



Published in final edited form as:

*Ann Surg.* 2019 March ; 269(3): 574–581. doi:10.1097/SLA.0000000000002478.

## Modeling surgical technical skill using expert assessment for automated computer rating

David P. Azari, MS<sup>(1)</sup>, Lane L. Frasier, MD<sup>(2)</sup>, Sudha R. Pavuluri Quamme, MD MS<sup>(1)</sup>, Caprice C. Greenberg, MD MPH<sup>(1),(2)</sup>, Carla Pugh, MD PhD<sup>(2)</sup>, Jacob A. Greenberg, MD EDM<sup>(2)</sup>, Robert G. Radwin, PhD<sup>(1),(3),\*</sup>

<sup>1</sup>Department of Industrial and Systems Engineering, University of Wisconsin-Madison, Madison, WI

<sup>2</sup>Wisconsin Surgical Outcomes Research (WiSOR) Program, Department of Surgery, University of Wisconsin-Madison, Madison, WI

<sup>3</sup>Department of Biomedical Engineering, University of Wisconsin-Madison, Madison, WI

### Structured Abstract

**Objective:** Computer vision was used to predict expert performance ratings from surgeon hand motions for tying and suturing tasks.

**Summary Background Data:** Existing methods, including the objective structured assessment of technical skills (OSATS) have proven reliable, but do not readily discriminate at the task level. Computer vision may be used for evaluating distinct task performance throughout an operation.

**Methods:** Open surgeries were videoed and surgeon hands were tracked without using sensors or markers. An expert panel of three attending surgeons rated tying and suturing video clips on continuous scales from 0 to 10 along three task measures adapted from the broader OSATS: motion economy, fluidity of motion, and tissue handling. Empirical models were developed to predict the expert consensus ratings based on the hand kinematic data records.

**Results:** The predicted v. panel ratings for suturing had slopes from 0.73 to 1, and intercepts from 0.36 to 1.54 (Average  $R^2 = 0.81$ ). Predicted v. panel ratings for tying had slopes from 0.39 to 0.88, and intercepts from 0.79 to 4.36 (Average  $R^2 = 0.57$ ). The mean square error among predicted and expert ratings were consistently less than the mean squared difference among individual expert ratings and the eventual consensus ratings.

**Conclusions:** The computer algorithm consistently predicted the panel ratings of individual tasks, and were more objective and reliable than individual assessment by surgical experts.

### Mini Abstract:

Expert panel consensus ratings of performance were modeled by tracking surgeon hand movements. Experts rated tying and suturing tasks on continuous scales, from 0 to 10 for motion economy, fluidity of motion, and tissue handling. The kinematic models objectively and reliably predicted panel ratings of individual tasks.

\*Corresponding Author: Robert G. Radwin, PhD, University of Wisconsin-Madison, 1550 Engineering Drive, Madison, WI 53706, radwin@engr.wisc.edu, Telephone: 608-263-6596.

## Keywords

computer vision; marker-less video tracking; surgical task analysis; objective structured assessment of technical skills; OSATS

---

## INTRODUCTION

Computer vision technologies utilize digital cameras and computer algorithms for tracking motions, classifying information, detecting features, and recognizing patterns in digital images and videos.<sup>1</sup> It has impacted a diverse field of applications, ranging from industrial robotics, intelligent and autonomous vehicles, security surveillance, manufacturing inspection, and human-computer interaction. Applications include face detection, expression and emotion recognition, video surveillance and monitoring, and motion tracking and analysis,<sup>2,3</sup> to name a few. This proof-of-concept paper compares models using computer motion tracking of surgeon hand movements against performance ratings made by experts.

A growing body of literature suggests that the technical skill of the surgeon can affect patient outcomes.<sup>4-6</sup> Technical skill contributes to three-quarters of adverse events in the operating room,<sup>7,8</sup> but common methods of measuring performance (i.e. in-training reports, mentor observations and evaluation) have been criticized for lacking reliability and being generally subjective<sup>9</sup> and are under increasing pressure to document proficiency.<sup>10</sup> The Objective Structured Assessment of Technical Skills (OSATS) is based on observing performance in real-time and rating candidates along a series of Likert-based scales in conjunction with a procedure-specific checklist.<sup>11</sup> Studies examining OSATS have demonstrated strong validity evidence in accordance with Kane's framework,<sup>12,13</sup> especially in providing formative feedback during training.<sup>14</sup> Despite this proven track record, correctly implementing OSATS is resource intensive and time consuming,<sup>15</sup> prompting exploration of more efficient assessment techniques.<sup>16,17</sup>

Motion capture and tracking of surgeon hand movements have the potential to fill these gaps.<sup>18-22</sup> The Imperial College Surgical Assessment Device (ICSAD) is an excellent example. Studies utilizing ICSAD have shown it is possible to identify consistent differences in trainees and experts by observing the position of small sensors placed on the hands,<sup>23,24</sup> and that these movements correlate strongly with OSATS global rating assessments.<sup>25</sup> Unfortunately, this self-contained system, is limited to benchtop models and consequently not easily applied in the operating room.<sup>26</sup> Assessment in the operating room requires a non-invasive and scalable means of observing surgical motion without relying on embedded sensor technology.

Video motion capture of the surgeon's hands in the surgical field, facilitated by increasing availability of cameras in the OR, offers an expedient alternative to integrated sensor-software systems. We have developed marker-less video processing methods using conventional single camera digital video to reliably track the motion trajectory of a selected region of interest over successive video frames without the need for sensors or markers.<sup>27-29</sup> We have also demonstrated how tracked hand motion can quantify kinematic properties of movements and exertions during specific tasks.<sup>27,30,31</sup> Previous studies by our group have

used this technology to isolate kinematic differences in surgical hand motion,<sup>20</sup> and identify statistically significant differences between surgeon roles (attending vs resident), tasks (tying vs suturing) and tissue types during open surgery.<sup>32</sup>

This study compares expert ratings of surgical skills (our current gold standard) to kinematic measures of surgeon hand motions to evaluate the potential use of video to automatically measure technical skill. The objective is to establish empirical models of expert-rated *in vivo* performance during operations by extracting kinematic features of tracked hand movements. We created task-specific rating scales and tracked hand motion records to predict subject matter expert ratings of surgical performance for a series of suturing and tying tasks. We hypothesize that the kinematic features of a surgeon's hand motion measured using marker-less tracking can be used to accurately model subjective performance ratings made by a panel of experts.

## METHODS

### Participants

This study utilized in-light OR video cameras (Figure 1) to capture hand motions of 9 surgeons (6 attendings and 3 residents) during 16 surgical cases. Cases included colorectal, complex upper gastrointestinal, hepatobiliary, surgical oncology, transplant, vascular, thoracic, and cardiac operations. Our institutional review board granted ethical approval for recording and analysis of video data for these operations. Written informed consent was obtained from participant surgeons in advance. Cases were recorded over nine months. Due to the positioning and limited field of view of the in-light OR camera, no patient details or protected health information were captured. Audio was not recorded.

### Video Selection

Recording of approved cases was initiated remotely after the operation began and stored on a secure hospital computer enabled with a video encoder (AXIS video encoder, Axis Communications, Lund, Sweden). We reviewed the videos using Multimedia Video Task Analysis (MVTA)<sup>TM</sup> software (Wisconsin Alumni Research Foundation, Madison, WI) developed by Yen and Radwin.<sup>33</sup> MVTA allows real-time review, labeling, and exporting of video events. A member of the research team familiar with operative technical tasks (LLF) screened and categorized the videos in MVTA for segments of tying or suturing tasks where the hands were clearly visible for at least five seconds. Suturing tasks were categorized as bowel anastomosis (SBA), complex anastomosis (SCA) such as hepaticojejunostomy, or suturing on the body wall (SBW) during closure or stoma formation, for example. Each tying task was categorized as follows: closing and external tying on the body wall (TBW), intra-abdominal superficial tying (TSF) or intra-abdominal deep tying (TDP). A full taxonomy and further information on the surgical procedures performed are described in Frasier et al.<sup>32</sup>

### Rating Scales

Subjective visual-analog rating scales<sup>34</sup> were created for evaluating surgical performance during short clips (5-30 seconds) for motion economy, fluidity of motion, and tissue

handling (Figure 2). We developed the scales utilizing the OSATS motion scales (i.e. respect for tissue, time and motion, and instrument handling) as assessment blueprints. The scales were framed to evaluate performance during a short clip rather than for an individual over a whole procedure (thereby omitting procedure-specific checklists), and defined to comprise the entire range (0-10) of possible behaviors observed over that segment.

Fluidity of motion is a measure of hesitancy, pauses, or changes in direction and “resets,” which may be a component of Moulton’s “slowing down,”<sup>35</sup> or contribute to time spent idle.<sup>36</sup> Tissue handling quantifies the appropriateness of the surgeon’s force and tension when manipulating the tissue,<sup>37,38</sup> and varies based on the tissue’s friability and fragility.<sup>18</sup> Motion economy is defined as efficiency of movement, or conservation of energy in any trajectory. Such behavior is consistently documented as a mark of expert psychomotor behavior, and increasingly studied as a measure of surgical skill.<sup>25,39</sup>

A consensus panel of three expert surgeons (CCG, JAG, CMP) viewed each clip in random order and independently rated hand motion across the scales. Each surgeon announced their ratings to the group. Discrepancies were discussed until consensus was achieved for each rating. If absolute consensus was unable to be achieved, the clip was scored according to the majority evaluation.

### **Motion Tracking**

Custom video tracking software,<sup>27</sup> developed in one of the authors’ (RGR) lab, was used to trace a region of interest (ROI) on a visible portion of the surgeon’s hands (generally the index finger or thumb) across successive video frames. We previously used this marker-less tracking approach to examine differences in hand motion between attending and resident surgeons,<sup>32</sup> to discriminate between dominant and non-dominant hand motion during reduction mammoplasty<sup>20</sup> and to evaluate hand-motion patterns during simulated clinical breast exams.<sup>31</sup> The x-y pixel location of the ROI for each frame (every  $\frac{1}{30}$  of a second) was recorded. An analyst selected the size and position of the ROI to track the surgeon’s hand for each video clip, and due to occasional occlusions of the hands and changes in lighting, supervised tracking of the ROI, making manual corrections as necessary.

### **Calibration**

In-frame visible measurements of the hands were used to calibrate each video clip from pixels to millimeters. We previously used hand dimensions to provide acceptable estimates of hand speed.<sup>30</sup> Observed proximal interphalangeal joint breadth was scaled to the population means of males (23.0 mm) or females (19.9 mm) depending on the gender of the surgeon. The proximal interphalangeal joint breadth was selected due to its small coefficient of variation for males (0.071) and females (0.064) as determined by the US Army.<sup>40</sup> The video was calibrated based on average hand measurements across three different frames, and recalibrated for any change in the position of the in-light camera or the surgical field.

### **Variable Selection**

The tracked record of the ROI location allowed quantification of instantaneous displacement, speed and acceleration of the surgeon’s hand, and several additional measures

including jerk<sup>41</sup> and spatiotemporal curvature.<sup>42</sup> Jerk is the third derivative of position with respect to time and generally represents how smooth a motion is, while the spatiotemporal curvature function is a measure of direction change based on multiple derivatives of the position signal and is used to indicate the number of discrete movements.

A second order Butterworth filter within empirically observed upper limits of hand frequency for mono-hand tasks<sup>43</sup> (pass band = 0.005 to 1.000 Hz) was applied to smooth the acceleration and curvature signals and a Fast Fourier Transform (FFT) function was applied to each signal set to isolate any consistently cyclic and repeated motion patterns, a growing avenue of research in surgical dexterity.<sup>44</sup> Moving averages were calculated for the original and smoothed signal types (speed, acceleration, curvature), as well as relative densities to see how much area is consistently and repeatedly traversed by the hand (this can essentially be interpreted as a proxy for path length, a common variable output by platforms such as the ICSAD.<sup>24</sup>) These predictor variables are distributed across the variable families shown in Table 1.

### Modeling Process

We developed a set of linear regression models to test whether the kinematic features could predict the expert ratings across each of the motion scales. The predictor variables were examined for collinear relationships, and subsets of these selected for regression analysis. In selecting variables to initiate the modeling process, we hypothesized that the average speed and consistent locations in the speed signal would give strong predictions of motion economy ratings, while the peak arrival rate in the acceleration signal would predict fluidity ratings and jerk or pauses would predict tissue handling ratings. In other words, motion economy might correspond to how fast the surgeons are moving in a consistent area, while fluidity might relate to how many changes in speed there are, and tissue handling would be sensitive to any sudden, repeated changes in direction. The models were optimized using the Akaike information criterion to balance the prediction against the number of variables.<sup>45</sup> Every consensus rating scale (fluidity of motion, motion economy, tissue handling) was first modeled for each distinct category of surgery (i.e. SBW, SBA, SCA, TBW, TSF, TDP). We then examined the combined predictions from suturing tasks and tying tasks for each rating scale (Figure 3).

### Validation

In order to assess the internal validity of the prediction models, we compared the sum of squared errors (SSE), a statistical measure of variance, to the leave-one-out predicted residual sum of squares (PRESS) statistic – a common approach in validating models where test data is not available.<sup>46,47</sup> The SSE measure assessed the overall fit, whereas the PRESS statistic penalizes the model if it depends on outliers. The closer these values are, the better the model matches the current data without relying on outliers, and the better prediction it provides. We calculated the root-mean SSE and root-mean PRESS statistics to normalize the comparison between models, and refer to this value as the model's "error" for ease of reference (see Table 2). Ideal prediction models would have several properties. The predicted and ground truth ratings should form a straight line from (0,0) to (10,10), such that any unitary increase in expert appraisal would be met with a similar predicted increase. We

arbitrarily defined acceptable models to have a slope between 0.5 and 1.5, an intercept within  $\pm 2$  of zero, and an  $R^2 > 0.75$ . While other thresholds could be defined, these bounds were chosen to ensure (1) a limited difference between the ground truth and the prediction, and (2) a prediction consistent throughout the domain of each scale.

## RESULTS

### Video Data

In total, 103 video clips (mean time = 11.72 seconds) were recorded of six attending and three resident surgeons performing suturing and tying tasks throughout 16 varied operations. These included SBA, SCA and SBW suturing clips (n = 44), and TBW, TSF and TDP tying clips (n = 59). The marker-less hand tracking frame-by-frame position data provided more than 1,500 kinematic variables, spread across the variable types in Table 1.

### Task Expert Rating Scales

Each clip was observed by three expert surgeons and rated along 0-10 analog scales for motion economy, fluidity of motion, and tissue handling. The expert raters achieved consensus after a maximum of three iterations for all the video clips. Observed ratings for suturing tasks ranged from 2-9 (mean = 5.91, sd = 1.62). Tying task ratings ranged from 3-9 (mean = 7.12, sd = 1.10) and were particularly skewed towards the upper end of the scale.

### Prediction Models of Expert Ratings

The linear regression slopes and intercepts for each of the rating scales and task types are seen in Table 2. The models for fluidity of motion for suturing (Figure 3a; slope = 0.86, intercept = 0.80,  $R^2 = 0.86$ ) and motion economy for suturing (Figure 3b; slope = 0.89, intercept = 0.64,  $R^2 = 0.88$ ) had the best predictions, and moderately better than tissue handling of suturing (Figure 3c; slope = 0.76, intercept = 1.1,  $R^2 = 0.69$ ).

The predicted vs. consensus ratings of motion economy for tying tasks (Figure 3e; slope = 0.65, intercept = 2.51,  $R^2 = 0.64$ ) had better predictions than tissue handling (Figure 3f; slope = 0.53, intercept = 3.33,  $R^2 = 0.52$ ) and fluidity of motion (Figure 3d; slope = 0.54, intercept = 3.28,  $R^2 = 0.54$ ).

### Prediction Model Validity

Linear prediction models were validated by comparing their difference, or “error”, between root-mean predicted residual (PRESS) and root-mean sum of square errors (SSE). Motion economy had the least internal error (0.17, 0.07) during bowel anastomosis and superficial tying, respectively, as well as the least average error (0.27) for suturing ratings overall. Fluidity of motion models similarly performed well for suturing and tying along the body wall with errors of 0.21 and 0.11. Tissue handling had the highest error for both suturing and tying tasks (0.52, 0.26). Low errors suggest strong positive relationships between hand kinematics and expert ratings – a necessary component of Kane and Messick’s modern approach to validity.<sup>13</sup>

## DISCUSSION

This study demonstrates that computer vision capturing kinematic data from marker-less motion tracking of video records offers an objective and scalable approach for measuring surgical skill, and establishes evidence consistent with Kane's validity framework<sup>13</sup>. We created subjective rating scales for fluidity of motion, motion economy and tissue handling using existing OSATS measures as an assessment blueprint. We used kinematic features of the hands to develop prediction models of the expert ratings and compared the model performance against the variance among experts. Models consistently had less variance than the individual experts exhibited prior to consensus (see Figure 4). In this process, we identified the kinematic measures most closely linked with expert surgeons' ratings of common surgical tasks. Because this approach relies only on non-invasive video tracking and not invasive markers for motion capture, it represents a critical step forward in a scalable and reproducible method for objective assessment of surgical technical skill during actual, open operations.

Numerous features were extracted from the video. Hand position, speed, acceleration, and curvature were measured and recorded. Based on prior expectations of psychomotor performance, a series of additional variables were calculated, including peak frequency and variance, path density, working area and moving averages. Peak arrival rates in both the unfiltered and smoothed speed and acceleration signals were consistently correlated and significant in the prediction models.

The fluidity of motion and motion economy models for all suturing tasks had slopes between 0.73 and 1, and intercepts between 0.30 and 1.54. While models of tissue handling underperformed the former, the tissue handling prediction for complex or bowel anastomosis suturing tasks had a slope of 0.93 and intercept of 0.36. Many models utilized the peak rate arrivals in the acceleration signal, and achieved slopes close to 1, with intercepts between 0.5 and 1.5, which indicates excellent fit. These results support the conclusion that kinematic features of a surgeon's hand motion, measured using marker-less tracking, can accurately model subjective performance ratings made by a panel of experts. Such measures are useful in providing objective and automatic feedback, and necessary in developing evidence of validity for future competency-based assessments.<sup>48</sup>

There are several limitations to consider. First, this study was contingent on access to video of live operations. Thus, not all scores were observed across all task sub-types. Tying tasks demonstrated skewed expert ratings or had a limited range of scores, reducing the variance and making these tasks more difficult to predict, despite low error estimates. Additionally, the task-specific rating scales themselves did not always address the surgical context. For example, a surgeon with excellent fluidity may exhibit periods of frequent pauses to gather information, avoid distraction, compensate for a hand tremor, or simply resolve confusion or nervousness. While the detection itself does not suggest a particular decision or course of action to improve, automatically identifying events such as "slowing down," as Moulton<sup>35</sup> describes, could streamline video analysis to target critical points in an operation and provide insight into surgical decision making. In other words, computer vision technology

might automatically detect difficulties in an operation, demonstrated by acute changes in hand motion.

This technology also ignores the consequences to a procedure or a patient, and would need to meet additional evidence requirements as outlined by Kane and Messick<sup>12,13</sup> before any deployable assessment routine takes shape. Such assessment would also need to incorporate whether the surgeon identifies the anatomy and adjusts their technique appropriately, commits errors, or generally performs the operation correctly. A suturing task with low fluidity of motion (high hesitancy, frequent pauses) may exhibit significant motion while repositioning the driver multiple times, without contacting tissue, despite a continuous motion tracking record. It is possible that such events may be isolated by combining the different scale predictions (i.e. sudden poor tissue handling during a period of low fluidity and high motion economy may hold valuable information about the current state of the surgery.) Future work may address these challenges by considering the implications or consequences of assessment scores for a whole surgical case. These investigations may consider the quality of observed ties or stitches,<sup>49</sup> and focus on more advanced relationships between kinematic features and the overall surgical state, utilizing patterns of motion across repeated cycles<sup>21</sup> and language models<sup>26,50</sup> to automatically classify surgeons and procedures.

We anticipated that fluidity would produce the best prediction outcomes based on recorded speed, and that tissue handling would be most difficult to model, given dependence on tissue type. Similarly, the more complex anastomoses and tying tasks were expected to have greater variability in hand movements – potentially making them more difficult to predict. The tissue handling model for anastomoses (bowel and complex) performed better than tissue handling along the body wall, supporting the hypothesis that tissue handling ratings may be more sensitive to the categorization process and exhibit higher variation. It is also possible that the cues the raters used to gauge tissue handling were not as readily identifiable in the motion tracking record for suturing along the body wall (i.e. hand shape, individual finger dexterity), and therefore currently unavailable to the automatic routine.

Arranging time for the expert panel to view and rate the videos was the most difficult aspect of the study. Surgeons reported achieving consensus relatively easily after discussion, but had general difficulty when the field of view was limited. How surgeons resolved disagreements for these clips may suggest additional contextual features to improve the predictions, and may help to establish evidence of validity in the response process itself. Collecting and applying motion tracking routines to the videos played in real-time was also resource intensive and time consuming. More complex procedures with multiple changes in patient or camera positions, shifting light conditions and occlusions required active input from human observers to recalibrate and instantiate a new tracking ROI. Future software could expedite this workflow, allowing observers to adjust tracking properties on the fly. Software may also allow surgeons to rate videos independently or remotely and resolve any discrepancies in a consensus rating over time. As video collection and editing programs advance, and video recording of operations becomes more widespread, we anticipate that these processing stages will require less manual effort.

In future studies, the overall difficulty with skewed distributions could be mitigated by collecting videos of specific clinically simulated scenarios, where variety and experience is intentionally controlled to accommodate the full range of performance ratings. In this setting, the rating scales and kinematics could incorporate different surgical approaches, as some surgeons throw one handed ties while others prefer two handed ties and the kinematic path for these are quite different, limiting the features available for comparison. An expert viewer may be able to distinguish between a successfully thrown knot and a dropped ligature, but depending on how the surgeon moves after the error, their kinematic path may look identical. Such situations continue to require subjective supervisory intervention and assessment, and future work is needed to (1) automatically identify these contextual factors, and (2), assess measures of reproducibility among different stations and raters in more controlled (i.e. benchtop) simulations.

Access to immediate, reproducible kinematic-based feedback based on video review can inform self-assessment, direct practice of specific tasks, and build overall surgical skill. This may help to provide a venue in which skill development is quantifiably traceable to training interventions. If used to track expert surgeons, the capacity to deconstruct surgical skill can provide a deeper understanding of the kinematics of surgical performance and aid in the development of novel approaches to skill acquisition. Understanding the components of performance and providing such feedback can potentially shorten the current learning curve and help detect skill decay using objective measures either during periods of inactivity or toward the end of surgical careers.

The findings in this study represent a measurable step forward in creating more objective, reproducible, and accessible assessments of surgical skill commensurate with direct video observation by panel of expert raters. The prediction models have the potential to be packaged into automatic, on-demand feedback in and out of training settings, providing a reliable measure of assessment and consistent feedback to facilitate direct practice of hand motions for specific tasks. Future work is needed to explore the capacity of motion capture alone to make decisions about trainee competence. As video capture is becoming increasingly common in the operating room, computer vision motion tracking allows for a uniquely scalable approach to surgical skill analysis.

## Acknowledgments

Sources of Support:

The project described was supported by the Clinical and Translational Science Award (CTSA) program, through the NIH National Center for Advancing Translational Sciences (NCATS), grant UL1TR000427. Lane Frasier is currently supported by AHRQ F32 HS022403 and also received support from NIH T32 CA90217 and the AAS Research Fellowship Award. The content is solely the responsibility of the authors and does not necessarily represent the official views of the NIH.

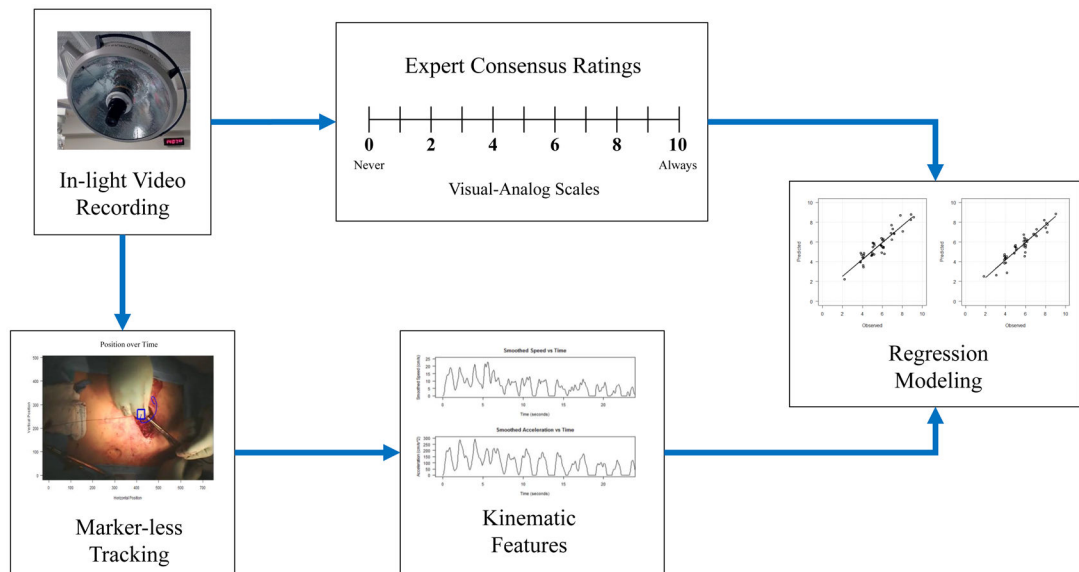
## References

1. Gavrilu D The Visual Analysis of Human Movement: A Survey. *Comput Vis Image Underst* [Internet]. 1999 1;73(1):82–98. Available from: <http://linkinghub.elsevier.com/retrieve/pii/S1077314298907160>

2. Poppe R Vision-based human motion analysis: An overview. *Comput Vis Image Underst.* 2007;108(1–2):4–18.
3. Wang L, Hu W, Tan T. Recent developments in human motion analysis. *Pattern Recognit.* 2003;36(3):585–601.
4. Birkmeyer JD, Finks JF, O'Reilly A, Oerline M, Carlin AM, Nunn AR, et al. Surgical skill and complication rates after bariatric surgery. *N Engl J Med [Internet]*. 2013;369(15):1434–42. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/24106936>
5. Reznick RK, MacRae HM. Teaching Surgical Skills — Changes in the Wind. *N Engl J Med.* 2006;355(25):2664–9. [PubMed: 17182991]
6. Darzi A, Smith S, Taffinder N. Assessing operative skill Needs to become more objective. *Bmj [Internet]*. 1999;318(7188):887–8. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/10102830>
7. Rogers SO, Gawande AA, Kwaan M, Puopolo AL, Yoon C, Brennan TA, et al. Analysis of surgical errors in closed malpractice claims at 4 liability insurers. *Surgery [Internet]*. 2006 7 [cited 2016 Dec 3];140(1):25–33. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/16857439>
8. Greenberg CC. Learning from adverse events and near misses. *J Gastrointest Surg.* 2009;13(1):3–5. [PubMed: 18797974]
9. Moorthy K, Munz Y, Sarker SK, Darzi A. Objective assessment of technical skills in surgery. *Br Med J [Internet]*. 2003;327(7422):1032–7. Available from: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC261663/>
10. Aggarwal R, Darzi A. Technical-skills training in the 21st century. *N Engl J Med.* 2006;355(25):2695–6. [PubMed: 17182997]
11. Martin J, Regehr G, Reznick R, MacRae H, Brown M, Murnaghan H, et al. Objective structured assessment of technical skill (OSATS) for surgical residents. *Br J Surg.* 1997;84:273–8. [PubMed: 9052454]
12. Kane M. Validation In: Brennan R, editor. *Educational measurement*. 4th ed. Westport, CT: American Council on Education and Praeger; 2006 p. 17–64.
13. Cook DA, Brydges R, Ginsburg S, Hatala R. A contemporary approach to validity arguments: A practical guide to Kane's framework. *Med Educ.* 2015;49(6):560–75. [PubMed: 25989405]
14. Hatala R, Cook DA, Brydges R, Hawkins R. Constructing a validity argument for the Objective Structured Assessment of Technical Skills (OSATS): a systematic review of validity evidence. *Adv Heal Sci Educ [Internet]*. 2015;20(5): 1149–75. Available from: 10.1007/s10459-015-9593-1
15. Reznick R, Regehr G, MacRae H, Martin J, McCulloch W. Testing technical skill via an innovative “bench station” examination. *Am J Surg.* 1997;173(3):226–30. [PubMed: 9124632]
16. Datta V, Bann S, Mandalia M, Darzi A. The surgical efficiency score: a feasible, reliable, and valid method of skills assessment. *Am J Surg.* 2006;192(3):372–8. [PubMed: 16920433]
17. White LW, Lendvay TS, Holst D, Borbely Y, Bekele A, Wright A. Using crowd-assessment to support surgical training in the developing world. *J Am Coll Surg [Internet]*. 2014;219(4):e40 Available from: <http://www.sciencedirect.com/science/article/pii/S1072751514010278>
18. D'Angelo ALD, Rutherford DN, Ray RD, Mason A, Pugh CM. Operative skill: Quantifying surgeon's response to tissue properties. *J Surg Res [Internet]*. 2015;198(2):294–8. Available from: 10.1016/j.jss.2015.04.078
19. Hu Y-Y, Peyre SE, Arriaga AF, Osteen RT, Corso KA, Weiser TG, et al. Postgame Analysis: Using Video-Based Coaching for Continuous Professional Development. *J Am Coll Surg [Internet]*. 2012 1;214(1):115–24. Available from: <http://linkinghub.elsevier.com/retrieve/pii/S1072751511011604>
20. Glarner CE, Hu YY, Chen CH, Radwin RG, Zhao Q, Craven MW, et al. Quantifying technical skills during open operations using video-based motion analysis. *Surg (United States) [Internet]*. 2014;156(3):729–34. Available from: <http://dx.doi.org/10.1016/j.surg.2014.04.054>
21. Watson RA. Use of a Machine Learning Algorithm to Classify Expertise: Analysis of Hand Motion Patterns During a Simulated Surgical Task. *Acad Med [Internet]*. 2014;89(8):1–5. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/24853195>
22. Mackenzie CF, Watts D, Patel R, Yang S, Hagegeorge G, Hu PF, et al. Sensor-free Computer-Vision hand-motion entropy and video-analysis of technical performance during open surgery on fresh cadavers: report of methodology and analysis. In: *Proceedings of the Human Factors and Ergonomics Society 2016 Annual Meeting*. 2016 p. 691–5.

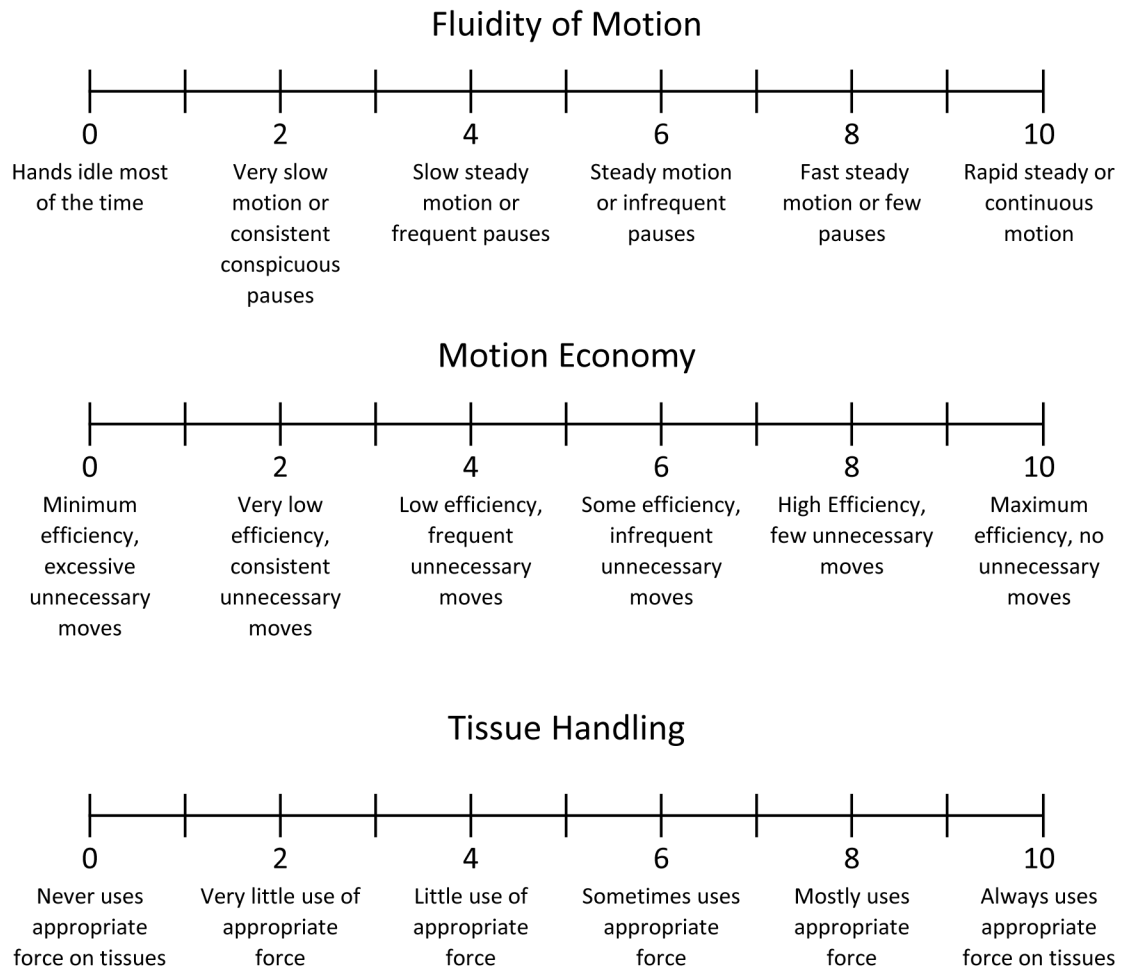
23. Hayter MA, Friedman Z, Bould MD, Hanlon JG, Katznelson R, Borges B, et al. Validation of the Imperial College Surgical Assessment Device (ICSAD) for labour epidural placement. *Can J Anesth.* 2009;56(6):419–26. [PubMed: 19340491]
24. Datta V, Mackay S, Mandalia M, Darzi A. The use of electromagnetic motion tracking analysis to objectively measure open surgical skill in the laboratory-based model. *J Am Coll Surg.* 2001;193(5):479–85. [PubMed: 11708503]
25. Datta V, Chang A, Mackay S, Darzi A. The relationship between motion analysis and surgical technical assessments. *Am J Surg.* 2002;184(1):70–3. [PubMed: 12135725]
26. Reiley CE, Lin HC, Yuh DD, Hager GD. Review of methods for objective surgical skill evaluation. *Surg Endosc Other Interv Tech.* 2011;25(2):356–66.
27. Chen CH, Hu YH, Radwin RG. A motion tracking system for hand activity assessment. 2014 IEEE China Summit Int Conf Signal Inf Process IEEE ChinaSIP 2014 - Proc 2014;320–4.
28. Chen C-H, Hu YH, Yen TY, Radwin RG. Automated Video Exposure Assessment of Repetitive Hand Activity Level for a Load Transfer Task. *Hum Factors J Hum Factors Ergon Soc* [Internet]. 2012;55(2):298–308. Available from: <http://www.scopus.com/inward/record.url?eid=2-s2.0-84875742600&partnerID=tZOTx3y1>
29. Chen C-H, Azari DP, Hu YH, Lindstrom MJ, Thelen D, Yen TY, et al. The accuracy of conventional 2D video for quantifying upper limb kinematics in repetitive motion occupational tasks [Internet] Vol. 0, *Ergonomics*. Taylor & Francis; 2015 p. 1–10. Available from: <http://www.tandfonline.com/doi/full/10.1080/00140139.2015.1051594>
30. Akkas O, Azari DP, Chen C-HE, Hu YH, Ulin SS, Armstrong TJ, et al. A hand speed – duty cycle equation for estimating the ACGIH hand activity level rating. *Ergonomics* [Internet]. 2014 24;58(2):184–94. Available from: <http://www.tandfonline.com/doi/abs/10.1080/00140139.2014.966155>
31. Azari DP, Pugh CM, Laufer S, Kwan C, Chen (Eric) C-H, Yen TY, et al. Evaluation of Simulated Clinical Breast Exam Motion Patterns Using Marker-Less Video Tracking. *Hum Factors* [Internet]. 2015;58(3):427–40. Available from: <http://hfs.sagepub.com/content/early/2015/10/31/0018720815613919.abstract>
32. Frasier LL, Azari DP, Ma Y, Quamme SRP, Radwin RG, Pugh CM, et al. A marker-less technique for measuring kinematics in the operating room. *Surgery* [Internet]. 2016; Available from: <http://linkinghub.elsevier.com/retrieve/pii/S0039606016301301>
33. Yen TY, Radwin RG. A video-based system for acquiring biomechanical data synchronized with arbitrary events and activities. *IEEE Trans Biomed Eng.* 1995;42(9):944–8. [PubMed: 7558070]
34. Annett J Subjective rating scales: science or art? *Ergonomics.* 2002;45(14):966–87. [PubMed: 12569049]
35. Moulton CA, Regehr G, Lingard L, Merritt C, MacRae H. Slowing down to stay out of trouble in the operating room: remaining attentive in automaticity. *Acad Med* [Internet]. 2010;85(10):1571–7. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/20881677>
36. D’Angelo ALD, Rutherford DN, Ray RD, Laufer S, Kwan C, Cohen ER, et al. Idle time: An underdeveloped performance metric for assessing surgical skill. *Am J Surg* [Internet]. 2015;209(4):645–51. Available from: <http://dx.doi.Org/10.1016/j.amjsurg.2014.12.013>
37. Pugh CM. Application of national testing standards to simulation-based assessments of clinical palpation skills. *Mil Med* [Internet]. 2013; 178(10 Suppl):55–63. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/24084306%5Chttp://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC4120121>
38. Laufer S, D’Angelo A-LD, Kwan C, Ray RD, Yudkowsky R, Boulet JR, et al. Rescuing the Clinical Breast Examination. *Ann Surg* [Internet]. 2016;XX(X):1 Available from: <http://content.wkhealth.com/linkback/openurl?sid=WKPTLP:landingpage&an=00000658-900000000-96384>
39. Grober ED, Roberts M, Shin EJ, Mahdi M, Bacal V. Intraoperative assessment of technical skills on live patients using economy of hand motion: establishing learning curves of surgical competence. *Am J Surg* [Internet]. 2010; 199(1):81–5. Available from: 10.1016/j.amjsurg.2009.07.033

40. Greiner TM. Hand Anthropometry of U.S. Army Personell. Tech Rep Natick. 1991;TR-92/011:434.
41. Hogan N, Sternad D. Sensitivity of smoothness measures to movement duration, amplitude, and arrests. *J Mot Behav* [Internet]. 2009;41(6):529–34. Available from: <http://www.tandfonline.com/doi/abs/10.3200/35-09-004-RC>
42. Rao C, Yilmaz A, Shah M. View-invariant representation and recognition of actions. *Int J Comput Vis.* 2002;50(2):203–26.
43. Radwin RG, Azari DP, Lindstrom MJ, Ulin SS, Armstrong TJ, Rempel D. A frequency–duty cycle equation for the ACGIH hand activity level. *Ergonomics* [Internet]. 2015 24;58(2): 173–83. Available from: <http://www.tandfonline.com/doi/abs/10.1080/00140139.2014.966154>
44. Watson RA. Computer-aided feedback of surgical knot tying using optical tracking. *J Surg Educ* [Internet]. 2012 [cited 2016 Dec 2];69(3):306–10. Available from: <http://fulltext.study/download/4298520.pdf>
45. Akaike H A new look at the statistical model identification. *IEEE Trans Automat Contr* [Internet]. 1974 12;19(6):716–23. Available from: <http://ieeexplore.ieee.org/document/1100705/>
46. Neter J, Wasserman W, Kutner M. *Applied Linear Statistical Models*. 3rd ed. Hercher RT, Shiell E, editors. Richard D. Irwin, INC.; 1990 1181 p.
47. Zumel N, Mount J. *Practical Data Science with R*. 1st ed. Shelter Island, NY: Manning INC; 2014 416 p.
48. Aggarwal R, Grantcharov T, Moorthy K, Milland T, Darzi A. Toward feasible, valid, and reliable video-based assessments of technical surgical skills in the operating room. *Ann Surg.* 2008;247(2):372–9. [PubMed: 18216547]
49. Frischknecht AC, Kasten SJ, Hamstra SJ, Perkins NC, Gillespie RB, Armstrong TJ, et al. The objective assessment of experts’ and novices’ suturing skills using an image analysis program. *Acad Med.* 2013;88(2):260–4. [PubMed: 23269303]
50. Lin HC, Shafran I, Yuh D, Hager GD. Towards automatic skill evaluation: detection and segmentation of robot-assisted surgical motions. *Comput aided Surg.* 2006;11(5):220–30. [PubMed: 17127647]

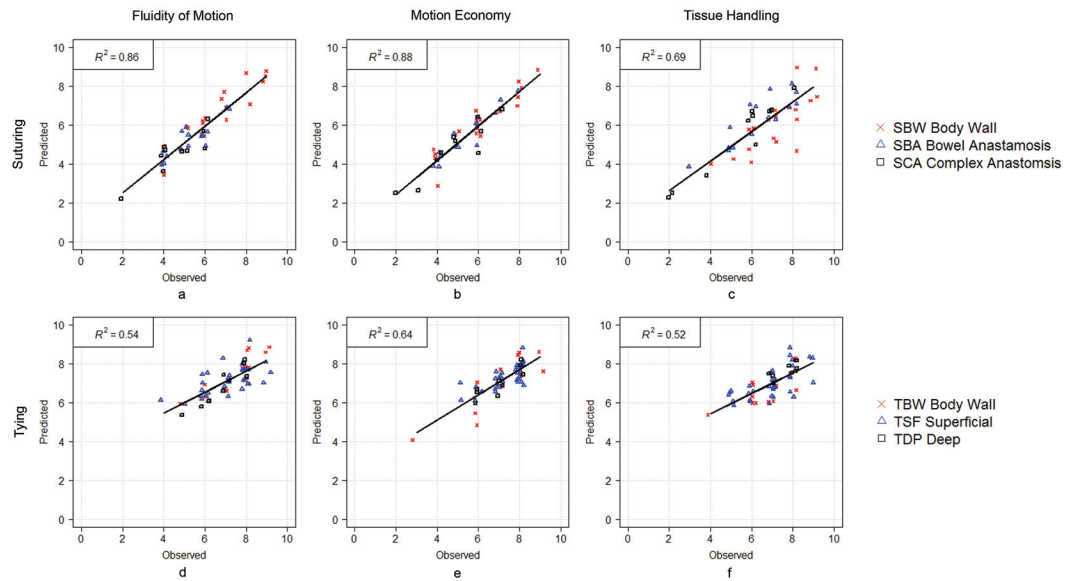


**Figure 1 –.**

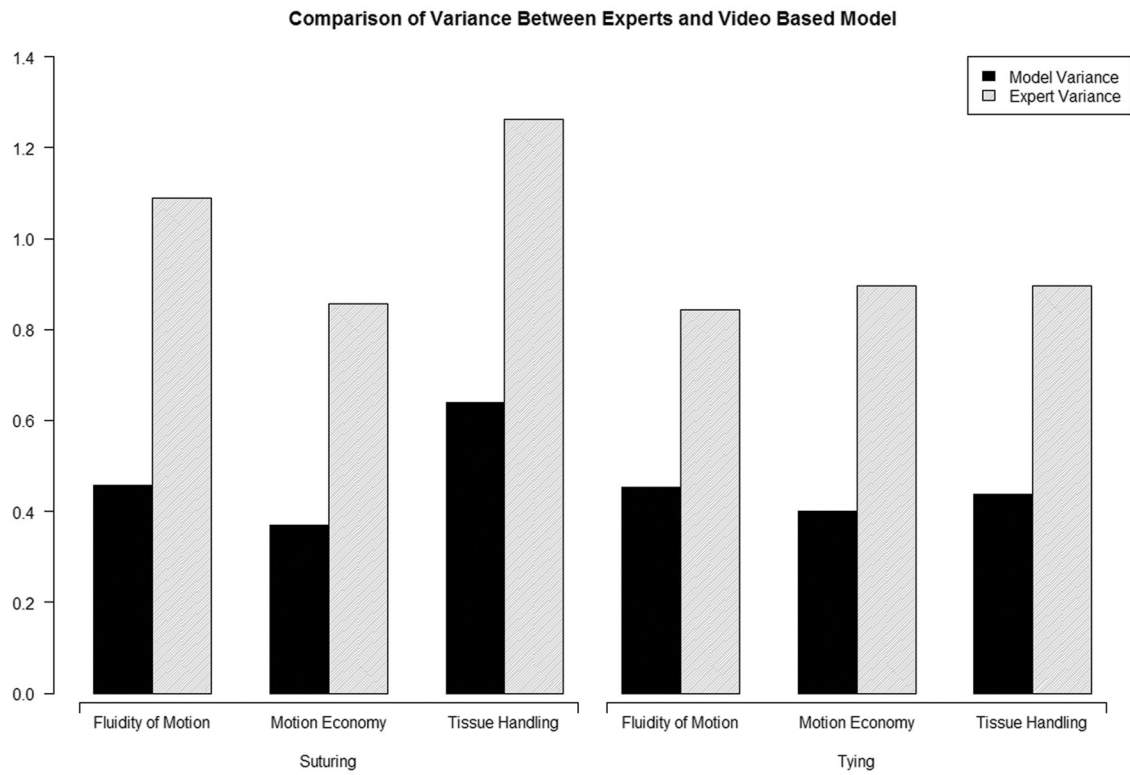
Study overview. Video from the in-light camera in the OR was collected. A panel of three expert surgeons rated task performance through a consensus process. Hand motion in the video was tracked and a series of kinematic feature characteristics were extracted. Regression models were used to predict the expert consensus ratings based on the extracted features for all the video clip segments.



**Figure 2 –.** Subjective rating scales were based on fluidity of motion, motion economy, and tissue handling for tying and suturing tasks.



**Figure 3 –.** Prediction models vs. observed results for suturing and tying tasks across all tissue types and scales for (a) suturing fluidity of motion, (b) suturing motion economy, (c) suturing tissue handling, (d) tying fluidity of motion, (e) tying motion economy, and (f) tying tissue handling.



**Figure 4 –.** Comparison of variance in video-based model (Mean Square Error) and individual prediction ratings (Mean Squared Difference) from the expert consensus ratings. The model consistently provides greater reliability than the individual experts in predicting the consensus ratings.

**TABLE 1.****Features Used for Modeling the Consensus Rating Scales**

Summary Kinematics	Mean, median and max of speed, acceleration, curvature and jerk
RMS Speed & Acceleration	Root mean-squared transformation of speed; acceleration (respectively)
Moving Averages	Simple, exponential and moving averages applied to original and smoothed speed, acceleration and curvature signals
Peak Counts	Number of peaks, applied to pure and filtered speed, acceleration, curvature and all moving averages at thresholds at intervals of 20% from 0% to 100% of maximum signal amplitude and intervals of 20 mm/s for speed (up to 100) and 200 mm/s <sup>2</sup> for acceleration (up to 1000)
Peak Frequencies	Number of peaks divided by length of video clip (s) for each type of peak count category
Peak Variation	Coefficient of variation in peak arrivals (standard deviation / mean) for each peak count category
Idle time (%)	Percent of time spent below speed threshold (set every 5 mm/s up to 50 mm/s)
Working area	Mean variance in x-y distance calibrated positions over all recorded frames
Path Density(Avg, Med, Sd, Max)	Ratio of recurrent (range 0-7) recorded x-y positions to total recorded x-y positions
Speed Density(Avg, Med, Sd, Max)	Ratio of instantaneous speed in recurrent (range 0-7) recorded x-y positions to total recorded x-y positions
Curvature Density(Avg, Med, Sd, Max)	Ratio of instantaneous curvature values in recurrent (range 0-7) recorded x-y positions to total recorded x-y positions

Kinematic variable family names (left) and associated descriptions (right).

TABLE 2.

Regression Model Summary Statistics and Predictor Variables\*\*

Task*	Fluidity of Motion			Motion Economy			Tissue Handling		
	Pred vs Obs	Variables	p	Pred vs Obs	Variables	p	Pred vs Obs	Variables	p
Suturing SBW (n = 19)	m = 0.85 b = 0.99 R <sup>2</sup> = 0.82	Peak acceleration (T=5000)	.000	m = 0.87 b = 0.78 R <sup>2</sup> = 0.85	Peak curvature (T=0)	.001	m = 0.82 b = 0.30 R <sup>2</sup> = 0.84	Peak variance of acceleration (T=3000)	.008
		Peak smoothed speed (T=100)	.025		Peak acceleration (T=6000)	.000		Smoothed peak curvature (T=0.5)	.000
		Peak curvature (T=0)	.006		Peak speed (T=100)	.005		RMS acceleration	.000
Suturing SBA (n = 15)	m = 0.73 b = 1.43 R <sup>2</sup> = 0.62	Peak acceleration (T=1000)	.021	m = 0.89 b = 0.59 R <sup>2</sup> = 0.85	Path density of median speed (R=7)	.061	m = 0.75 b = 1.54 R <sup>2</sup> = 0.66	Path density of mean speed (V=7)	.036
		RMS speed	.008		Peak acc. (T=9000)	.013		Peak acceleration (T=9000)	.017
		Peak rate speed (T=1000)	.027		Peak acc. (T=1000)	.005		Median acceleration	.118
Suturing SCA (n = 10)	m = 0.75 b = 1.13 R <sup>2</sup> = 0.63	Median speed	.149	m = 0.85 b = 0.75 R <sup>2</sup> = 0.77	FFT of Speed	.106	m = 0.93 b = 0.36 R <sup>2</sup> = 0.90	Peak curvature (T=0)	.191
		Weighted avg. of smooth speed (W = 18s)	.061		Path density of max. curvature (T=0)	.184		Weighted average of smoothed curvature [0:1] (W = 53s)	.007
		Path density of average curvature (T=0)	.074		Peak variance of speed (T=300)	.050		Peak rate of speed (T=100)	.001
Tying TBW (n = 15)	m = 0.73 b = 2.05 R <sup>2</sup> = 0.65	Peak variance of speed (T=300)	.582	m = 0.76 b = 1.65 R <sup>2</sup> = 0.70	Peak acceleration (T=8000)	.020	m = 0.57 b = 2.88 R <sup>2</sup> = 0.46	Path density of standard deviation in curvature (R=2)	.050
		Weighted average of acceleration	.001		Weighted average of curvature	.000		Peak variance of speed (T=700)	.004
		Fast Fourier transform (FFT) of speed	.004		Peaks of curvature [0:1], (T=80% of Maximum)	.011		Maximum speed	.040
Tying TSF (n = 35)	m = 0.39 b = 4.36 R <sup>2</sup> = 0.33	RMS speed	.000	m = 0.46 b = 3.96 R <sup>2</sup> = 0.40	Idle Time (T=50%)	.005	m = 0.45 b = 3.83 R <sup>2</sup> = 0.40	Peak acceleration (T=1000)	.092
		Peak acceleration (T=1000)	.073		Peak acceleration (T=9000)	.009		Weighted average of smoothed curvature	.0005
		Peak speed (T=500)	.727		Maximum Jerk	.005		Path density of position (R=7)	.017
Tying TDP (n = 9)	m = 0.88 b = 0.79 R <sup>2</sup> = 0.85	Idle time (T=30%)	.047	m = 0.74 b = 1.81 R <sup>2</sup> = 0.65	Working area	.015	m = 0.56 b = 3.37 R <sup>2</sup> = 0.30	FFT of acceleration	.040
		Peak speed (T=600)	.028		FFT of speed	.018		Peak variance of speed (T=0)	.132
		Peak speed (T=300)	.199		Working area	.334		Path density of med. speed (R=7)	.834
							FFT of acceleration	.131	

Acceptable models (slope between 0.5 and 1.5, with an intercept of 2 or less and R<sup>2</sup> > 0.75) are indicated in darker boxes and bold text.

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

\* SBW = Suturing on the body wall; SBA = Suturing as bowel anastomosis; SCA = Suturing as complex anastomosis; TBW = Tying on the body wall; TSF = intra-abdominal superficial tying; TDP = intra-abdominal deep tying.

\*\* Pred = Predicted, Obs = Observed, m = Slope, and b = Intercept, with Predicted =  $m(\text{Observed}) + b$ ;  $R^2$  = Adjusted coefficient of determination. Avg. = Average; Med. = Median; Max. = Maximum; RMS = Root Mean Square; FFT = Fast Fourier Transform; T = Threshold in mm/s for speed or  $\text{mm/s}^2$  for acceleration; R = Number of recurrent position traversal as threshold for density calculation; W = Window.