skills (GOALS) quarterly to assess the skills of fellows performing laparoscopic fundoplication. GOALS is a tool that assesses generic laparoscopic skills [16]. However, when determining if a surgeon is competent to perform LFS, assessing generic skills might not be enough. An assessment tool specific to the procedure where specific skills required or important steps of the procedure are being deliberately and consistently assessed may be more relevant. This way, both the surgeon and the programs can determine aspects of the procedure that require further practice, which would not be possible to do with generic items alone. Furthermore, if surgical competence in hiatal surgery is to be evaluated, a high-stakes assessment, the tool for hiatal surgery should have robust validity evidence supporting its use and intended consequences.

The Society of Gastrointestinal and Endoscopic Surgeons (SAGES) is committed to surgical education and innovation to improve patient outcomes. In this context, there may be a

role for the use of video-based assessments for LFS, to determine competence and readiness for independent practice. As a high-stakes assessment, the tool would need to have robust validity evidence. Therefore, the purpose of this review is to identify assessment tools that have been used to assess performance of LFS and to evaluate the validity evidence for each instrument.

## Methods

### Search strategy

This is a systematic review of articles published between 1990 and July 2017 according to the preferred reporting items for systematic reviews and meta-analyses (PRISMA) guidelines. The search strategies were developed with the assistance of a health sciences librarian (T.L). A search was done in MEDLINE, Embase, Scopus and the Central registry of controlled trials via The Cochrane Library. Only English articles were used and no geographical limits were applied. Reference lists were searched to identify additional articles. A detailed search strategy is provided in the "Appendix".

### Study selection and extraction

Eligible studies reported data on assessment tools that have been used for LFS skill assessment in the operating room (OR) through video or direct observation. Studies using assessment tools for pediatric surgery, simulated settings, as well as reviews, meeting abstracts, editorials, and letters were excluded. Data were extracted by two independent reviewers (E.B and M.A). Differences in data from included articles were resolved through consensus adjudication.

### Validity evidence

Validity is defined as appropriate interpretation of assessment results. A validation study is the process of collecting evidence supporting the interpretations of assessment results [17]. Five sources of validity were evaluated according to the Standards established by the American Educational Research Association, the American Psychological Association, and the National Council on Measurement in Education [18]. These sources are: Content (relationship between the content of the tool and construct being measured), response process (scoring accuracy of the tool), internal structure (consistency of the score between raters, items etc), relationship to other variables (comparison of score to other measures), and consequences (implications of incorporating the tool into training programs) [17]. The validity evidence from various studies were summarized using Criteria for Rating Validity Evidence by Ghaderi et al. where each source is

scored from 0 to 3, with a total score ranging from 1 to 15 (1–5 limited, 6–10 moderate, 11–15 strong evidence) [19].

## Assessment tool characteristics

The extracted information included assessment tool type and the context the tools were used in. For performance assessments, three types of tool are used most commonly: global rating scales (assessing general abilities required to complete a task, using scales such as a Likert scale), checklists (assessing specific actions, using a dichotomous scale), and error rating scales (measuring the number of times an error has been made). Sometimes these tools are combined or used concomitantly. In addition, tools can be generic (assessing general skills required for laparoscopy) or specific (assessing specific skills required for LFS procedures).

In terms of the context, we identified whether the assessment was done using video-recorded performances or by direct observation. In addition, if it was done through direct observation, we recorded whether the raters were the attending surgeon, observer, or individuals assessing themselves. Finally, we evaluated whether the assessments were done for formative (feedback, low stakes), or summative (decision-making, high stakes) purposes.
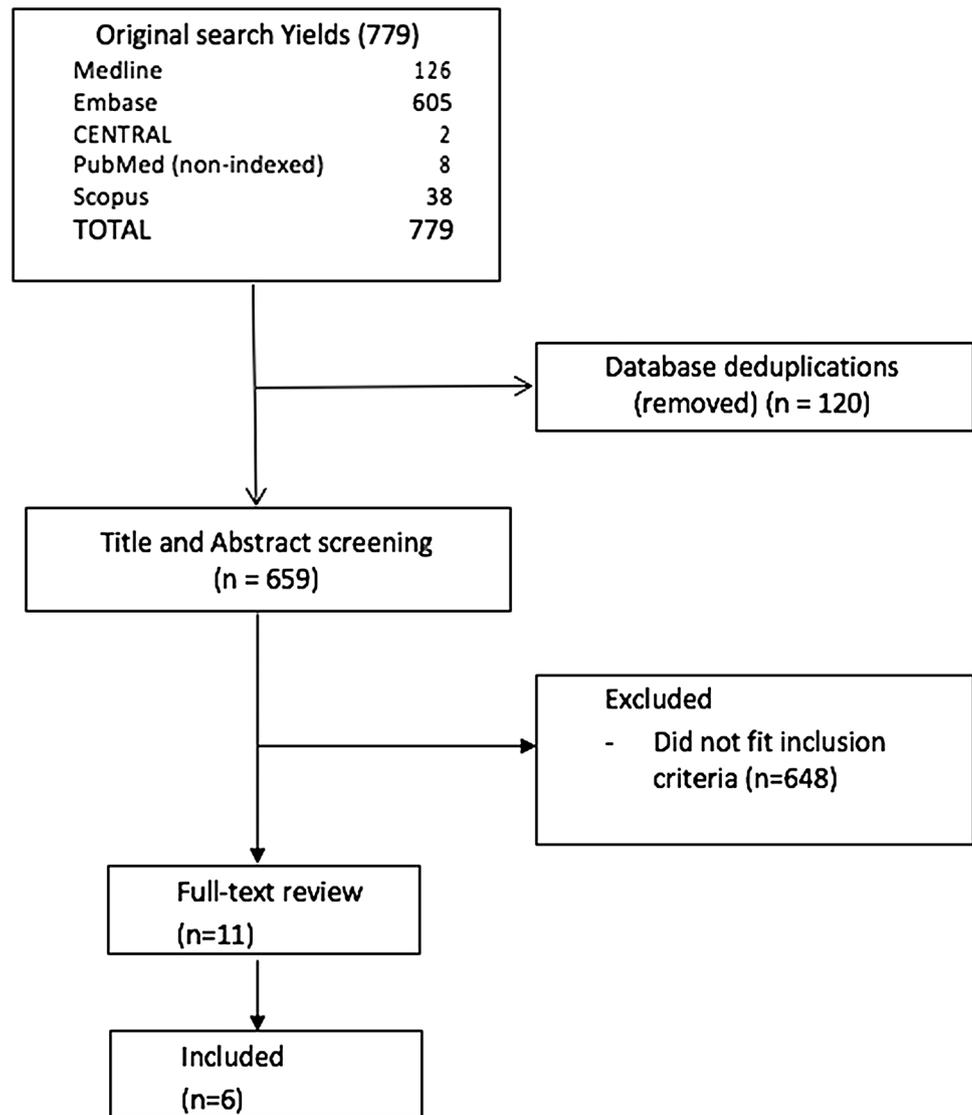
## Results

The primary search identified 659 studies. After title and abstract screening, 11 underwent full-text review, of which 6 met our inclusion criteria (Fig. 1) [20–25].

### Assessment tool characteristics

Six separate tools were identified. Two tools assessed generic laparoscopic performance during LFS, while four assessed a specific procedure [two for Nissen fundoplication, one for Heller myotomy and one for paraesophageal hernia repair (Table 1)]. Five tools were global rating scales, and one was a checklist. Four tools (two generic and two LFS-specific) were used for assessment through direct observation by the attending surgeon and self-assessment, whereas the remaining two tools (LFS-specific) were used for video assessments.

### Validity evidence

The majority of the evidence came from internal structure with all tools reporting rater reliability and/or item analysis (Table 2). Within this source of validity, the Global Operative Assessment of Laparoscopic Skills (GOALS) used three methods (item analysis, generalizability theory, and item response theory), while the others used one

**Fig. 1** Study selection flow chart



method. The second most frequently addressed source of validity was content where two studies referenced previous papers for tool development in the context of clinical assessment, and the other four reported the items without specifying how they developed the items. For the other three sources of validity, there was little or no evidence. For response process, two studies reported rater training. For relationship to other variables, two studies compared scores in subjects with different clinical experience and one of those two studies compared scores to operative data. Finally, no study reported on consequences of the assessment.

Using the validity rating criteria by Ghaderi et al., among the four tools used for direct observation, only GOALS had moderate evidence, while others had limited evidence. The two tools used for video assessment both had limited evidence (Table 3).

## Discussion

Through this review, we identified six assessment tools that have been used to assess laparoscopic benign hiatal and foregut surgery in the clinical setting: four through direct observation and two using video assessments. In addition, we evaluated the validity evidence that has been provided for the tools using a contemporary framework of validity. Most of the evidence came from internal structure in the form of rater reliability and item analysis. Overall, only one of the six assessment tools had moderate validity evidence (GOALS for direct observation; score between 6 and 10 in the validity rating criteria), while others had limited evidence. None of the assessments came close to meeting the established criteria for high-stakes examinations, and none can be recommended for determinations of competence or decisions about readiness for independent practice.

**Table 1** Characteristics of the assessment tools

| | Type of scale | Type of items | Number of items | Total score | Setting | Assessor | | | | Procedures |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | Reviewer | Observer | Attending | Self | |
| GOALS [20–23] | Global rating scale | Generic | 5 | 25 | Direct | | | + | + | HM, LNF, PEH |
| OSATS [23] | Global rating scale | Generic | 5 | 25 | Direct | | | + | + | HM |
| LNF-GRS [24] | Global rating scale | Specific to LNF | 7 | 21 | Video | + | | | | LNF |
| OCRS-HM [23] | Global rating scale | Specific to HM | 10 | 50 | Direct | | | + | + | HM |
| OCRS-PEH [23] | Global rating scale | Specific to PEH | 8 | 40 | Direct | | | + | + | PEH |
| LNF-CL [25] | Checklist | Specific to LNF | 65 | 65 | Video | + | | | | LNF |

*GOALS* Global Operative Assessment of Laparoscopic Skills, *OSATS* Objective Structured Assessment of Technical Skills, *LNF-GRS* Laparoscopic Nissen Fundoplication-Global Rating Scale, *OCRS* Objective Component Rating Scale, *LNF-CL* Laparoscopic Nissen Fundoplication-Checklist, *HM* Heller myotomy, *PEH* paraesophageal hernia, *LNF* laparoscopic Nissen fundoplication

GOALS is a generic assessment tool that evaluates general technical skills required to perform a laparoscopic procedure, such as bimanual dexterity and depth perception [16]. These are certainly important for surgeons; however, if someone obtains a high GOALS score while performing LFS, does that mean they are competent in performing LFS independently? A tool that is specific to LFS and that assesses the crucial steps of that procedure, would provide Program and Fellowship directors with a much better sense of how well their trainees can perform these operations. Such a tool could also be used to determine areas of weakness and for formative purposes. In addition, GOALS mostly focuses on technical skills while surgical expertise requires more complex and higher-order cognitive skills, that could potentially be incorporated into a LFS assessment to have a more comprehensive understanding of readiness for practice and how trainees integrate all of the critical subskills [11].

In addition, validity is context specific, meaning that where and how an assessment tool is used will determine the types of evidence that need to be provided [19]. For example, an assessment tool may have evidence for low-stakes assessment through direct observation by the attending surgeon, perhaps to guide the delivery of feedback to a trainee after a case. But this evidence may not be applicable when the same tool is used for high-stakes assessment through video-recorded performances. Therefore, validity is not about the tool itself, but about the evidence that we provide for various assessment conditions, and building an argument regarding the interpretation of the assessment scores. None of the identified studies clearly defined their purpose of assessment.

The majority of the studies investigated two of the five sources of validity, without much emphasis on the other three sources (response process, relationship to other variables, and consequences). Within the validity framework that we are using, the robustness of the validity evidence with regards to the tool's potential to assess LFS performance increases by addressing more sources of validity, meaning that each of the five sources adds something to the overall strength of the evidence [17].

An important limitation of the tools identified in the search is the lack of validity evidence for the content of the tools. Most of the studies listed the items without explaining how they generated the items. The content of the assessment tool is obviously very important since the items should reflect the range of competencies required when performing LFS. The process for content development was only reported for two tools using the previous literature. This is a problem since development of a robust assessment tool starts with content, which includes ensuring that the items of the tool are appropriate for measuring LFS skills in a high-stakes, video-based context. Recommendations for developing test content could include

**Table 2** Detailed analysis of the validity evidence of the assessment tools

| | GOALS [20–23] | OSATS [23] | LNF-GRS [24] | OCRS-HM [23] | OCRS-PEH [23] | LNF-CL [25] |
|---|---|---|---|---|---|---|
| *Content* | | | | | | |
| Expert judgment | | | | | | |
| Task analysis | | | | | | |
| Consensus method | | | | | | |
| *Response process* | | | | | | |
| Rater training | | | | | | + |
| Score interpretation and meaning | | | | | | |
| *Internal structure* | | | | | | |
| Rater reliability | | | + | | | + |
| Item analysis | + | + | | + | + | |
| Generalizability theory | + | | | | | |
| Item response theory | + | | | | | |
| Others | | | | | | |
| *Relationship to other variables* | | | | | | |
| Training level or case experience | + | | + | | | |
| Other performance assessment tool scores | | | | | | |
| Time | | | | | | |
| Operative data | + | | | | | |
| Motion analysis | | | | | | |
| Others | | | | | | |
| *Consequences* | | | | | | |
| Applications to residency program | | | | | | |
| Criterion-referenced score (benchmark or pass/fail) | | | | | | |

*GOALS* Global Operative Assessment of Laparoscopic Skills, *OSATS* Objective Structured Assessment of Technical Skills, *LNF-GRS* Laparoscopic Nissen Fundoplication-Global Rating Scale, *OCRS* Objective Component Rating Scale, *LNF-CL* Laparoscopic Nissen Fundoplication-Checklist

**Table 3** Validity evidence for the assessment tools

| | Content | Response process | Internal structure | Relations to other variables | Consequences | Total |
|---|---|---|---|---|---|---|
| GOALS [20–23] | 2 | 0 | 3 | 2 | 0 | 7 |
| OSATS [23] | 1 | 0 | 1 | 0 | 0 | 2 |
| LNF-GRS [24] | 1 | 0 | 1 | 1 | 0 | 3 |
| OCRS-HM [23] | 1 | 0 | 1 | 0 | 0 | 2 |
| OCRS-PEH [23] | 1 | 0 | 1 | 0 | 0 | 2 |
| LNF-CL [25] | 2 | 1 | 1 | 0 | 0 | 4 |

Each evidence category score is out of 3, with a total score of 15 [19]

*GOALS* Global Operative Assessment of Laparoscopic Skills, *OSATS* Objective Structured Assessment of Technical Skills, *LNF-GRS* Laparoscopic Nissen Fundoplication-Global Rating Scale, *OCRS* Objective Component Rating Scale, *LNF-CL* Laparoscopic Nissen Fundoplication-Checklist

creating an item bank through the literature review and expert opinion, followed by a Delphi method to get expert consensus on the items and the rating scale [26]. An example could be the Lapco assessment tool (high-stakes, video-based), where the development included evaluation of the selected literature to generate the items and scoring system, which was followed by expert consensus and finalization of the tool [10].

# Conclusion

There is limited evidence supporting the validity of existing intraoperative assessment tools for laparoscopic benign hiatal surgery. This precludes their use for high-stakes video-based assessment to verify competency. Further research is needed to develop an assessment tool for this purpose.

## Compliance with ethical standards

## Appendix

### Search strategy for medline

1. exp Esophageal Diseases/su [Surgery] (33036)
2. exp Esophagus/su [Surgery] (9601)
3. Hernia, Hiatal/su [Surgery] (1903)
4. exp Diverticulum, Esophageal/su [Surgery] (1098)
5. (esophag* or oesophag* or gastro?esophag* or para?esophag*).tw,kf. (167048)
6. Fundoplication/ (4057)
7. fundoplicat*.tw,kf. (5356)
8. (nissen adj3 (operat* or procedure*)).tw,kf. (362)
9. (heller adj3 myotom*).tw,kf. (687)
10. epiphrenic-diverticulectom*.tw,kf. (8)
11. mediastinal-dissect*.tw,kf. (283)
12. or/1-11 (175735)
13. exp Laparoscopy/ (84511)
14. laparoscop*.tw,kf. (107840)
15. 13 or 14 (120628)
16. 12 and 15 (7394)
17. Clinical Competence/ (80679)
18. (competen* or skill or skills).tw,kf. (235018)
19. 17 or 18 (287948)
20. Educational Measurement/ (33649)
21. "Task Performance and Analysis"/ (28752)
22. Checklist/(4544)
23. ((operativ* or intraoperativ* or perform*) adj5 (assess* or evaluat* or measur*)).tw,kf. (294777)
24. "Surveys and Questionnaires"/ (386261)
25. ed.fs. (254539)
26. (tool or tools or grade or grading or checklist* or check-list* or questionnaire* or form or rating or score* or scoring or scale or scaling).tw,kf. (2951343)
27. or/20-26 (3530579)
28. 19 and 27 (128912)
29. 16 and 28 (116)
30. ("24414454" or "20103075" or "18417086").ui. (3)
31. 29 and 30 (3)
32. 15 and 28 (2962)
33. ("26679826" or "25454951" or "24414454" or "20103075" or "18417086" or "27776756").ui. (6)
34. 32 and 33 (6)

## References

1. Beard JD, Marriott J, Purdie H, Crossley J (2011) Assessing the surgical skills of trainees in the operating theatre: a prospective observational study of the methodology. Health Technol Assess. https://doi.org/10.3310/hta15010
2. Hamstra SJ, Dubrowski A (2005) Effective training and assessment of surgical skills, and the correlates of performance. Surg Innov 12:71–77
3. Vassiliou MC, Feldman LS (2011) Objective assessment, selection, and certification in surgery. Surg Oncol 20:140–145
4. Birkmeyer JD, Finks JF, O'Reilly A, Oerline M, Carlin AM, Nunn AR, Dimick J, Banerjee M, Birkmeyer NJO (2013) Surgical skill and complication rates after bariatric surgery. N Engl J Med 369:1434–1442
5. Touchie C, Ten Cate O (2016) The promise, perils, problems and progress of competency-based medical education. Med Educ 50(1):93–100
6. Feldman LS et al (2004) Relationship between objective assessment of technical skills and subjective in-training evaluations in surgical residents. J Am Coll Surg 198:105–110
7. Larson JL, William RG, Ketchum J, Boehler ML, Dunnington GL (2005) Feasibility, reliability and validity of an operative performance rating system for evaluating surgery residents. Surgery 138:640–649
8. Streiner DL (1985) Global rating scales. In: Neufeld VR, Norman GR (eds) Assessing clinical competence. Springer, New York, pp 119–141
9. Mackenzie H, Ni M, Miskovic D, Motson RW, Gudgeon M, Khan Z, Longman R, Coleman MG, Hanna GB (2015) Clinical validity of consultant technical skills assessment in the English National Training Programme for laparoscopic colorectal surgery. Br J Surg 102:991–997
10. Miskovic D, Wyles SM, Carter F, Coleman MG, Hanna GB (2011) Development, validation and implementation of a monitoring tool for training in laparoscopic colorectal surgery in the English National Training Program. Surg Endosc 25:1136–1142
11. Madani A, Vassiliou MC, Watanabe Y, Al-Halabi B, Al-Rowais MS, Deckelbaum DL, Fried GM, Feldman LS (2017) What are the principles that guide behaviors in the operating room?: creating a framework to define and measure performance. Ann Surg 265:255–267
12. Medina M (2001) Formidable challenges to teaching advanced laparoscopic skills. J Soc Laparoendosc Surg 5:153–158
13. fellowshipcouncil.org
14. Broeders JAJL, Draaisma WA, Jong HGRd, Smout AJPM, Lanschot JJBv, Broeders IAMJ, Gooszen HG (2011) Impact of surgeon experience on 5-year outcome of laparoscopic Nissen fundoplication. Arch Surg 146:340–346
15. Birkmeyer JD, Stukel TA, Siewers AE, Goodney PP, Wennberg DE, Lucas FL (2003) Surgeon volume and operative mortality in the United States. N Engl J Med 349:2117–2127
16. Vassiliou MC, Feldman LS, Andrew CG, Bergman S, Leffondre K, Stanbridge D, Fried GM (2005) A global assessment tool for evaluation of intraoperative laparoscopic skills. Am J Surg 190:107–113
17. Downing SM (2003) Validity: on the meaningful interpretation of assessment data. Med Educ 37:830–837
18. Kogan JR, Holmboe ES, Hauer KE (2009) Tools for direct observation and assessment of clinical skills of medical trainees: a systematic review. JAMA 302:1316–1326
19. Ghaderi I, Manji F, Park YS, Juul D, Ott M, Harris I, Farrell TM (2015) Technical skills assessment toolbox: a review using the unitary framework of validity. Ann Surg 261:251–262

20. Hogle NJ, Liu Y, Ogden RT, Fowler DL (2014) Evaluation of surgical fellows' laparoscopic performance using Global Operative Assessment of Laparoscopic Skills (GOALS). Surg Endosc 28:1284–1290
21. Watanabe Y, Madani A, Ito YM, Bilgic E, McKendy KM, Feldman LS, Fried GM, Vassiliou MC (2016) Psychometric properties of the Global Operative Assessment of Laparoscopic Skills (GOALS) using item response theory. Am J Surg. https://doi.org/10.1016/j.amjsurg.2016.09.050
22. Bilgic E, Watanabe Y, McKendy K, Munshi A, Ito YM, Fried GM, Feldman LS, Vassiliou MC (2016) Reliable assessment of operative performance. Am J Surg 211:426–430
23. Ghaderi I, Auvergne L, Park YS, Farrell TM (2015) Quantitative and qualitative analysis of performance during advanced laparoscopic fellowship: a curriculum based on structured assessment and feedback. Am J Surg 209:71–78
24. Ahlberg G, Kruuna O, Leijonmarck CE, Ovaska J, Rosseland A, Sandbu R, Stromberg C, Arvidsson D (2005) Is the learning curve for laparoscopic fundoplication determined by the teacher or the pupil? Am J Surg 189:184–189
25. Peyre SE, Peyre CG, Hagen JA, Sullivan ME (2010) Reliability of a procedural checklist as a high-stakes measurement of advanced technical skill. Am J Surg 199:110–114
26. Awad M, Awad F, Carter F, Jervis B, Buzink S, Foster J, Jakimowicz J, Francis NK (2018) Consensus views on the optimum training curriculum for advanced minimally invasive surgery: a Delphi study. Int J Surg 53:137–142

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

## Affiliations

Elif Bilgic[1] · Mohammed Al Mahroos[1] · Tara Landry[2] · Gerald M. Fried[1] · Melina C. Vassiliou[1] · Liane S. Feldman[1]

✉ Liane S. Feldman
liane.feldman@mcgill.ca

1 Steinberg-Bernstein Centre for Minimally Invasive Surgery and Innovation, McGill University Health Centre, Montreal, QC, Canada

2 Montreal General Hospital Medical Library, McGill University Health Centre, 1650, Cedar Avenue, L9. 309, Montréal, QC H3G 1A4, Canada