




# Machine learning methods for automated technical skills assessment with instructional feedback in ultrasound-guided interventions

Matthew S. Holden<sup>1</sup>  · Sean Xia<sup>1</sup> · Hillary Lia<sup>1</sup> · Zsuzsanna Keri<sup>1</sup> · Colin Bell<sup>2</sup> · Lindsey Patterson<sup>3</sup> · Tamas Ungi<sup>1</sup> · Gabor Fichtinger<sup>1</sup>

Received: 31 October 2018 / Accepted: 9 April 2019  
© CARS 2019

## Abstract

**Objective** Currently, there is a worldwide shift toward competency-based medical education. This necessitates the use of automated skills assessment methods during self-guided interventions training. Making assessment methods that are transparent and configurable will allow assessment to be interpreted into instructional feedback. The purpose of this work is to develop and validate skills assessment methods in ultrasound-guided interventions that are transparent and configurable.

**Methods** We implemented a method based upon decision trees and a method based upon fuzzy inference systems for technical skills assessment. Subsequently, we validated these methods for their ability to predict scores of operators on a 25-point global rating scale in ultrasound-guided needle insertions and their ability to provide useful feedback for training.

**Results** Decision tree and fuzzy rule-based assessment performed comparably to state-of-the-art assessment methods. They produced median errors (on a 25-point scale) of 1.7 and 1.8 for in-plane insertions and 1.5 and 3.0 for out-of-plane insertions, respectively. In addition, these methods provided feedback that was useful for trainee learning. Decision tree assessment produced feedback with median usefulness 7 out of 7; fuzzy rule-based assessment produced feedback with median usefulness 6 out of 7.

**Conclusion** Transparent and configurable assessment methods are comparable to the state of the art and, in addition, can provide useful feedback. This demonstrates their value in self-guided interventions training curricula.

**Keywords** Ultrasound-guided needle insertion · Simulation-based training · Medical education · Objective skill assessment

## Introduction

Globally, skills training for medical interventions is transitioning from a time-based model to a competency-based model. Under the old time-based model, trainees would practice an intervention for a fixed amount of time, at which point they would be deemed competent and graduate, or they would be deemed incompetent and have to undertake significant remediation. Under the new competency-based model,

trainees practice until they achieve a predefined competency benchmark. This scheme allows each trainee to practice the precise amount of time they need to achieve competency. The drawback of this method is that trainees' competency needs to be continually monitored.

Expert-based methods for skills assessment include checklists, global rating scales, and entrustment scores. Checklists are application-specific rubrics which assess whether the operator performs each step in the intervention correctly [1]. Global rating scales (GRS) offer application-independent assessment of interventions across several different facets [2]. Entrustment scores assess to what degree a supervisor trusts the trainee to complete each face of the intervention [3]. While these methods provide reliable assessment, in particular when combined [4], they rely on experts. With increasing medical class sizes and demands on expert time, it is not feasible to implement expert-based assessment on a wide scale. Instead, skill assessment should be automated.

✉ Matthew S. Holden  
72mh@queensu.ca

<sup>1</sup> Laboratory for Percutaneous Surgery, School of Computing, Queen's University, Kingston, ON, Canada

<sup>2</sup> Department of Emergency Medicine, School of Medicine, Queen's University, Kingston, ON, Canada

<sup>3</sup> Department of Anesthesiology and Perioperative Medicine, School of Medicine, Queen's University, Kingston, ON, Canada

Automated skills assessment can be applied to many different interventions (e.g., laparoscopy, open surgery, needle insertion, etc.), and can use data from many sources (e.g., instrument tracking, video, surgeon status, patient monitors) [5], [6]. Perhaps, the most common method for automated skills assessment is metrics-based assessment. Under this paradigm, clinical experts specify what aspects of the intervention are relevant to operator skill. Subsequently, these can be implemented into a set of performance metrics: quantities that are understandable to trainees and clinicians and can be readily computed from measurable data. From these performance metrics, an overall skill level can be derived using pattern recognition or machine learning approaches.

Metrics-based overall skills assessment was initially addressed as an optimization problem, where each metric is treated as a cost and the most skillful operator is the one who best minimized the weighted sum of costs [7, 8]. Since, pattern recognition approaches have been used to achieve improved reliability in assessment. Chmarra et al. showed that linear discriminant analysis reliably distinguishes novices from intermediates and from experts in laparoscopic training tasks [9]. Likewise, Allen et al. showed that support vector machines outperform cost-based approaches for skill classification in laparoscopic training tasks [10]. Oropesa et al. also demonstrated that support vector machines outperform linear discriminant analysis and adaptive neuro-fuzzy inference systems for laparoscopic training tasks [11]. Ahmidi et al. use support vector machines for skill classification for several difference types of performance metrics in septoplasty [12]. Fard et al. contrasted support vector machines with k-nearest neighbors and logistic regression for identifying novices and experts in robotic suturing tasks on real patients [13]. Kramer et al. have suggested learning vector quantization and self-organizing maps for assessment in simulated vascular surgery [14]. Neural network-based approaches have seen some success [15]. Fuzzy pattern recognition approaches have also gained some traction, including rule-based methods [16, 17] and adaptive fuzzy inference systems [18].

In consultation with clinical experts, we suggest two criteria which metrics-based skills assessment methods should meet in order to be clinically useful: transparency and configurability. A machine learning approach is considered transparent if both the models are easy to interpret and the principal of the method is easily understood [19, 20]. A machine learning approach is considered configurable if it has parameters which can be configured to improve performance based on domain knowledge from a domain expert (Chiticariu et al. consider this a component of transparency [20]). In interventional skills assessment, transparency allows both the supervisor and trainee to understand why the trainee received a particular score and to interpret their results into actionable strategies to improve

performance. Configurability allows the expert to adjust the assessment to their particular training scenario or to emphasize particular skills.

Of course, there are other methods for interventional skills assessment that are not based on performance metrics. In particular, temporal modeling [21], process monitoring [22], and end-to-end deep learning approaches [23] have shown some promise for skills assessment. Unfortunately, these methods do not provide adequate transparency to allow trainees and supervisors to interpret results into actionable feedback to improve performance. Crowdsourcing can also provide accurate skills assessment and is effectively automated [24], but cannot provide immediate feedback.

The objective of this work is to develop and validate methods for overall skills assessment in percutaneous interventions. The methods should be transparent, configurable, and conducive to self-guided training. Subsequently, we evaluated (1) the accuracy of the proposed methods compared to state-of-the-art computer-assisted assessment and (2) the usefulness of the feedback provided by our proposed methods.

A preliminary version of this work has been reported [25].

## Methods

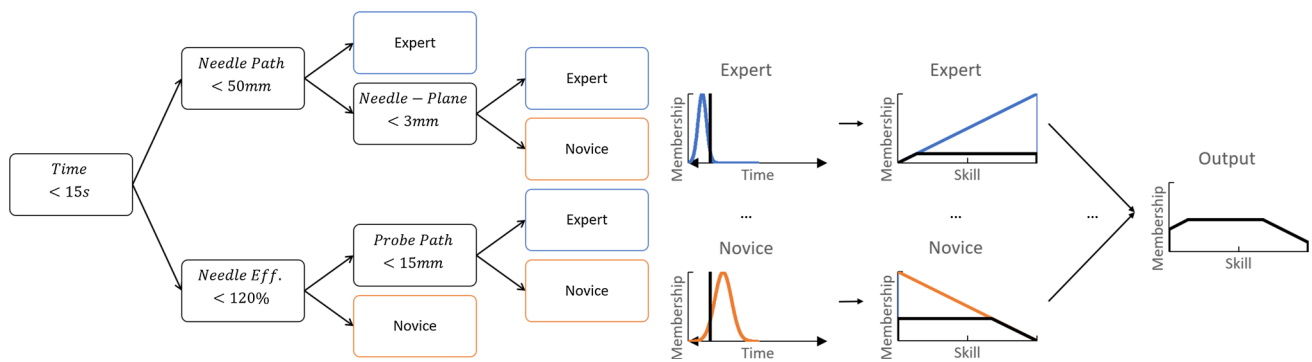
### Skills assessment algorithms

We aim to implement skills assessment algorithm which are transparent and configurable. Transparency and configurability are inherently subjective and fuzzy. In a review of common machine learning techniques, Kotsiantis identified three techniques as highly transparent: decision trees, Naïve Bayes, and rule-based learners [19]. Naïve Bayes is more applicable for classification and performs poorly for regression tasks [26]. This leaves decision trees and rule-based learners as transparent and configurable machine learning approaches for interventional skills assessment. Each skills assessment method takes a set of performance metrics as input (i.e., feature vector) and computes an overall skill level as output.

### Decision tree assessment

For transparent and configurable assessment using decision trees, we use importance-aided decision trees [27] (Fig. 1). This method is intended to incorporate domain knowledge into decision trees, especially in lower levels of the decision tree when training data is limited.

In decision tree learning, a split is made based on the attribute and value which optimizes some measure of purity of each branch. In our case, for regression, we choose the within-branch variance as the attribute selection score. Following Al Iqbal et al., for an attribute  $x$ , we create a new



**Fig. 1** Illustration of importance-aided decision tree assessment (left) and fuzzy rule-based assessment (right) in ultrasound-guided needle insertion assessment using performance metrics

attribute selection score  $S$  based on a linear combination of the within-branch variance score  $S_v$  and the attribute's weight  $W$  [27].

$$S(x) = (1 - \rho)S_v(x) + (\rho)W(x)$$

We select the attribute and split point which optimizes this new attribute selection score  $S$ . The coefficient in the linear combination  $\rho$  grows inversely with the proportion of remaining training samples in the branch [27]. The splitting is stopped once the within-branch variance decreases beyond a certain threshold. At this point, all training instances in the branch will effectively have the same skill level. We observe that in the case of equal attribute weights, this functions in the same way as a classical decision tree.

As identified by Kotsiantis, this method is transparent [19]. The user is presented with the traversal of the decision tree and the splitting criteria. As actionable feedback, we can identify the metrics associated with splits in the traversal where the branch center changed the most. The feedback “well done” is provided when all splits in the traversal result in positive change in the branch center. This method is configurable in that the weights associated with each attribute can be adjusted. As demonstrated by Al Iqbal et al., incorporating this domain knowledge into the decision tree can improve the accuracy of assessment [27].

### Fuzzy rule-based assessment

For rule-based assessment that is transparent and configurable, we use a set of fuzzy inference rules (Fig. 1). In particular, we choose to use rules of the form: IF <metric> is <skill level> THEN operator is <skill level>. For example, IF elapsed time is expert THEN operator is expert. Such a rule is defined for each metric and skill level pair.

In practice, this requires us to define a membership function for each skill class and a membership function for each metric for each skill class. We define the skill class mem-

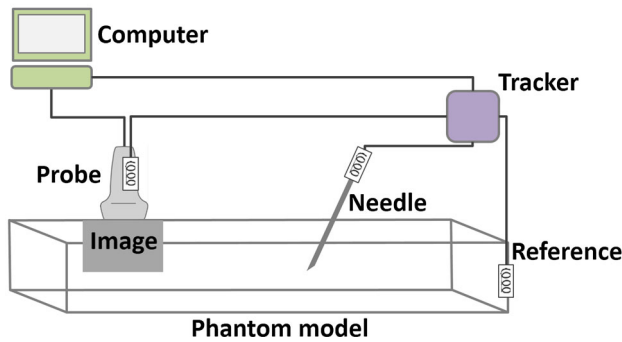
bership functions as symmetrical triangular functions on the range [0, 1], overlapping such that membership over all classes sums to one [16, 17]. The metric membership function for each skill class is computed empirically from the training data by Gaussian kernel density estimation, using the Silverman's rule-of-thumb to estimate the bandwidth [28]. Importantly, each training instance may have membership in multiple skill classes and contribute with different weights to multiple metric membership functions.

We use clipping based on the membership in the input function to compute the output membership function for each rule. The set of fuzzy rules is combined and defuzzified by computing the mean of the maximum of the output membership functions.

This rule-based assessment method is transparent [19]. The user is presented with the rules that were applied and their strengths. As actionable feedback, we can identify the metrics for which the net influence of all associated fuzzy rules is the strongest. The feedback “well done” is provided when for each metric, the net influence of all rules associated with that metric is positive. This method can be configured by allowing the weights associated with each rule to be adjusted or fuzzy rules to be added or removed. In particular, experts can add more sophisticated rules based on their domain knowledge for improved accuracy.

### Validation of assessment accuracy

We validated our assessment methods on both in-plane and out-of-plane needle insertions on a vascular access phantom (CAE Healthcare), following the setup used in Xia et al. [29]. We recorded 19 trainees and five experts performing in-plane insertions and 19 trainees and five experts performing out-of-plane insertions. Trainees were recorded at two points during a training curriculum (Fig. 2). Experts were each recorded once. Operators used a Telemed MicrUs linear ultrasound probe (Telemed Medical Systems). We recorded videos of participants' hands and tracked the needle and



**Fig. 2** Photograph of a trainee participant performing an ultrasound-guided in-plane needle insertion (top) and schematic diagram of setup (bottom). The electromagnetic pose trackers used are attached to the base of the needle, base of the ultrasound probe, and exterior of the phantom

ultrasound probe. Tools were tracked with the Ascension trakStar (Northern Digital Inc.), and data was recorded using the PLUS Toolkit ([www.plustoolkit.org](http://www.plustoolkit.org)) [30] and Perk Tutor ([www.perktutor.org](http://www.perktutor.org)) [31]. We computed eight metrics for in-plane insertions and seven metrics for out-of-plane insertions (Table 1) [29]. These metrics were designed based on consultation with clinical experts, and they are intended to cover all relevant aspects of the ultrasound-guided needle insertion tasks.

As ground-truth assessment, we did not use participants' level of training. Instead, we recruited three clinical experts to assess participants' performance via anonymized hand motion videos using a previously validated global rating scale [32, 33]. The mean overall expert assessment provides a ground-truth skill level out of 25.

To determine the weight associated with each metric, we interviewed the same three clinical experts who provided ratings on the global rating scale, and we asked them to rate the importance of each metric for skills assessment on a seven-point Likert scale. We linearly scaled these ratings onto the interval [0, 1].

Subsequently, we validated the performance of our proposed assessment methods using leave-one-user-out cross-validation. We computed difference in the output of the proposed assessment methods with the mean expert rating. We then compared these results with the results achieved

**Table 1** Description of performance metrics for in-plane and out-of-plane insertions

<i>In-plane metrics</i>	
Elapsed time (s)	Total time from the start of the insertion to the end of the insertion
Needle path length (mm)	Total distance travelled by the tip of the needle
Probe path length (mm)	Total distance travelled by the foot of the ultrasound probe
Needle path efficiency (%)	Ratio of the needle's path length to the distance between the needle's start and end points
Average needle to ultrasound plane distance (mm)	Average orthogonal distance between the needle tip and the ultrasound plane
Maximum needle to ultrasound plane distance (mm)	Maximum orthogonal distance between the needle tip and the ultrasound plane
Average needle to ultrasound plane angle (°)	Average angle between the needle and the ultrasound plane
Maximum needle to ultrasound plane angle (°)	Maximum angle between the needle and the ultrasound plane
<i>Out-of-plane metrics</i>	
Elapsed time (s)	Total time from the start of the insertion to the end of the insertion
Needle path length (mm)	Total distance travelled by the tip of the needle
Probe path length (mm)	Total distance travelled by the foot of the ultrasound probe
Needle path efficiency (%)	Ratio of the needle's path length to the distance between the needle's start and end points
Maximum distance needle is past ultrasound plane (mm)	Maximum orthogonal distance the needle tip travels past the ultrasound plane
Total time needle is past ultrasound plane (s)	Total time spent with the needle tip past the ultrasound plane
Average rotation from needle to ultrasound plane normal (°)	Average angle between the needle and the plane orthogonal to the ultrasound marked-unmarked vector

from several standard methods: (1) zero-rule regression (i.e., always guessing the mean scores), (2) linear regression, an empirically optimal version of the sum of z-scores method [8], (3) support vector machine regression, which has been shown to achieve state-of-the-art results in several assessment tasks [10–13], (4) nearest neighbor regression with sequential forward feature selection, which achieves highly accurate assessment in suturing and knot tying [34, 35], and (5) random regression forests, a generalization on decision tree regression. To compare the methods, we used a Fried-

**Table 2** Plain-language feedback associated with each performance metric for in-plane and out-of-plane insertions

<i>In-Plane Metrics</i>	
Elapsed Time (s)	F1. Keep practicing with proper technique to improve your time efficiency.
Needle path length (mm)	F2. Look at the depth of your target, and try to estimate the correct angle of needle insertion.
Probe path length (mm)	F3. Get a longitudinal ultrasound image of the middle of the vessel and stabilize the probe using your hand/finger against the gel surface.
Needle path efficiency (%)	F4. Try to focus on a smooth, straight needle path while inserting the needle as close to the ultrasound plane as possible.
Average needle to ultrasound plane distance (mm)	F5. Start with the needle in the middle of the ultrasound probe and try to keep it aligned with the ultrasound plane during needle insertion.
Maximum needle to ultrasound plane distance (mm)	
Average needle to ultrasound plane angle (°)	F6. Do not change the angle between the needle and the ultrasound plane during needle insertion. This will make sure that there is perfect alignment.
Maximum needle to ultrasound plane angle (°)	
	F7. Well done.
<i>Out-of-Plane Metrics</i>	
Elapsed Time (s)	F1. Keep practicing with proper technique to improve your time efficiency.
Needle path length (mm)	F2. Look at the depth of your target, and try to estimate the correct angle of needle insertion.
Probe path length (mm)	F3. Do not move the probe when advancing the needle. Advance the probe very slightly when the needle appears in the ultrasound image.
Needle path efficiency (%)	F4. Insert the needle in a straight, smooth path.
Maximum distance needle is past ultrasound plane (mm)	F5. Keep the ultrasound plane slightly ahead of the needle. If you see the needle tip on the screen, move the ultrasound slightly ahead until the needle disappears and then continue needle insertion until the needle appears again.
Total time needle is past ultrasound plane (s)	
Average rotation from needle to ultrasound plane normal (°)	F6. Start with the target in the middle of the ultrasound screen, with the needle in the middle of the probe at 90° to the probe and 45° to the gel surface. Do not change this angle during the needle insertion.
	F7. Well done.

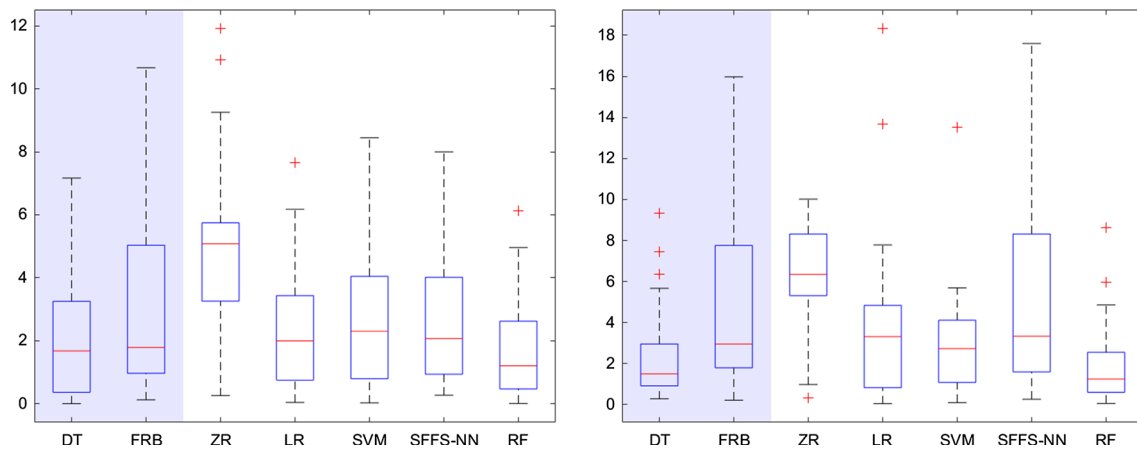
man test with pairwise Dunn's post hoc tests with Bonferroni correction ( $\alpha = 0.05$ ). To determine whether our assessment methods are comparable to these other methods, we performed non-inferiority signed-rank tests ( $\alpha = 0.05$ ) with the pooled standard deviation in expert ratings as the non-inferiority margin.

To determine the added value of expert-configured assessment, we used Bonferroni-corrected signed-rank tests ( $\alpha = 0.05$ ) to compare unconfigured assessment with expert-configured assessment. We tested: (1) assessing the mean expert-assigned score using the mean expert configuration and (2) assessing each expert-assigned score using each expert's respective configuration.

### Validation of feedback accuracy

To assess the quality of feedback provided by our methods, we mapped each metric to a plain-language description (Table 2). This was done in consultation with our clinical experts to ensure the vocabulary covers all possible feedbacks an expert might provide to trainees during in a typical training scenario. Subsequently, we asked one expert to review each trainee's post-training video (as this was identified by experts as the most useful stage for feedback). At the end of each video, we showed the expert all the different feedbacks and asked them to rate the usefulness of each one on a seven-point Likert scale (1 = strongly disagree that feedback was useful; 4 = neutral; 7 = strongly agree that feedback was useful).





**Fig. 3** Error in assessment for the decision tree (DT), fuzzy rule-based (FRB), zero-rule (ZR), linear regression (LR), support vector machine (SVM), nearest neighbor with sequential forward feature selection

(SFFS-NN), and random forest (RF) assessment methods for in-plane insertions (left) and out-of-plane insertions (right). Data from 24 users over 43 trials

We compared the usefulness of the feedback provided by the proposed methods with the usefulness of the  $k$ th most useful feedback by signed-rank test ( $\alpha = 0.05$ ), for all  $k$ . We report the smallest  $k$  for which the predicted feedback is significantly more useful than the  $k$ th most useful feedback provided by the expert. This provides evidence of the usefulness of the proposed methods relative to expert feedback, without being skewed by the fact that experts found the majority of feedbacks to be useful. We also report confusion matrices for the truly most useful feedback compared to the predicted feedback.

## Results

### Assessment accuracy

For ground-truth skill, the average measures intraclass correlation coefficient was 0.90 for the in-plane insertions and 0.93 for the out-of-plane insertions, indicating good reliability. For decision tree assessment and fuzzy rule-based assessment, respectively, the median errors were 1.7 and 1.8 for in-plane insertions and 1.5 and 3.0 for out-of-plane insertions (Fig. 3). Post hoc tests revealed decision tree assessment significantly outperformed all methods except support vector machine assessment (Table 3).

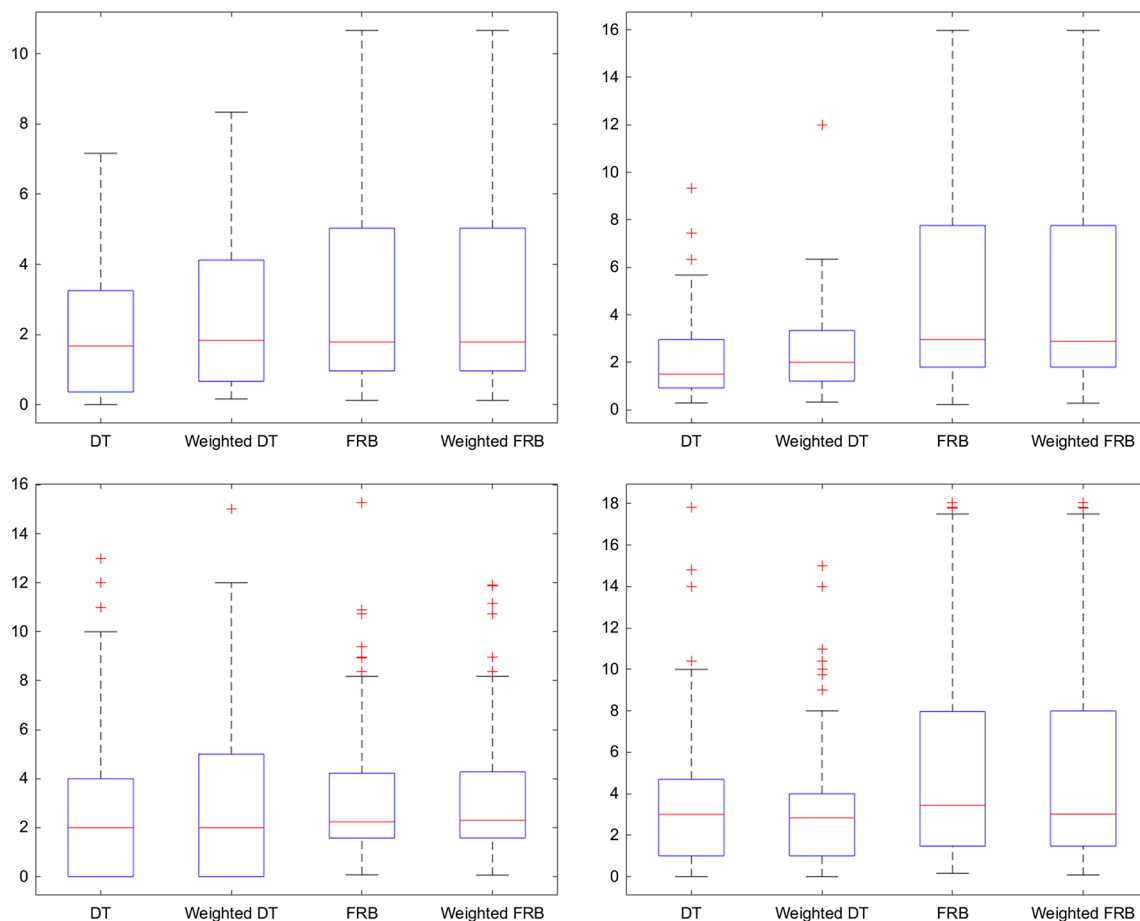
Decision tree assessment was non-inferior to all other assessment methods for both in-plane and out-of-plane insertions. Fuzzy rule-based assessment was non-inferior to all other assessment methods for in-plane insertions. For out-of-plane insertions, however, significance was not achieved. In fact, for out-of-plane insertions, fuzzy rule-based assessment was non-inferior to only zero-rule and nearest neighbor with sequential forward feature selection.

**Table 3** Results of post hoc testing for differences in decision tree (DT), fuzzy rule-based (FRB), zero-rule (ZR), linear regression (LR), support vector machine (SVM), nearest neighbor with sequential forward feature selection (SFFS-NN), and random forest (RF) assessment

Assessment method	Mean rank	Significance
<i>In-plane</i>		
Decision tree	3.09	< ZR
Fuzzy rule-based	4.40	> RF
Zero-rule	5.74	>DT, LR, SVM, SFFS-NN, RF
Linear regression	3.63	<ZR
Support vector machine	4.09	<ZR
SFFS-nearest neighbor	4.21	<ZR
Random forest	2.84	<FRB, ZR
<i>Out-of-plane</i>		
Decision tree	3.19	<FRB, ZR, SFFS-NN
Fuzzy rule-based	4.63	>DT, SVM, RF
Zero-rule	5.81	>DT, LR, SVM, RF
Linear regression	3.88	<ZR
Support vector machine	3.21	<FRB, ZR
SFFS-nearest neighbor	4.60	>DT, RF
Random forest	2.67	<FRB, ZR, SFFS-NN

Mean ranks indicate the mean rank of accuracy for the method when compared to the other methods. Significance indicates which methods were significantly different, and whether the method was more accurate (<) or less accurate (>)

Reliability in the mean expert-defined weights was poor. The average measures intraclass correlation coefficient was 0.49 for in-plane insertions and 0.32 for out-of-plane insertions. When we used the expert-defined weights in the configurable assessment methods, the change in accuracy was insignificant (Fig. 4).



**Fig. 4** Error in assessment for decision tree (DT) and fuzzy rule-based (FRB) assessment methods with or without expert-defined weights for in-plane insertions (left) and out-of-plane insertions (right). Top row shows results from predicting the mean expert-assigned score using

the mean expert configuration; bottom row shows results from predicting each individual expert-assigned score with each expert's respective configuration. Data from 24 users over 43 trials

## Feedback accuracy

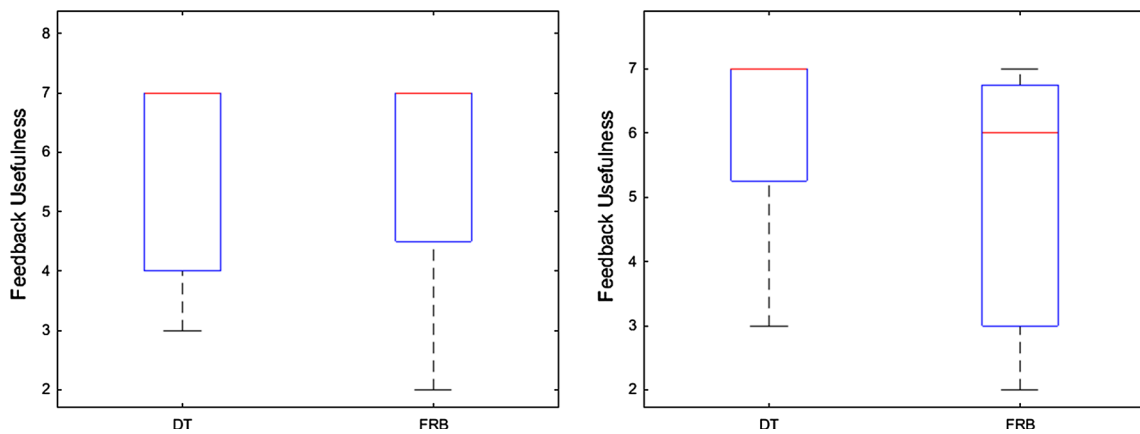
The usefulness of the feedback was rated a median 7 out of 7 and a mean 5.8 out of 7 on a Likert scale for decision tree assessment. Likewise, the usefulness of the feedback was rated a median 6 out of 7 and a mean 5.3 out of 7 on a Likert scale for fuzzy rule-based assessment (Fig. 5). Decision tree assessment produced useful feedback 74% of the time, and fuzzy rule-based assessment produced useful feedback 63% of the time (feedback rated 5, 6, or 7 out of 7 on a Likert scale). Furthermore, both methods produced significantly better than neutral feedback. Confusion matrices illustrate the most commonly misclassified feedback (Fig. 6).

Compared to expert feedback, we found that for in-plane insertions, both decision tree assessment and fuzzy rule-based assessment produced significantly better feedback than the fifth best expert feedback. For out-of-plane insertions, decision tree assessment produced significantly better feedback than the third best expert feedback, and fuzzy rule-based assessment produced significantly better feedback than the

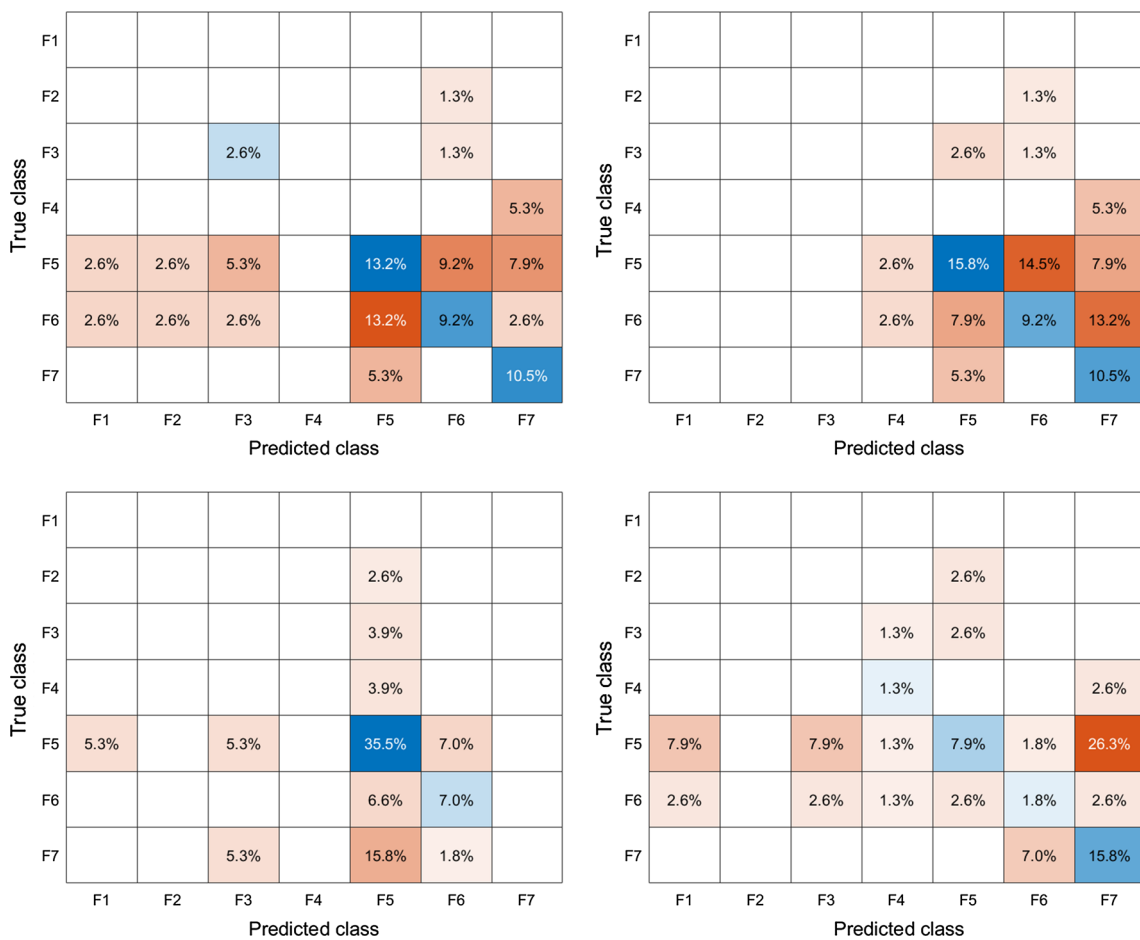
fifth best expert feedback. In all cases, the feedback produced by the proposed methods was better than the median expert feedback, but this was significant only for decision tree assessment in out-of-plane insertions.

## Discussion

The results show that transparent and configurable assessment methods (1) perform comparably to state-of-the-art methods and (2) provide useful feedback for training. In particular, decision tree assessment performed most accurately and provided the most useful feedback for our dataset. We did not observe a significant change in assessment accuracy when experts configured the proposed methods based on their domain knowledge [36, 37]. We believe the lack of significant improvement in the presented results may be due to our definition of ground-truth skill as a sum of global rating scale scores, without considering the importance of each aspect. Furthermore, the experts found all the metrics that



**Fig. 5** Usefulness of feedback produced by the decision tree (DT) and fuzzy rule-based (FRB) assessment methods for in-plane insertions (left) and out-of-plane insertions (right). Red line indicates median. Data from 24 users over 43 trials



**Fig. 6** Confusion matrices illustrating errors in predicted feedback for in-plane (top) and out-of-plane (bottom) ultrasound-guided needle insertions using importance-aided decision tree (left) and fuzzy rule-based (right) assessment. Ties are distributed across all tied labels.

Labels correspond to feedback vocabulary. Blue shading indicates correct predictions; red shading indicates incorrect predictions. Intensity of shading indicates higher concentration. Data from 19 users over 19 trials

were defined to be useful on average (rated as 5 or higher out of 7), and thus, the expert configurations are not substantially different from the default configuration.

We have identified our methods as transparent and identified methods such as support vector machines as opaque. While we have followed the work of Kotsiantis in identifying



machine learning techniques as transparent [19], these classifications are inherently fuzzy. Although there is ongoing work in introspection in deep learning [38] allowing users to gain some understanding of how the deep model reached the result, it is unclear how well such methods will be accepted into practice [39].

One of the limitations of our feedback is the finite nature of our feedback vocabulary. While our feedback vocabulary was generated in consultation with experts to cover every aspect of the intervention, it does not allow feedback to be tailored to a particular trainee, as preceptors would do in practice. We observed that experts rated the top feedback as a median 7 out of 7, indicating they agreed that the feedback from the vocabulary was indeed useful.

Another challenge of this work was determining the weights for each feature. We used a Likert scale to capture experts' opinion about the importance of each aspect of the intervention, and linearly scaled these responses to weights. But we observed that there was poor consistency between experts. This indicates that each expert may value different aspects of ultrasound-guided insertions. Our methods would allow the assessment to be tailored to each expert individually.

Although there are five experts and nineteen trainees for each of the in-plane and out-of-plane insertions, we observe that ground-truth global rating scores cluster toward the higher end of the scale. This creates a problem of unbalanced regression and may affect the reported results.

We observed that in 8% and 29% of cases, "well done" was incorrectly predicted as the most useful feedback, for decision tree assessment and fuzzy rule-based assessment, respectively. But the feedback "well done" may not be the most instructive for trainees. The proposed methods could be adapted to provide this feedback less frequently. For decision tree assessment, this could be achieved by requiring all splits in the traversal to have a change in branch center above a certain threshold. Analogously, for fuzzy rule-based assessment, this could be achieved by requiring the net influence of all rules associated with each metric to be above a certain threshold. The threshold value could be tuned to optimize a sensitivity and specificity criterion.

We have shown that the proposed methods work effectively for skills assessment and feedback in both in-plane and out-of-plane ultrasound-guided needle insertions. Our setup has shown evidence for face and content validity [29]. Because the overwhelming majority of ultrasound-guided interventions use one of these approaches, we suggest the results will apply to most ultrasound-guided interventions. Recent work has shown that it takes approximately 85 practice attempts to reach proficiency in ultrasound-guided needle insertions [40]. Our experts believe that the feedback provided by our system will be most applicable after ten prac-

tice attempts when the trainee has fully understood the basics of the intervention.

Our results are consistent with other work demonstrating the utility of metrics-based assessment of interventional skills [7, 8]. In the context of ultrasound-guided needle insertions, we have shown that transparent and configurable methods are comparable to state-of-the-art methods for assessment but, in addition, can provide useful feedback.

In the future, we suggest further study into how the proposed methods perform in specific ultrasound-guided interventions (e.g., biopsy, epidural, central line) and how they may be extended to other types of interventions. It has been previously shown that generic performance metrics may not be equally applicable to all interventions [41] and application-specific metrics provide added value over generic metrics [42]. Thus, in order to extend these methods to other interventions, it is necessary to develop application-specific performance metrics and a feedback vocabulary in consultation with expert clinicians. These are the only places where we have infused application-specific knowledge into the proposed methods.

We also suggest a future longitudinal study examining the effect of providing the proposed computer-generated feedback on trainee learning. Such a study could better identify the added value of the proposed feedback methods over self-guided training without feedback. Prior work has shown that feedback through 3D visualization can improve ultrasound-guided interventions learning [43], but has not evaluated the impact of targeted feedback.

We make the proposed methods available to the community through Perk Tutor ([www.perktutor.org](http://www.perktutor.org)) [31].

## Conclusion

We have demonstrated that transparent and configurable skills assessment methods are comparably accurate to state-of-the-art methods. In contrast to state-of-the-art methods, however, transparent and configurable methods were shown to provide useful feedback for training. Importance-aided decision tree assessment provided the most accurate assessment with feedback.

Thus, transparent and configurable assessment methods can be adopted into practice to provide feedback without compromising accuracy. We have also demonstrated that they can be customized by experts to suit the particular application or emphasize particular skills.

We envision that these methods could be employed in an ultrasound-guided interventions training curriculum. They would monitor trainee learning curves and provide automated instructions during self-directing learning. This would serve to supplement supervision and assessment from expert preceptors.

**Acknowledgements** Matthew S. Holden is supported by the Link Foundation Fellowship in Modeling, Simulation, and Training. Gabor Fichtinger is supported by a Canada Research Chair in Computer-Integrated Surgery.

## Compliance with ethical standards

**Conflict of interest** All authors declare that they have no conflict of interest.

**Ethical approval** All procedures in this study involving human participants were performed in accordance with the ethical standards of the institution and were approved by the research ethics board at Queen's University. This study does not contain any procedures involving animals.

**Informed consent** All participation was voluntary, and written informed consent was obtained from all participants.

## References

- Harden RM, Stevenson M, Downie WW, Wilson GM (1975) Assessment of clinical competence using objective structured examination. *Br Med J* 1(5955):447–451
- Winckel CP, Reznick RK, Frcsc M, Cohen R (1994) Reliability and construct validity of a structured technical skills assessment form. *Am J Surg* 167:423–427
- Gofton WT, Dudek NL, Wood TJ, Balaa F, Hamstra SJ (2012) The Ottawa surgical competency operating room evaluation (O-SCORE): a tool to assess surgical competence. *Acad Med* 87(10):1401–1407
- Martin JA, Regehr G, Reznick R, Macrae H, Murnaghan J, Hutchison C, Brown M (1997) Objective structured assessment of technical skill (OSATS) for surgical residents. *Br J Surg* 84(2):273–278
- Reiley CE, Lin HC, Yuh DD, Hager GD (2011) Review of methods for objective surgical skill evaluation. *Surg Endosc* 25(2):356–366
- Vedula SS, Ishii M, Hager GD (2017) Objective assessment of surgical technical skill and competency in the operating room. *Annu Rev Biomed Eng* 19(1):301–325
- Fraser SA, Klassen DR, Feldman LS, Ghitulescu GA, Stanbridge D, Fried GM (2003) Evaluating laparoscopic skills, setting the pass/fail score for the MISTELS system. *Surg Endosc Other Interv Tech* 17(6):964–967
- Stylopoulos N, Cotin S, Maithel SKK, Ottensmeyer M, Jackson PGG, Bardsley RSS, Neumann PFF, Rattner DWW, Dawson SLL, Ottensmeyer M, Jackson PGG, Bardsley RSS, Neumann PFF, Rattner DWW, Dawson SLL (2004) Computer-enhanced laparoscopic training system (CELTS): bridging the gap. *Surg Endosc* 18(5):782–789
- Chmarra MK, Klein S, de Winter JCF, Jansen F-WW, Dankelman J (2010) Objective classification of residents based on their psychomotor laparoscopic skills. *Surg Endosc Other Interv Tech* 24(5):1031–1039
- Allen B, Nistor V, Dutson E, Carman G, Lewis C, Faloutsos P (2010) Support vector machines improve the accuracy of evaluation for the performance of laparoscopic training tasks. *Surg Endosc* 24(1):170–178
- Oropesa I, Sánchez-González P, Chmarra MK, Lamata P, Pérez-Rodríguez R, Jansen FW, Dankelman J, Gómez EJ (2014) Supervised classification of psychomotor competence in minimally invasive surgery based on instruments motion analysis. *Surg Endosc Other Interv Tech* 28(2):657–670
- Ahmidi N, Poddar P, Jones JD, Vedula SS, Ishii L, Hager GD, Ishii M (2015) Automated objective surgical skill assessment in the operating room from unstructured tool motion in septoplasty. *Int J Comput Assist Radiol Surg* 10(6):981–991
- Fard MJ, Ameri S, Darin Ellis R, Chinnam RB, Pandya AK, Klein MD (2017) Automated robot-assisted surgical skill evaluation: Predictive analytics approach. *Int J Med Robot Comput Assist Surg* 14(1):e1850
- Kramer BD, Losey DP, O'Malley MK, O'Malley MK (2016) SOM and LVQ classification of endovascular surgeons using motion-based metrics. In: Merényi E, Mendenhall MJ, O'Driscoll P (eds) *Advances in self-organizing maps and learning vector quantization: proceedings of the 11th international workshop WSOM 2016*, Houston, Texas, USA, January 6–8, 2016, vol. 428. Springer, Cham, pp 227–237
- Uemura M, Tomikawa M, Miao T, Souzaki R, Ieiri S, Akahoshi T, Lefor AK, Hashizume M (2018) Feasibility of an AI-based measure of the hand motions of expert and novice surgeons. *Comput Math Methods Med*. <https://doi.org/10.1155/2018/9873273>
- Hajshirmohammadi I, Payandeh S (2007) Fuzzy set theory for performance evaluation in a surgical simulator. *Presence Teleoper Virtual Environ* 16(6):603–622
- Riojas M, Feng C, Hamilton A, Rozenblit J (2011) Knowledge elicitation for performance assessment in a computerized surgical training system. *Appl Soft Comput J* 11(4):3697–3708
- Huang J, Payandeh S, Doris P, Hajshirmohammadi I (2005) Fuzzy classification: towards evaluating performance on a surgical simulator. *Stud Health Technol Inform* 111:194–200
- Kotsiantis SB (2007) Supervised machine learning: a review of classification techniques. *Informatica* 31:249–268
- Chiticariu L, Li Y, Reiss F (2015) Transparent machine learning for information extraction: state-of-the-art and the future. In: *Conference on empirical methods in natural language processing*, pp 4–6
- Rosen J, Brown JD, Chang L, Barreca M, Sinanan M, Hannaford B (2002) The BlueDRAGON—a system for measuring the kinematics and dynamics of minimally invasive surgical tools in-vivo. *IEEE Int Conf Robot Autom* 2:1876–1881
- Forestier G, Lalys F, Riffaud L, Trelhu B, Jannin P (2012) Classification of surgical processes using dynamic time warping. *J Biomed Inform* 45(2):255–264
- Doughty H, Damen D, Mayol-Cuevas W (2018) Who's better? Who's best? pairwise deep ranking for skill determination. In: *2018 IEEE/CVF conference on computer vision and pattern recognition*, pp 6057–6066. <https://doi.org/10.1109/CVPR.2018.00634>
- Kowalewski TM, Comstock B, Sweet R, Schaffhausen C, Menhadji A, Averch T, Box G, Brand T, Ferrandino M, Kaouk J, Knudsen B, Landman J, Lee B, Schwartz BF, McDougall E, Lendvay TS (2016) Crowd-sourced assessment of technical skills for validation of basic laparoscopic urologic skills tasks. *J Urol* 195(6):1859–1865
- Holden MS, Lia H, Xia S, Keri Z, Ungi T, Fichtinger G (2018) Configurable overall skill assessment in ultrasound-guided needle insertion. In: *16th annual imaging network Ontario symposium (ImNO)*
- Frank E, Trigg L, Holmes G, Witten IH (2000) Technical note: Naive Bayes for regression. *Mach Learn* 41(1):5–25
- Al Iqbal MR, Rahman S, Nabii SI, Chowdhury IUA (2012) Knowledge based decision tree construction with feature importance domain knowledge. In: *2012 7th international conference on electrical and computer engineering*, pp 659–662
- Silverman BW (1986) *Density estimation for statistics and data analysis*, no. 1951
- Xia S, Keri Z, Holden MS, Hisey R, Lia H, Ungi T, Mitchell CH, Fichtinger G (2018) A learning curve analysis of ultrasound-guided in-plane and out-of-plane vascular access training with Perk Tutor.

- In: *Medical imaging 2018: image-guided procedures, robotic interventions, and modeling*, vol 10576, p 66
30. Lasso A, Heffter T, Rankin A, Pinter C, Ungi T, Fichtinger G (2014) PLUS: open-source toolkit for ultrasound-guided intervention systems. *IEEE Trans Biomed Eng* 61(10):2527–2537
  31. Ungi T, Sargent D, Moulton E, Lasso A, Pinter C, McGraw RC, Fichtinger G (2012) Perk tutor: an open-source training platform for ultrasound-guided needle insertions. *IEEE Trans Biomed Eng* 59(12):3475–3481
  32. Domuracki K, Wong A, Olivieri L, Grierson LEM (2015) The impacts of observing flawed and flawless demonstrations on clinical skill learning. *Med Educ* 49(2):186–192
  33. Ma IWY, Zalunardo N, Pachev G, Beran T, Brown M, Hatala R, McLaughlin K (2012) Comparing the use of global rating scale with checklists for the assessment of central venous catheterization skills using simulation. *Adv Health Sci Educ* 17(4):457–470
  34. Zia A, Sharma Y, Bettadapura V, Sarin EL, Clements MA, Essa I (2015) Automated assessment of surgical skills using frequency analysis. In: *Medical image computing and computer-assisted interventions—MICCAI 2015, Pt I*, vol 9349, pp 430–438
  35. Zia A, Sharma Y, Bettadapura V, Sarin EL, Essa I (2018) Video and accelerometer-based motion analysis for automated surgical skills assessment. *Int J Comput Assist Radiol Surg* 13(3):443–455
  36. Stumpf S, Rajaram V, Li L, Burnett M, Dietterich T, Sullivan E, Drummond R, Herlocker J (2007) Toward harnessing user feedback for machine learning. In: *Proceedings of the 12th international conference on Intelligent user interfaces—IUI'07*, p 82
  37. Talbot J, Lee B, Kapoor A, Tan DS (2009) EnsembleMatrix: interactive visualization to support machine learning with multiple classifiers. *Learning*. <https://doi.org/10.1145/1518701.1518895>
  38. Hendricks LA, Akata Z, Rohrbach M, Donahue J, Schiele B, Darrell T (2016) Generating visual explanations. In: *Lecture notes in computer science (including subseries lecture notes in artificial intelligence and lecture notes in bioinformatics)*, vol 9908. LNCS, pp 3–19
  39. Muir BM (1987) Trust between humans and machines, and the design of decision aids. *Int J Man Mach Stud* 27(5–6):527–539
  40. McGraw R, Chaplin T, McKaigney C, Rang L, Jaeger M, Redfearn D, Davison C, Ungi T, Holden M, Yeo C, Keri Z, Fichtinger G (2016) Development and evaluation of a simulation-based curriculum for ultrasound-guided central venous catheterization. In: *CJEM*, pp 1–9
  41. Datta V, Mackay S, Mandalia M, Darzi A (2001) The use of electromagnetic motion tracking analysis to objectively measure open surgical skill in the laboratory-based model. *J Am Coll Surg* 193(5):479–485
  42. Holden MS, Keri Z, Ungi T, Fichtinger G (2017) Overall proficiency assessment in point-of-care ultrasound interventions: the stopwatch is not enough. In: Cardoso MJ, Arbel T, Tavares JMRS, Aylward S, Li S, Boctor E, Fichtinger G, Cleary K, Freeman B, Kohli L, Shipley Kane D, Oetgen M, Pujol S (eds) *Imaging for patient-customized simulations and systems for point-of-care ultrasound: international workshops, BIVPCS 2017 and POCUS 2017, held in conjunction with MICCAI 2017, Québec City, QC, Canada, September 14, 2017*. Springer International Publishing, Cham, pp 146–153
  43. Lia H, Keri Z, Holden MS, Harish V, Mitchell CH, Ungi T, Fichtinger G (2017) Training with Perk Tutor improves ultrasound-guided in-plane needle insertion skill. In: *SPIE medical imaging, 2017*, p 101350T

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