

1 **Expert Surgeons and Deep Learning Models Can Predict the Outcome of Surgical Hemorrhage**  
2 **from One Minute of Video**

3  
4 Dhiraj J Pangal BS<sup>1</sup>, Guillaume Kugener MEng<sup>1</sup>, Yichao Zhu BS<sup>2</sup>, Aditya Sinha BS<sup>1</sup>, Vyom Unadkat  
5 BS<sup>2</sup>, David J Cote MD, PhD<sup>1</sup>, Ben Strickland MD<sup>1</sup>, Martin Rutkowski MD<sup>3</sup>, Andrew Hung MD<sup>4</sup>,  
6 Animashree Anandkumar PhD<sup>5,6</sup>, X.Y. Han MS<sup>7</sup>, Vardan Papyan PhD<sup>8</sup>, Bozena Wrobel MD<sup>9</sup>, Gabriel  
7 Zada MD MS<sup>1</sup>, Daniel A Donoho MD\*<sup>10</sup>

8  
9 <sup>1</sup>Department of Neurosurgery, Keck School of Medicine of the University of Southern California, Los  
10 Angeles CA

11 <sup>2</sup>Viterbi School of Engineering, University of Southern California, Los Angeles, CA

12 <sup>3</sup>Department of Neurosurgery, Medical College of Georgia, Augusta, GA

13 <sup>4</sup>Center for Robotic Simulation and Education, USC Institute of Urology, Keck School of Medicine of  
14 the University of Southern California, Los Angeles CA

15 <sup>5</sup>Department of Computer Science + Mathematics, California Institute of Technology, Pasadena CA

16 <sup>6</sup>Nvidia Corp., Santa Clara CA

17 <sup>7</sup>Department of Operations Research and Information Engineering, Cornell University, Ithaca NY

18 <sup>8</sup>Department of Mathematics, University of Toronto, Ontario, Canada

19 <sup>9</sup>Department of Otolaryngology, Keck School of Medicine of the University of Southern California,  
20 Los Angeles CA

21 <sup>10</sup>Division of Neurosurgery, Center for Neuroscience, Children's National Hospital, Washington DC

22

23

24

25 Corresponding Author\*:

26 Daniel A. Donoho MD

27 Division of Neurosurgery

28 Center for Neuroscience, Children's National Hospital

29 Washington DC 20010

30 danieldonohomd@gmail.com

31 Office Phone: (202)-476-3020

32

33

34 **Word Count (Manuscript): 3065**

35 **Competing Interests:** The authors declare that there are no competing interests

36 **Author Contributions:**

37 Study Design: DJP, GK, AS, GZ, DAD

38 Data Acquisition: DJP, GK, BS, MR, GZ, DAD

39 Model Development: DJP, GK, AS, VU, XH, VP, DAD

40 Statistical Analysis: DJP, GK, DAD

41 Writing- Original Draft: DJP, GK, DAD

42 Writing- Revisions: All Authors

43 Final Approval: All Authors

44 Study Supervision: GZ, DAD

45

**NOTE:** This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

47 **Abstract**

48

49 **Background:** Major vascular injury resulting in uncontrolled bleeding is a catastrophic and often fatal  
50 complication of minimally invasive surgery. At the outset of these events, surgeons do not know how  
51 much blood will be lost or whether they will successfully control the hemorrhage (achieve  
52 hemostasis). We evaluate the ability of a deep learning neural network (DNN) to predict hemostasis  
53 control ability using the first minute of surgical video and compare model performance with human  
54 experts viewing the same video.

55

56 **Methods:** The publicly available SOCAL dataset contains 147 videos of attending and resident  
57 surgeons managing hemorrhage in a validated, high-fidelity cadaveric simulator. Videos are labeled  
58 with outcome and blood loss (mL). The first minute of 20 videos was shown to four, blinded,  
59 fellowship trained skull-base neurosurgery instructors, and to SOCALNet (a DNN trained on  
60 SOCAL videos). SOCALNet architecture included a convolutional network (ResNet) identifying  
61 spatial features and a recurrent network identifying temporal features (LSTM). Experts independently  
62 assessed surgeon skill, predicted outcome and blood loss (mL). Outcome and blood loss predictions  
63 were compared with SOCALNet.

64

65 **Results:** Expert inter-rater reliability was 0.95. Experts correctly predicted 14/20 trials (Sensitivity:  
66 82%, Specificity: 55%, Positive Predictive Value (PPV): 69%, Negative Predictive Value (NPV):  
67 71%). SOCALNet correctly predicted 17/20 trials (Sensitivity 100%, Specificity 66%, PPV 79%,  
68 NPV 100%) and correctly identified all successful attempts.

69

70 Expert predictions of the highest and lowest skill surgeons and expert predictions reported with  
71 maximum confidence were more accurate. Experts systematically underestimated blood loss (mean  
72 error -131 mL, RMSE 350 mL,  $R^2$  0.70) and fewer than half of expert predictions identified blood  
73 loss > 500mL (47.5%, 19/40). SOCALNet had superior performance (mean error -57 mL, RMSE  
74 295mL,  $R^2$  0.74) and detected most episodes of blood loss > 500mL (80%, 8/10).

75

76 In validation experiments, SOCALNet evaluation of a critical on-screen surgical maneuver and  
77 high/low-skill composite videos were concordant with expert evaluation.

78

79 **Conclusion:** Using only the first minute of video, experts and SOCALNet can predict outcome and  
80 blood loss during surgical hemorrhage. Experts systematically underestimated blood loss, and  
81 SOCALNet had no false negatives. DNNs can provide accurate, meaningful assessments of surgical  
82 video. We call for the creation of datasets of surgical adverse events for quality improvement  
83 research.

84

## 85 **Expert Surgeons and Deep Learning Models Can Predict the Outcome of Surgical Hemorrhage** 86 **from One Minute of Video**

### 87 88 **Introduction:**

89 Major bleeding complications during minimal access, endoscopic or robotic-assisted surgery can impair  
90 visualization and requires immediate action to control.<sup>1,2</sup> Despite maximal efforts, including the  
91 conversion from minimally invasive to ‘open’ surgery, 13-60% of major vascular injuries result in  
92 patient death.<sup>2-6</sup> Surgeon assessments of the likelihood of achieving hemostasis and the need for blood  
93 transfusion should be made immediately; however, inexperience, inability<sup>7-11</sup> and stress<sup>1,3,12,13</sup> impair  
94 decision-making, and surgeon self-assessments of the likelihood of controlling an unexpected vascular  
95 complication are uncorrelated with their actual performance.<sup>14</sup> Inaccurate predictions of blood loss and  
96 task outcome risk patient harm by delaying changes in technique, aid from surgical colleagues, or  
97 transfusion of blood products. Rather than waiting for a patient’s clinical deterioration, early prediction  
98 of difficulty at achieving hemostasis and high-volume blood loss using computer vision (CV)  
99 techniques could optimize patient outcomes.

100

101 We created SOCAL (Simulated Outcomes following Carotid Artery Laceration), a video dataset of  
102 attending and resident surgeons (otorhinolaryngologists and neurosurgeons) controlling life-threatening  
103 internal carotid artery injury (ICAI) in a validated, high-fidelity bleeding cadaveric simulator.<sup>14-18</sup>  
104 Carotid injury is a catastrophic complication of endonasal surgery and results in up to 30% mortality,  
105 similar to vascular injuries during minimally-invasive abdominal and thoracic surgery.<sup>5,19,20</sup> In prior  
106 work, we applied artificial intelligence (AI) methods to SOCAL video and developed tools that quantify  
107 blood loss and measure surgeon performance metrics from video.<sup>21,22</sup> Using these tools, we showed that  
108 video contains signals of surgical task outcome, but we do not know whether the model can detect  
109 predictive signals early in a bleeding episode, nor its performance compared to gold-standard human  
110 experts

111

112 We provided human experts (fellowship trained skull-base neurosurgeons) with the first minute of 20  
113 videos from SOCAL (‘Test Set’) and collected predictions of blood loss and task success over the entire  
114 unseen task. Experts’ predictions of outcome and blood loss established a benchmark of human  
115 performance. We then built a deep learning neural network (DNN) trained on the SOCAL video dataset  
116 (excluding the Test Set), called SOCALNet, and compared model performance on the Test Set to expert  
117 benchmarks. We validated SOCALNet predictions in subsequent experiments. To the authors  
118 knowledge this is the first comparison of DNN-derived surgical video outcome prediction to human  
119 experts viewing the same video.

120

121

122 **Methods:**

123 Experimental Design:

124 Experimental setup, data collection, consent and implementation parameters for the dataset are found  
125 in Appendix 1. Seventy-five surgeons ranging from junior trainees to world experts on endoscopic  
126 endonasal approaches (EEA) were recorded in a nationwide, validated, high-fidelity training exercise.  
127 Surgeons attempted to control an ICAI in a cadaveric head perfused with blood substitute. Performance  
128 data and intraoperative video was used to develop the SOCAL database.<sup>14-18,23</sup> The SOCAL database  
129 was developed in concordance with previously published methods, and is publicly available.<sup>23-25</sup> The  
130 SQUIRE reporting guidelines were followed.<sup>26</sup> The study was approved by the IRB of the University  
131 of Southern California. All research was performed in accordance with relevant regulations/guidelines.  
132 No patient data was utilized therefore patient-level informed consent was waived. Participating  
133 surgeons' consent was obtained for intraoperative video recording. Surgeon-expert consent was  
134 obtained.

135

136 Datasets:

137 The 147 videos in SOCAL were divided into a training set of 127 videos and a separate test set of 20  
138 videos. Ten videos depicting successes and 10 of failure were initially chosen at random for the test set;  
139 ultimately, 11 success videos (and 9 failures) were used due to ease of video formatting. Videos were  
140 truncated after 60 seconds. Only videos in the test set were shown to experts for grading.

141

142 SOCALNet Model Development:

143 See eSuppl for model code. Video was sampled at 1 frame-per-second (fps) and input into two layers,  
144 a feature generating layer and a temporal analysis algorithm (**Figure 1**). The output of the model was a  
145 binary prediction of surgical ability (trial success or failure) and estimated blood loss over the entire  
146 trial (in milliliters).

147 For the feature generator, we utilized a Residual Learning Neural Network (ResNet) model pretrained  
148 on the ImageNet 2012 classification dataset.<sup>27,28</sup> ResNet is a single-stage convolutional neural network  
149 (CNN) which uses skip connections to allow for large networks with many layers to skip layers that  
150 hurt overall performance. ResNet has become ubiquitous for object detection and classification in  
151 computer vision (CV).<sup>28</sup> The final three layers of the ResNet were retrained on SOCAL images to detect  
152 features indicative of blood loss or task success. Features from the penultimate layer of the ResNet and  
153 manual instrument annotations were passed into a bi-layer Long Short-Term Memory (LSTM) recurrent  
154 neural network.<sup>29</sup> LSTM cells contain an input, output and forget gate, allowing the network to regulate  
155 the flow of information across cells. Instrument annotations alone are inadequate for outcome  
156 prediction; successful detectors incorporate instrument data and image features.<sup>21</sup>

157

158 Expert Assessment:

159 Experts were four skull base fellowship-trained neurosurgeon instructors in ICAI management. Experts  
160 watched the 20, one-minute test videos and provided: blood loss estimates (in mL), outcome predictions  
161 (success/failure), and surgeon grades (1-5 Likert scale, 1 represents novice and 5 represents master).  
162 Experts also reported self-confidence in their outcome prediction (1-5 Likert scale; 5 represents most  
163 confident). Prior to grading, experts watched anchoring videos of novice, average, and master  
164 performances with respective outcomes data. Anchoring videos were not contained in the Test-Set, and  
165 were chosen as representative videos of each skill level by adjudication by the study team. Grading  
166 sessions were conducted in double-blinded fashion by the lead author (DJP) and individual experts (BS,  
167 MR, GZ, DAD, referred to as S1-S4). Given high concordance, mean and mode are reported for experts  
168 ('S').

169

170 Validation Analysis:

171 We conducted two experiments to evaluate model and expert concordance. In experiment one, two  
172 videos were identified in the Test-Set which where a critical error occurred shortly after the 1-minute  
173 video sample concluded (i.e., not shown to the model or surgeons). The model and all surgeons  
174 predicted, incorrectly, that both videos were successes. A new, one minute clip was generated showing  
175 the critical error and its aftermath. These new clips were evaluated by one of the human experts and  
176 SOCALNet.

177

178 In a second experiment, the three best (least blood loss, successes) and worst (most blood loss, failures)  
179 videos were identified from within the Test-Set. Composite 'best' and 'worst' videos were constructed  
180 by combining the first 20 seconds of each of the three best and worst trials in each possible order  
181 permutation (6 'best', 6 'worst' videos). The twelve composite videos were then presented to  
182 SOCALNet.

183

184 Statistical Analysis:

185 Blood loss prediction was reported using mean error, root mean square error (RMSE), and Pearson's  
186 correlation coefficients. Categorical inter-rater reliability was calculated using Cohen's Kappa and  
187 Krippendorff's alpha for more than two raters. Continuous inter-rater reliability was calculated using  
188 Pearson's correlation coefficient and an inter-rater correlation coefficient (ICC) (>2 groups; using a  
189 two-way random effects ICC model ).<sup>30</sup> We used Fisher's exact test for categorical comparisons. We  
190 performed analysis in Python with SciPy.<sup>31</sup>

191

192

193 **Results:**

194 **Table 1** lists predictions and ground truth data. There were 11 successful trials and 9 failed trials in the  
195 Test Set, with mean blood loss of 568mL (range 20-1640 mL, mean success=323 mL, mean failure=868

196 mL). Experts correctly predicted outcome in 55/80 predictions (69%, Sensitivity: 79%, Specificity:  
197 56%). Expert predictions were concordant, with one dissent in 80 ratings (Fleiss' kappa = 0.95). The  
198 average root mean square error (RMSE) for blood loss prediction of surgeons was 351 mL (mean  
199 error=-131mL, average  $R^2 = .70$ ). Expert ICC was high at 0.72.

200

201 **Figure 2, and Supplemental Table 1** demonstrates the relationship between prediction confidence,  
202 surgeon skill and prediction accuracy. Experts were most accurate when maximally confident (5/5  
203 confidence, accuracy 88%) or viewing a surgeon they rated as having minimal (Likert scale 1, accuracy  
204 92%) or maximal skill (Likert scale 5, accuracy 79%). Predictions with non-maximal confidence (levels  
205 2-4,) were only marginally better than chance (53%,  $p=0.02$  compared to maximal confidence).  
206 Predictions of intermediate skill surgeons were also less accurate (levels 2-4, 63%,  $p=0.04$  compared to  
207 composite 1/5 and 5/5 skill).

208

209 SOCALNet correctly predicted outcome in 17/20 trials (85%, Sensitivity: 100%, Specificity: 66%),  
210 noninferior to surgeons ( $p=0.12$ ). The model predicted blood loss with a RMSE of 295 mL (mean  
211 error=-57mL,  $R^2=.74$ ) (**Figure 3**). The model and experts all predicted outcome correctly in 13/20 trials.  
212 In four trials, the model was correct and all experts incorrect, in one trial the model was incorrect, and  
213 all experts correct, and two trials all were incorrect (**Figure 4**). Correlation ( $R^2$ ) between blood loss  
214 estimates for the model, experts and ground truth are shown in **Supplemental Figure 1**, and range from  
215 0.53-0.93. Correlation between the model and the average surgeon blood loss estimate was 0.73,  
216 ranging from 0.53 to 0.74 for individual surgeons (**Table 1**).

217

218 We then evaluated trials above the 50<sup>th</sup> percentile for blood loss, where blood loss exceeded 500mL and  
219 transfusion might be needed. The model predicted a blood loss estimate above 500 mL in 80% (8/10)  
220 compared to experts 47.5% (19/40); this difference was not statistically significant ( $p=0.09$ ).

221

222 Exploratory Model-Validation:

223 **Supplemental Table 2** reports model-validation experiments. In two trials, experts and SOCALNet  
224 predicted success, but the surgeon failed due to a critical error shortly after the end of the one-minute  
225 clip (therefore unseen by experts and SOCALNet). When we included the critical error, the model  
226 accurately predicted 'failure', as did an expert. In a second experiment, SOCALNet viewed six  
227 composite 'Best' trials and uniformly predicted success with low blood loss (328-473 mL); conversely,  
228 in six composite 'Worst' videos the model uniformly predicted failure with high blood loss (792-  
229 794mL).

230

231

232

233 **Discussion:**

234 To address the need for datasets depicting surgical adverse events we created SOCAL, a public video  
235 dataset of 147 attempts to control carotid injury in high-fidelity perfused cadavers. In this work we  
236 compared human expert predictions of outcome using one minute of video from 20 trials in the dataset  
237 to those of a DNN (SOCALNet). Compared to expert benchmarks, SOCALNet met or surpassed expert  
238 prediction performance, despite its relatively primitive architecture and small training data size relative  
239 to CV tasks. We synthesized counterfactual videos of excellent and poor surgeon performance to  
240 challenge SOCALNet, and it correctly predicted the outcomes in these challenges. SOCALNet and  
241 other CV methods can aid surgeons by quantifying and predicting outcome during surgical events, and  
242 in automatic video review. The absence of video datasets containing adverse events is a critical unmet  
243 need preventing the development of predictive models to improve surgical care.

244

245 **Benchmark Performance of Human Experts:**

246 Expert predictions were highly concordant, indicating that experts detected similar signals of blood loss  
247 and outcome (cross-correlation:  $R^2 = 0.74$  -0.93, Kappa for success prediction=0.95). Experts had  
248 uniform definitions of success (hemostasis) and were familiar with the stepwise progression of a well-  
249 described technique.<sup>18,32</sup> Thus, it is reasonable to conclude that using the first minute of video of a  
250 bleeding event, human experts detect signals predictive of blood loss and task outcome.

251

252 Although experts had reasonably accurate outcome and blood loss predictions (69% accuracy,  $R^2 = 0.7$ ),  
253 experts systematically overestimate surgeon success and underestimate bleeding: 4/6 of expert errors  
254 were false ‘success’ predictions, experts systematically underestimated blood loss by 131 mL and  
255 experts failed to identify 52% of high blood loss (above 500 mL) events. This post-hoc cutoff of 500mL  
256 represents a potential clinical marker of need for transfusion. The tendency for human experts to  
257 underestimate blood loss is well documented,<sup>33–36</sup> corroborated by our findings, and may result in  
258 delayed recognition of life-threatening hemorrhage.

259

260 To validate individual ratings, we asked experts to provide their confidence in each prediction, and  
261 perceived skill rating of the participating surgeon. Maximally confident predictions were more likely to  
262 be correct, as expected from prior work.<sup>33,34,37</sup> Similarly, predictions were most accurate when  
263 evaluating highest and lowest-skilled surgeons (skill rating 1 or 5), but scarcely better than chance when  
264 evaluating intermediate surgeons. Intermediate skill surgeons comprised half of all surgeons and may  
265 benefit greatly from performance assessments.

266

267 During a real vascular injury, estimation ability of the average surgeon is likely to be inferior to our  
268 experts calmly rating a single stereotyped task after training with videos of known blood loss. Experts’  
269 systematic underestimation of blood loss and struggle to assess performance of intermediate surgeons

270 represents a chasm in surgeon-assessment proficiency. Surgical patients may benefit from novel  
271 methods that improve on these benchmarks.

272

273 **SOCALNet Performance Compared to Experts:**

274 We designed a primitive deep-learning architecture containing a standard CNN and a recurrent neural  
275 network, which we call SOCALNet. We provided SOCALNet with short videos from a much smaller  
276 training dataset than is customary in CV. Despite these disadvantages, SOCALNet made statistically  
277 non-inferior (and numerically superior) outcome predictions and superior blood loss predictions  
278 compared to human experts. SOCALNet's predictions of blood loss had a smaller mean underestimation  
279 and standard error. Unlike experts, SOCALNet predictions were accurate for intermediate-skill  
280 surgeons.

281

282 The advantages of SOCALNet support the development of computer vision tools for surgical video  
283 review and as potential teammates for surgeons.<sup>38</sup> SOCALNet demonstrates that CV models can  
284 provide accurate, clinically meaningful analyses of surgical outcome from video. Future models could  
285 leverage the vast but largely untapped collections of surgical videos. Workflows developed in building  
286 SOCALNet can guide model deployment for other surgical adverse events. Human-AI teaming is a  
287 validated concept in other domains.<sup>39-41</sup> A SOCALNet-and-expert combined team (with model as a  
288 tiebreaker, particularly when expert confidence was low) would have generated 18/20 correct  
289 predictions. Furthermore, the only two inaccurate predictions from this teaming occurred when a critical  
290 error was made after the video ceased, and these errors were detected by the model and experts. If  
291 utilized at scale, AI-driven video analysis may quantify comparisons of surgical technique, provide real-  
292 time feedback for trainees, or provide guidance during rare scenarios a surgeon may not have  
293 encountered (e.g. vascular injury) but the model has been trained on.<sup>38</sup>

294

295 SOCALNet has room for improvement. For adverse events, the 1) accurate estimation of high-volume  
296 blood-loss and 2) detection of task failures may be prioritized as exsanguination is life-threatening.  
297 SOCALNet blood loss predictions exhibited more robust central tendency than experts, resulting in  
298 better predictions for typical performances. However, when grading edge cases of the two worst  
299 surgeons in the Test Set, SOCALNet underestimated blood loss (absolute error of 790-800 mL on  
300 videos exceeding 1.5L of blood loss). In predicting failure (specificity), both experts and SOCALNet  
301 showed limitations (Specificity= 0.56, 0.66 respectively); however, improving expert predictions are  
302 challenging, and most surgeons are non-experts. Accordingly, applying CV optimization techniques to  
303 AI models (e.g. cost-sensitive classification, oversampling) may be preferred.<sup>42,43</sup>

304

305 **Surgical Adverse Event Video Datasets: An Unmet Need in Surgical Safety:**



306 A growing body of evidence supports the quantitative analysis of surgical video.<sup>22,44–47</sup> One fundamental  
307 discovery has been the detection of signals in surgical video that predict patient outcome: surgeons have  
308 heterogeneous skill resulting in heterogeneous outcomes.<sup>14,44,45,48</sup> Although low-skill surgeons are more  
309 likely to have adverse intraoperative events, video of these events has not been systematically studied.  
310 Instead of studying surgical video, studies describe adverse events using textual medical records,  
311 radiography, and laboratory results. Analysis of these extra-operative records and correlations with pre-  
312 operative risk factors and post-operative management can be useful.<sup>49–53</sup> However, this research omits  
313 a crucial determinant of the outcome of the surgical patient: the surgical event itself. This omission  
314 limits root-cause analysis to only the extra-operative universe and prevents evaluation of the technical  
315 maneuvers and patient anatomic conditions that make adverse events more likely. Unlike textual  
316 records, surgical video depicts all visualized surgeon movements and patient anatomy, making video  
317 uniquely suited for the study of operative events. The results of the present study begin to demonstrate  
318 the value of studying video of surgical adverse events.

319

320 We propose the creation of large, multi-center datasets of surgical videos that includes adverse  
321 events.<sup>54,55</sup> Video datasets of surgical adverse events can be leveraged using predictive models (e.g.,  
322 SOCALNet) which can detect intraoperative events, evaluate performance and quantify technique. This  
323 study was supported the North American Skull Base Society, whose mission is to promote scientific  
324 advancement, share outcomes data for education and to advance outcomes research. Groups such as the  
325 Michigan Bariatric Surgery Collaborative and the Michigan Urologic Surgery Improvement  
326 Consortium have conducted similar work and we hope to call their attention to adverse events in  
327 addition to routine procedures.<sup>56,57</sup> National organizations capable of soliciting large bodies of data  
328 should prioritize collecting adverse event videos and apply technical innovations adopted by other  
329 medical fields to ensure privacy and confidentiality.<sup>58–60</sup> National organizations can also facilitate the  
330 scaling of expert labeling. Small groups face long delays in accruing sufficient cases and labeling video.  
331 In this study, despite a long term track record of collaboration amongst our team, it required two months  
332 for our experts to review 20 minutes of aggregated video.<sup>61</sup> Collaborative efforts may be able to require  
333 video review as a condition of membership.

334

335 Finally, high-fidelity simulation enables analysis of rare surgical events. Curating 150 videos of real  
336 carotid injuries would require tens of thousands of cases, an impossible task without streamlined data-  
337 sharing mechanisms; using perfused cadavers and real instruments we collected hundreds of  
338 observations of this otherwise rare event. Videos in the simulated environment can complement surgical  
339 video datasets that otherwise depict thousands of uncomplicated cases and only a few rare  
340 events.<sup>14,15,17,18,62–65</sup> As more surgical video datasets are developed, we can follow the ‘sim-to-real’  
341 process where models are trained on virtual data and then fine-tuned and validated in the real  
342 environment.<sup>66–68</sup>

343

344 **Limitations:**

345 Our study has several limitations. First, validation on clinical video is a clear next step, although  
346 accruing a corpus of carotid injury video would likely require substantial national efforts. Second,  
347 results from carotid injuries may not transfer to other vascular injuries, and vascular injuries differ from  
348 other adverse events. Rather than diminishing our results, these complementary challenges showcase  
349 the depth of unmet need within surgical-video data science. Separately from these study design  
350 limitations, SOCALNet ingests ground truth tool annotations as input, which requires pre-processing  
351 of data and is thus not fully automated.<sup>69–71</sup> The lack of curated surgical video datasets remain a major  
352 limitation for future work.

353

354 **Conclusion:**

355 Experts and a neural network can predict the outcome of surgical hemorrhage from the first minute of  
356 video of the adverse event. Neural network-based architectures can already achieve human or supra-  
357 human performance at predicting clinically relevant outcomes from video. To improve outcomes of  
358 surgical patients, advances in quantitative and predictive methods should be applied to newly collected  
359 video datasets containing adverse events.

360

361 **Data Availability:**

362 The datasets generated during and/or analyzed during the current study are available in the *figshare*  
363 repository, link: <https://doi.org/10.6084/m9.figshare.15132468.v1>

364

365

366 **References**

367

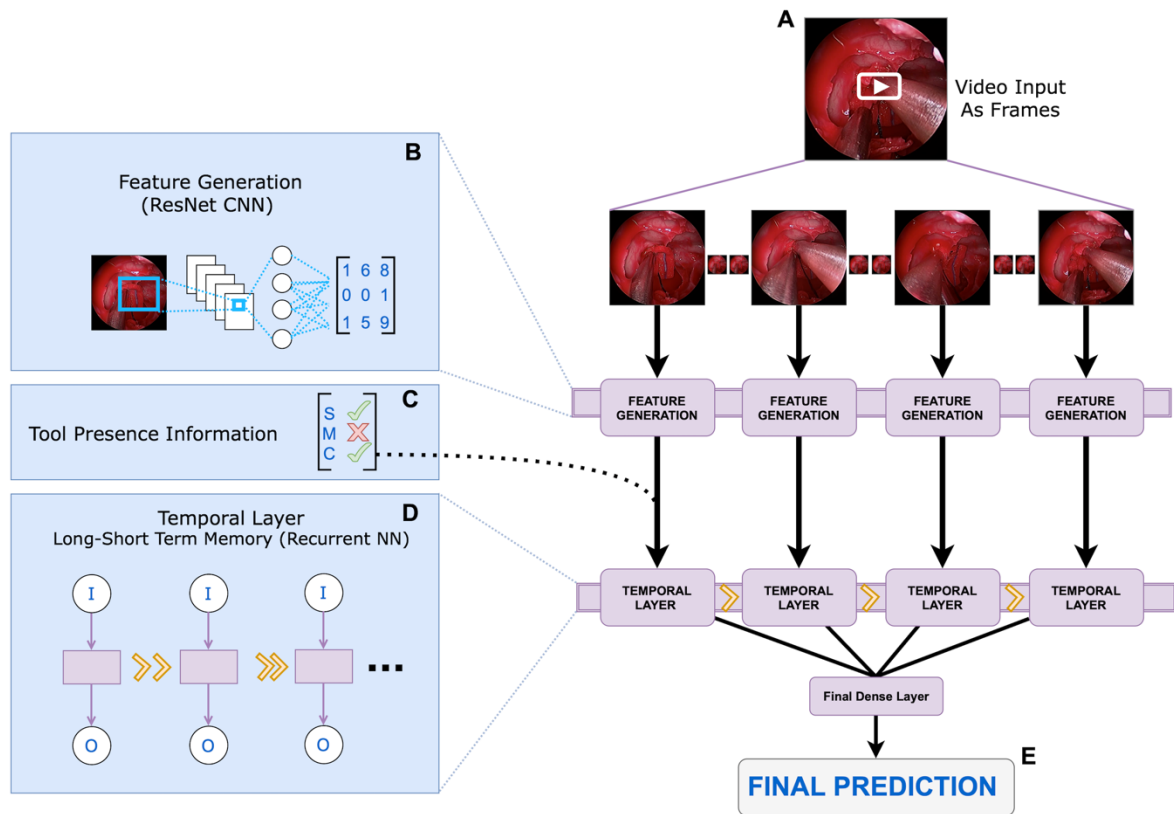
- 368 1. Lee, Y. F. *et al.* Unplanned Robotic-Assisted Conversion-to-Open Colorectal Surgery is  
369 Associated with Adverse Outcomes. *J Gastrointest Surg* **22**, 1059–1067 (2018).
- 370 2. England, E. C. *et al.* REBOA as a rescue strategy for catastrophic vascular injury during  
371 robotic surgery. *J Robot Surg* **14**, 473–477 (2020).
- 372 3. Sandadi, S. *et al.* Recognition and management of major vessel injury during  
373 laparoscopy. *J Minim Invasive Gynecol* **17**, 692–702 (2010).
- 374 4. Hemingway, J. F. *et al.* Intraoperative consultation of vascular surgeons is increasing at a  
375 major American trauma center. *J Vasc Surg* **74**, 1581–1587 (2021).
- 376 5. Laws, E. R. Vascular complications of transsphenoidal surgery. *Pituitary* **2**, 163–170  
377 (1999).
- 378 6. Beekley, A. C. Damage control resuscitation: a sensible approach to the exsanguinating  
379 surgical patient. *Crit Care Med* **36**, S267–274 (2008).
- 380 7. Tisherman, S. A. Management of Major Vascular Injury: Open. *Otolaryngol Clin North*  
381 *Am* **49**, 809–817 (2016).
- 382 8. Melnic, C. M., Heng, M. & Lozano-Calderon, S. A. Acute Surgical Management of  
383 Vascular Injuries in Hip and Knee Arthroplasties. *J Am Acad Orthop Surg* **28**, 874–883  
384 (2020).

- 385 9. Quasarano, R. T., Kashef, M., Sherman, S. J. & Hagglund, K. H. Complications of  
386 gynecologic laparoscopy. *J Am Assoc Gynecol Laparosc* **6**, 317–321 (1999).
- 387 10. Asfour, V., Smythe, E. & Attia, R. Vascular injury at laparoscopy: a guide to  
388 management. *J Obstet Gynaecol* **38**, 598–606 (2018).
- 389 11. Filis, K. *et al.* Iatrogenic Vascular Injuries of the Abdomen and Pelvis: The Experience at  
390 a Hellenic University Hospital. *Vasc Endovascular Surg* **53**, 541–546 (2019).
- 391 12. Arora, S. *et al.* Stress impairs psychomotor performance in novice laparoscopic surgeons.  
392 *Surg Endosc* **24**, 2588–2593 (2010).
- 393 13. Jukes, A. K. *et al.* Stress response and communication in surgeons undergoing training in  
394 endoscopic management of major vessel hemorrhage: a mixed methods study. *Int Forum*  
395 *Allergy Rhinol* **7**, 576–583 (2017).
- 396 14. Donoho, D. A. *et al.* Improved surgeon performance following cadaveric simulation of  
397 internal carotid artery injury during endoscopic endonasal surgery: training outcomes of a  
398 nationwide prospective educational intervention. *Journal of Neurosurgery* **1**, 1–9 (2021).
- 399 15. Shen, J. *et al.* Objective Validation of Perfusion-Based Human Cadaveric Simulation  
400 Training Model for Management of Internal Carotid Artery Injury in Endoscopic  
401 Endonasal Sinus and Skull Base Surgery. *Oper Neurosurg (Hagerstown)* **15**, 231–238  
402 (2018).
- 403 16. Zada, G. *et al.* Development of a Perfusion-Based Cadaveric Simulation Model  
404 Integrated into Neurosurgical Training: Feasibility Based On Reconstitution of Vascular  
405 and Cerebrospinal Fluid Systems. *Oper Neurosurg (Hagerstown)* **14**, 72–80 (2018).
- 406 17. Donoho, D. A. *et al.* Costs and training results of an objectively validated cadaveric  
407 perfusion-based internal carotid artery injury simulation during endoscopic skull base  
408 surgery. *Int Forum Allergy Rhinol* **9**, 787–794 (2019).
- 409 18. Pham, M. *et al.* A Perfusion-based Human Cadaveric Model for Management of Carotid  
410 Artery Injury during Endoscopic Endonasal Skull Base Surgery. *J Neurol Surg B* **75**,  
411 309–313 (2014).
- 412 19. Ciric, I., Ragin, A., Baumgartner, C. & Pierce, D. Complications of transsphenoidal  
413 surgery: results of a national survey, review of the literature, and personal experience.  
414 *Neurosurgery* **40**, 225–236; discussion 236-237 (1997).
- 415 20. AlQahtani, A. *et al.* Assessment of Factors Associated With Internal Carotid Injury in  
416 Expanded Endoscopic Endonasal Skull Base Surgery. *JAMA Otolaryngol Head Neck*  
417 *Surg* (2020) doi:10.1001/jamaoto.2019.4864.
- 418 21. Kugener, G. *et al.* Deep Neural Networks Can Accurately Detect Blood Loss and  
419 Hemorrhage Control Task Success from Intraoperative Video. *Neurosurgery (Accepted)*.
- 420 22. Pangal, D. J. *et al.* Surgical Video-Based Automated Performance Metrics Predict Blood  
421 Loss and Success of Simulated Vascular Injury Control in Neurosurgery: A Pilot Study.  
422 *Journal of Neurosurgery (Accepted)*.
- 423 23. Pangal, D. J. *et al.* Technical Note: A Guide to Annotation of Neurosurgical  
424 Intraoperative Video for Machine Learning Analysis and Computer Vision. *World*  
425 *Neurosurg* (2021) doi:10.1016/j.wneu.2021.03.022.
- 426 24. Kugener, G., Pangal, D. J. & Zada, G. Simulated Outcomes following Carotid Artery  
427 Laceration. (2021) doi:10.6084/m9.figshare.15132468.v1.
- 428 25. Paper Information / Code Submission Policy.  
429 <https://nips.cc/Conferences/2021/PaperInformation/CodeSubmissionPolicy>.
- 430 26. Squire 2.0 (Standards for Quality Improvement Reporting Excellence): Revised  
431 Publication Guidelines From a Detailed Consensus Process | American Journal of Critical  
432 Care | American Association of Critical-Care Nurses.  
433 [https://aacnjournals.org/ajconline/article-abstract/24/6/466/4045/Squire-2-0-Standards-](https://aacnjournals.org/ajconline/article-abstract/24/6/466/4045/Squire-2-0-Standards-for-Quality-Improvement)  
434 [for-Quality-Improvement](https://aacnjournals.org/ajconline/article-abstract/24/6/466/4045/Squire-2-0-Standards-for-Quality-Improvement).

- 435 27. He, K., Zhang, X., Ren, S. & Sun, J. Deep Residual Learning for Image Recognition.  
436 *arXiv:1512.03385 [cs]* (2015).
- 437 28. Deng, J. *et al.* ImageNet: A large-scale hierarchical image database. in *2009 IEEE*  
438 *Conference on Computer Vision and Pattern Recognition* 248–255 (2009).  
439 doi:10.1109/CVPR.2009.5206848.
- 440 29. Yengera, G., Mutter, D., Marescaux, J. & Padoy, N. Less is More: Surgical Phase  
441 Recognition with Less Annotations through Self-Supervised Pre-training of CNN-LSTM  
442 Networks. *arXiv:1805.08569 [cs]* (2018).
- 443 30. Koo, T. K. & Li, M. Y. A Guideline of Selecting and Reporting Intraclass Correlation  
444 Coefficients for Reliability Research. *J Chiropr Med* **15**, 155–163 (2016).
- 445 31. Virtanen, P. *et al.* SciPy 1.0: fundamental algorithms for scientific computing in Python.  
446 *Nat Methods* **17**, 261–272 (2020).
- 447 32. Kassir, Z. M., Gardner, P. A., Wang, E. W., Zenonos, G. A. & Snyderman, C. H.  
448 Identifying Best Practices for Managing Internal Carotid Artery Injury During  
449 Endoscopic Endonasal Surgery by Consensus of Expert Opinion. *Am J Rhinol Allergy*  
450 19458924211024864 (2021) doi:10.1177/19458924211024864.
- 451 33. Thomas, S. *et al.* Measured versus Estimated Blood Loss: Interim Analysis of a  
452 Prospective Quality Improvement Study. *Am Surg* **86**, 228–231 (2020).
- 453 34. Lopez-Picado, A., Albinarrate, A. & Barrachina, B. Determination of Perioperative Blood  
454 Loss: Accuracy or Approximation? *Anesth Analg* **125**, 280–286 (2017).
- 455 35. Saoud, F. *et al.* Validation of a new method to assess estimated blood loss in the obstetric  
456 population undergoing cesarean delivery. *Am J Obstet Gynecol* **221**, 267.e1-267.e6  
457 (2019).
- 458 36. Rubenstein, A. F., Zamudio, S., Douglas, C., Sledge, S. & Thurer, R. L. Automated  
459 Quantification of Blood Loss versus Visual Estimation in 274 Vaginal Deliveries. *Am J*  
460 *Perinatol* (2020) doi:10.1055/s-0040-1701507.
- 461 37. Serapio, E. T., Pearlson, G. A., Drey, E. A. & Kerns, J. L. Estimated versus measured  
462 blood loss during dilation and evacuation: an observational study. *Contraception* **97**,  
463 451–455 (2018).
- 464 38. Ward, T. M. *et al.* Computer vision in surgery. *Surgery* **169**, 1253–1256 (2021).
- 465 39. Maia Chess. <https://maiachess.com>.
- 466 40. Zhang, R., McNeese, N. J., Freeman, G. & Musick, G. ‘An Ideal Human’: Expectations  
467 of AI Teammates in Human-AI Teaming. *Proc. ACM Hum.-Comput. Interact.* **4**, 246:1-  
468 246:25 (2021).
- 469 41. Human–AI collaboration inspires tyre innovation.
- 470 42. Elkan, C. The foundations of cost-sensitive learning. in *Proceedings of the 17th*  
471 *international joint conference on Artificial intelligence - Volume 2* 973–978 (Morgan  
472 Kaufmann Publishers Inc., 2001).
- 473 43. Teh, K., Armitage, P., Tesfaye, S., Selvarajah, D. & Wilkinson, I. D. Imbalanced  
474 learning: Improving classification of diabetic neuropathy from magnetic resonance  
475 imaging. *PLOS ONE* **15**, e0243907 (2020).
- 476 44. Birkmeyer, J. D. *et al.* Surgical Skill and Complication Rates after Bariatric Surgery. *New*  
477 *England Journal of Medicine* **369**, 1434–1442 (2013).
- 478 45. Brajcich, B. C. *et al.* Association Between Surgical Technical Skill and Long-term  
479 Survival for Colon Cancer. *JAMA Oncol* (2020) doi:10.1001/jamaoncol.2020.5462.
- 480 46. Chhabra, K. R., Thumma, J. R., Varban, O. A. & Dimick, J. B. Associations Between  
481 Video Evaluations of Surgical Technique and Outcomes of Laparoscopic Sleeve  
482 Gastrectomy. *JAMA Surg* **156**, e205532 (2021).
- 483 47. Greenberg, C. C., Dombrowski, J. & Dimick, J. B. Video-Based Surgical Coaching: An  
484 Emerging Approach to Performance Improvement. *JAMA Surg* **151**, 282–283 (2016).

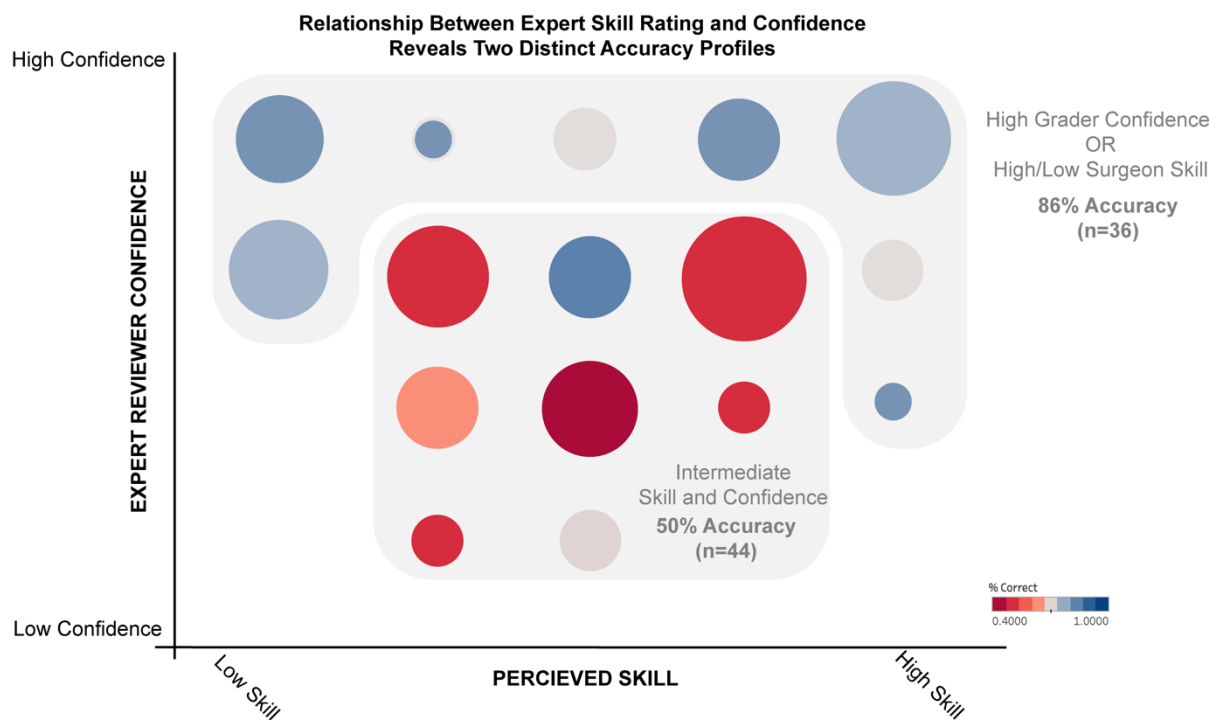
- 485 48. Stulberg, J. J. *et al.* Association Between Surgeon Technical Skills and Patient Outcomes.  
486 *JAMA Surg* (2020) doi:10.1001/jamasurg.2020.3007.
- 487 49. Elsamadicy, A. A. *et al.* Reduced Impact of Obesity on Short-Term Surgical Outcomes,  
488 Patient-Reported Pain Scores, and 30-Day Readmission Rates After Complex Spinal  
489 Fusion ( $\geq 7$  Levels) for Adult Deformity Correction. *World neurosurgery* **127**, e108–  
490 e113 (2019).
- 491 50. Jones, D. *et al.* Multicentre, prospective observational study of the correlation between  
492 the Glasgow Admission Prediction Score and adverse outcomes. *BMJ Open* **9**, e026599  
493 (2019).
- 494 51. Arango-Lasprilla, J. C. *et al.* Predictors of Extended Rehabilitation Length of Stay After  
495 Traumatic Brain Injury. *Archives of Physical Medicine and Rehabilitation* **91**, 1495–1504  
496 (2010).
- 497 52. Giannini, A. *et al.* Predictors of postoperative overall and severe complications after  
498 surgical treatment for endometrial cancer: The role of the fragility index. *Int J Gynaecol*  
499 *Obstet* **148**, 174–180 (2020).
- 500 53. Simpson, A. M., Donato, D. P., Kwok, A. C. & Agarwal, J. P. Predictors of  
501 complications following breast reduction surgery: A National Surgical Quality  
502 Improvement Program study of 16,812 cases. *J Plast Reconstr Aesthet Surg* **72**, 43–51  
503 (2019).
- 504 54. NEUROSURGERY Journal. *Carotid Injury in Endonasal Surgery*. (2013).
- 505 55. NEUROSURGERY Journal. *Managing Arterial Injury in Endoscopic Skull Base*  
506 *Surgery*. (2015).
- 507 56. Home | MBSC Coordinating Center. *Michigan Bariatric S* <https://www.mbscsurgery.org>.
- 508 57. Michigan Urological Surgery Improvement Collaborative (MUSIC).  
509 <https://musicurology.com/>.
- 510 58. Rieke, N. *et al.* The future of digital health with federated learning. *npj Digit. Med.* **3**, 1–7  
511 (2020).
- 512 59. Dou, Q. *et al.* Federated deep learning for detecting COVID-19 lung abnormalities in CT:  
513 a privacy-preserving multinational validation study. *NPJ Digit Med* **4**, 60 (2021).
- 514 60. Willemink, M. J. *et al.* Preparing Medical Imaging Data for Machine Learning.  
515 *Radiology* **295**, 4–15 (2020).
- 516 61. Lendvay, T. S., White, L. & Kowalewski, T. Crowdsourcing to Assess Surgical Skill.  
517 *JAMA Surg* **150**, 1086–1087 (2015).
- 518 62. Winer, J. L. *et al.* Cerebrospinal fluid reconstitution via a perfusion-based cadaveric  
519 model: feasibility study demonstrating surgical simulation of neuroendoscopic  
520 procedures. *J Neurosurg* **123**, 1316–1321 (2015).
- 521 63. Christian, E. A. *et al.* Perfusion-based human cadaveric specimen as a simulation training  
522 model in repairing cerebrospinal fluid leaks during endoscopic endonasal skull base  
523 surgery. *J Neurosurg* **129**, 792–796 (2018).
- 524 64. Strickland, B. A. *et al.* The Use of a Novel Perfusion-Based Human Cadaveric Model for  
525 Simulation of Dural Venous Sinus Injury and Repair. *Oper Neurosurg (Hagerstown)* **19**,  
526 E269–E274 (2020).
- 527 65. Bakhsheshian, J. *et al.* The use of a novel perfusion-based cadaveric simulation model  
528 with cerebrospinal fluid reconstitution comparing dural repair techniques: a pilot study.  
529 *The Spine Journal* **17**, 1335–1341 (2017).
- 530 66. Closing the Simulation-to-Reality Gap for Deep Robotic Learning. *Google AI Blog*  
531 <http://ai.googleblog.com/2017/10/closing-simulation-to-reality-gap-for.html>.
- 532 67. Christiano, P. *et al.* Transfer from Simulation to Real World through Learning Deep  
533 Inverse Dynamics Model. (2016).

- 534 68. Bissonnette, V. *et al.* Artificial Intelligence Distinguishes Surgical Training Levels in a  
535 Virtual Reality Spinal Task. *The Journal of Bone and Joint Surgery* **101**, (2019).
- 536 69. Kranzfelder, M. *et al.* Real-time instrument detection in minimally invasive surgery using  
537 radiofrequency identification technology. *Journal of Surgical Research* **185**, 704–710  
538 (2013).
- 539 70. Du, X. *et al.* Articulated Multi-Instrument 2-D Pose Estimation Using Fully  
540 Convolutional Networks. *IEEE Transactions on Medical Imaging* **37**, 1276–1287 (2018).
- 541 71. Staartjes, V. E., Volokitin, A., Regli, L., Konukoglu, E. & Serra, C. Machine Vision for  
542 Real-Time Intraoperative Anatomic Guidance: A Proof-of-Concept Study in Endoscopic  
543 Pituitary Surgery. *Oper Neurosurg (Hagerstown)* opab187 (2021)  
544 doi:10.1093/ons/opab187.  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555



556  
 557 **Figure 1. SOCALNet Architecture.** Deep learning model used to predict blood loss and task success  
 558 in critical hemorrhage control task. A) Video is snapshotted into individual frames. B) A pretrained  
 559 ResNet convolutional neural network (CNN) is fine-tuned on SOCAL images from (A), to find  
 560 features predictive of blood loss and task success in each individual frame. Output matrix from (B)  
 561 and tool presence information (C) [e.g. Is suction (S) present? Yes (check); is Muscle (M) present?  
 562 No (X), etc] is input into a temporal layer. D) Temporal layer: Long-short-term memory (LSTM)  
 563 modified recurrent neural network allowing for temporal analysis across all frames. All LSTM  
 564 predictions are consolidated in one dense layer and E) a final prediction of success/failure, and blood  
 565 loss (in mL) is output  
 566

567



568

569

570

571 **Figure 2. Association between expert confidence, surgeon skill level and accuracy of prediction.**

572 Experts are most accurate when viewing trials of surgeons with low or high skill, or where they

573 (experts) are maximally confident. For those with moderate skill or when experts have moderate

574 confidence, prediction accuracy is lower. Size of circle denotes number of trials. Color denotes

575 accuracy.

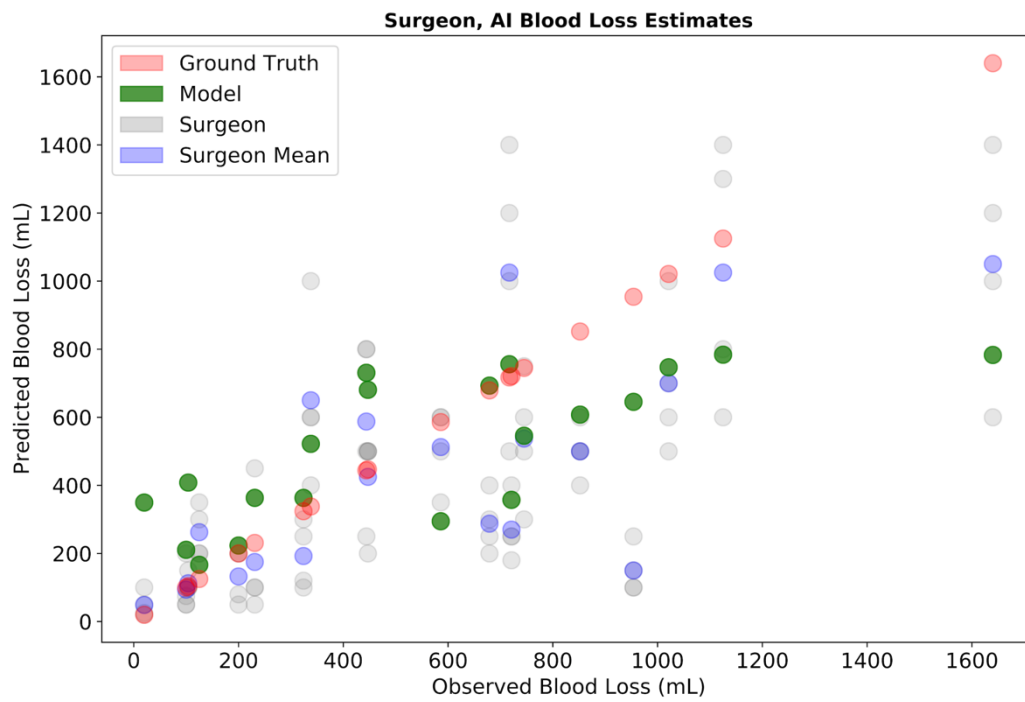
576

577

578



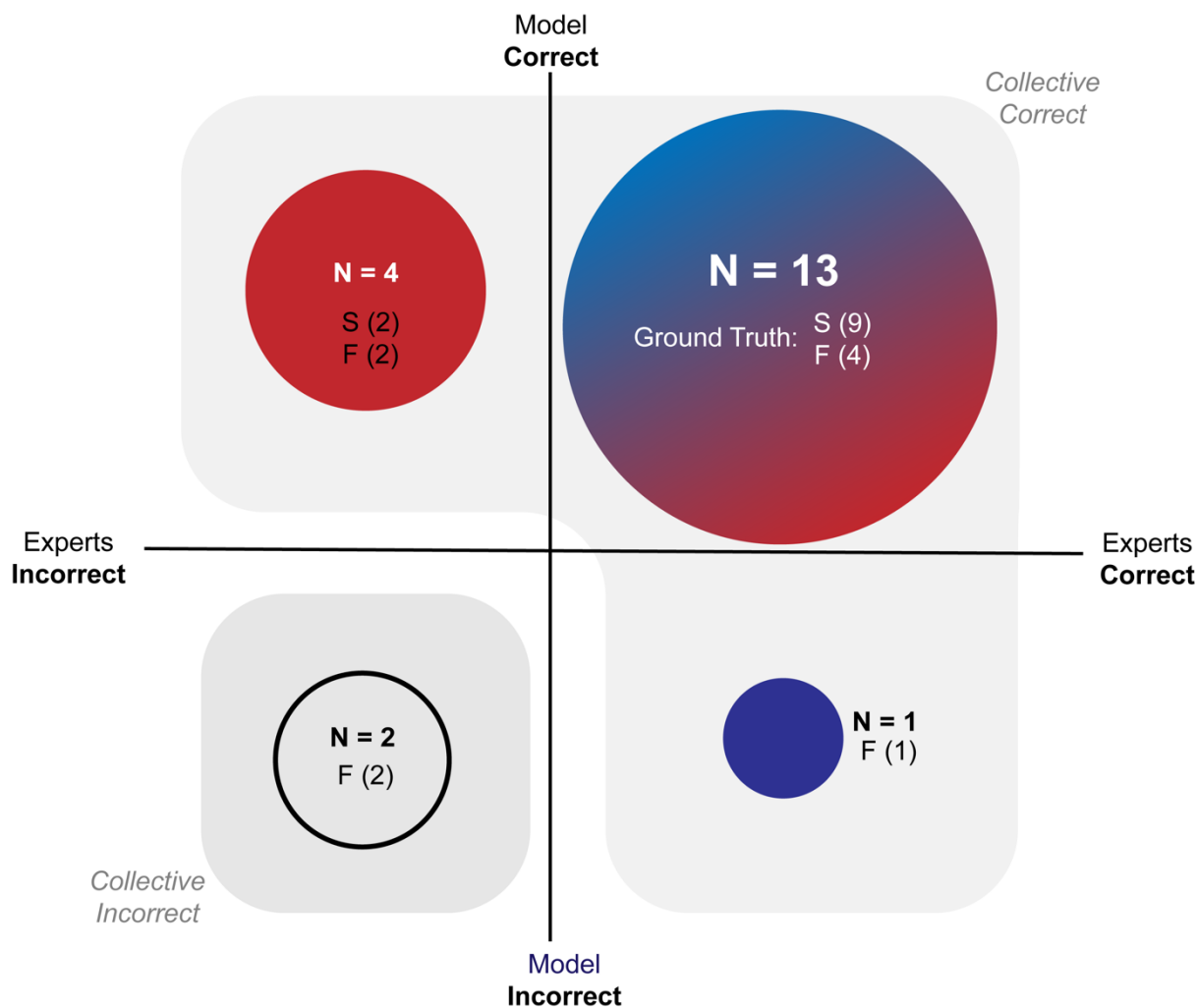
579  
580



581  
582 **Figure 3. Expert and SOCALNet Blood Loss Quantification.** Predicted versus observed blood loss  
583 estimations by individual surgeons (grey), surgeon mean (blue), and model (green). Red points  
584 represent measured blood loss (ground truth).  
585

586

### Comparison of Model and Expert Prediction Accuracy



587  
588  
589  
590  
591  
592  
593  
594  
595

**Figure 4. Outcome Predictions of Experts and SOCALNet.** Outcomes of experts (Blue) and model (Red) in predicting task success using one minute of video. Circle size denotes number of trials (N). Success (S) and failure (F) denoted underneath each N. When the union of successful predictions is taken, the model+expert grouping would successfully predict outcome in 18/20 cases. In the 2 remaining cases (bottom left quadrant), a critical error took place following the cessation of the video and was evaluated in subsequent counterfactual experiments.

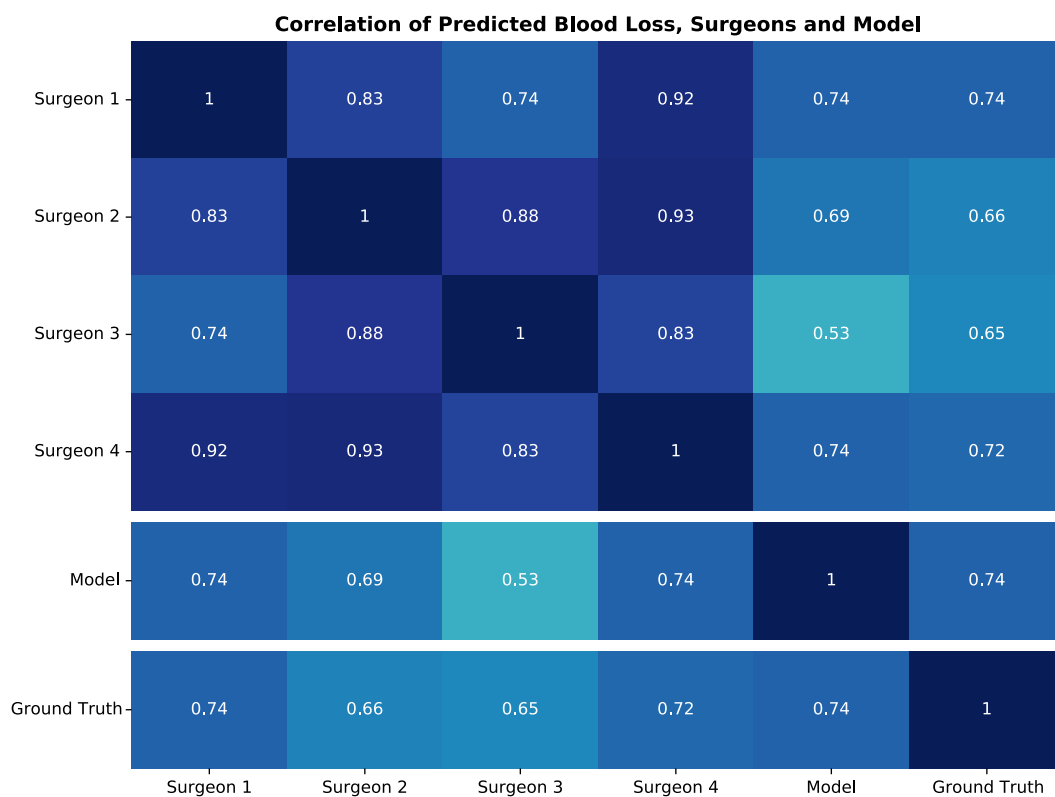
596

	<b>Accuracy (SN %, SP %)</b>	<b>RMSE (R<sup>2</sup>)</b>	<b>M-S Agreement:<sup>*</sup> <i>Success/Failure</i></b>	<b>M-S Agreement:<sup>†</sup> <i>Blood Loss</i></b>
Ground Truth	11 Success 9 Failures	-	-	Avg Blood Loss: 568 (Range:20-1640 )
Model	17/20 (85%) (100, 66)	295 (.74)	-	-
<b>Expert Cohort</b>	<b>55/80 (68.75) (79, 56)</b>	<b>351 (.70)</b>	<b>.43<sup>‡</sup></b>	<b>0.73<sup>‡</sup></b>
Surgeon 1	13/20 (65%) (73, 55)	306 (.73)	.34	.74
Surgeon 2	14/20 (65%) (81, 55)	335 (.66)	.43	.66
Surgeon 3	14/20 (65%) (81, 55)	423 (.65)	.43	.65
Surgeon 4	14/20 (65%) (81, 55)	329 (.74)	.43	.72

597  
598  
599  
600  
601  
602  
603  
604  
605

Table 1. Results comparing Deep Learning Model with Expert Surgeons. SN: Sensitivity, SP: Specificity, M-S: Model-Surgeon. \*: Kappa coefficient; †:inter-class coefficient; ‡: Inter-Surgeon Agreement: Success/Failure= 0.95, Blood-Loss: 0.72

606  
607



608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632

**Supplemental Figure 1.** Correlation ( $R^2$ ) between blood loss prediction from all 4 expert surgeon graders, model, and ground truth data.

633  
634  
635  
636  
637