

Anatomy Segmentation in Laparoscopic Surgery: Comparison of Machine Learning and Human Expertise – An Experimental Study

Fiona R. Kolbinger, MD^{1,2,3,✉}, Franziska M. Rinner¹, Alexander C. Jenke⁴, Matthias Carstens¹, Stefanie Krell⁴, Stefan Leger, PhD^{3,4}, Marius Distler, MD^{1,2}, Jürgen Weitz, MD^{1,2,3,5}, Stefanie Speidel, PhD^{3,4,5}, Sebastian Bodenstedt, PhD^{4,5,✉}

¹ Department of Visceral, Thoracic and Vascular Surgery, University Hospital and Faculty of Medicine Carl Gustav Carus, Technische Universität Dresden, Fetscherstraße 74, 01307 Dresden, Germany

² National Center for Tumor Diseases (NCT/UCC), Dresden, Germany; German Cancer Research Center (DKFZ), Heidelberg, Germany; Faculty of Medicine and University Hospital Carl Gustav Carus, Technische Universität Dresden, Dresden, Germany; Helmholtz-Zentrum Dresden-Rossendorf (HZDR), Fetscherstraße 74, 01307 Dresden, Germany

³ Else Kröner Fresenius Center for Digital Health (EKfZ), Technische Universität Dresden, Fetscherstraße 74, 01307 Dresden, Germany

⁴ Division of Translational Surgical Oncology, National Center for Tumor Diseases (NCT), Partner Site Dresden, Fetscherstraße 74, 01307 Dresden, Germany

⁵ Cluster of Excellence “Centre for Tactile Internet with Human-in-the-Loop” (CeTI), Technische Universität Dresden, 01062, Dresden, Germany

✉ Corresponding authors:

Dr. Fiona Kolbinger, Department of Visceral, Thoracic and Vascular Surgery, University Hospital and Faculty of Medicine Carl Gustav Carus, Technische Universität Dresden, Fetscherstraße 74, 01307 Dresden, Germany

☎ +49 (0) 351 458 19624

✉ fiona.kolbinger@uniklinikum-dresden.de

Dr. Sebastian Bodenstedt, Division of Translational Surgical Oncology, National Center for Tumor Diseases (NCT/UCC), Partner Site Dresden, Fetscherstrasse 74, 01307 Dresden, Germany

☎ +49 (0) 351 5413

✉ sebastian.bodenstedt@nct-dresden.de

STRUCTURED ABSTRACT

Background: Lack of anatomy recognition represents a clinically relevant risk in abdominal surgery. Machine learning (ML) methods can help identify visible patterns and risk structures, however, their practical value remains largely unclear.

Materials and Methods: Based on a novel dataset of 13195 laparoscopic images with pixel-wise segmentations of eleven anatomical structures, we developed specialized segmentation models for each structure and combined models for all anatomical structures using two state-of-the-art model architectures (DeepLabv3 and SegFormer), and compared segmentation performance of algorithms to a cohort of 28 physicians, medical students, and medical laypersons using the example of pancreas segmentation.

Results: Mean Intersection-over-Union for semantic segmentation of intraabdominal structures ranged from 0.28 to 0.83 and from 0.23 to 0.77 for the DeepLabv3-based structure-specific and

45 combined models, and from 0.31 to 0.85 and from 0.26 to 0.67 for the SegFormer-based structure-
46 specific and combined models, respectively. Both the structure-specific and the combined
47 DeepLabv3-based models are capable of near-real-time operation, while the SegFormer-based
48 models are not. All four models outperformed at least 26 out of 28 human participants in pancreas
49 segmentation.

50 **Conclusions:** These results demonstrate that ML methods have the potential to provide relevant
51 assistance in anatomy recognition in minimally-invasive surgery in near-real-time. Future research
52 should investigate the educational value and subsequent clinical impact of respective assistance
53 systems.

54

55 HIGHLIGHTS

- 56 • Machine learning models to reduce surgical risks that precisely identify 11 anatomical
57 structures: abdominal wall, colon, intestinal vessels (inferior mesenteric artery and inferior
58 mesenteric vein with their subsidiary vessels), liver, pancreas, small intestine, spleen,
59 stomach, ureter and vesicular glands
- 60 • Large training dataset of 13195 real-world laparoscopic images with high-quality anatomy
61 annotations
- 62 • Similar performance of individual segmentation models for each structure and combined
63 segmentation models in identifying intraabdominal structures, and similar segmentation
64 performance of DeepLabv3-based and SegFormer-based models
- 65 • DeepLabv3-based models are capable of near-real-time operation while SegFormer-
66 based models are not, but SegFormer-based models outperform DeepLabv3-based
67 models in terms of accuracy and generalizability
- 68 • All models outperformed at least 26 out of 28 human participants in pancreas
69 segmentation, demonstrating their potential for real-time assistance in recognizing
70 anatomical landmarks during minimally-invasive surgery.

71

72 KEYWORDS

73 Minimally-invasive surgery, laparoscopy, surgical data science, surgical anatomy, surgical
74 innovation, artificial intelligence

75

76 INTRODUCTION

77 Computer vision describes the computerized analysis of digital images aiming at the automation
78 of human visual capabilities, most commonly using machine learning methods, in particular deep
79 learning. This approach has transformed medicine in recent years, with successful applications
80 including computer-aided diagnosis of colonic polyp dignity in endoscopy (1,2), detection of
81 clinically actionable genetic alterations in histopathology (3), and melanoma detection in
82 dermatology (4). Availability of large amounts of training data is the defining prerequisite for
83 successful application of deep learning methods. With the establishment of laparoscopy as the
84 gold standard for a variety of surgical procedures (5–8) and the increasing availability of computing
85 resources, these concepts have gradually been applied to abdominal surgery. The overwhelming
86 majority of research efforts in the field of Artificial Intelligence (AI)-based analysis of intraoperative
87 surgical imaging data (i.e. video data from laparoscopic or open surgeries) has focused on
88 classifying images with respect to the presence and/or location of previously annotated surgical
89 instruments or anatomical structures (9–13) or on analysis of surgical proficiency (14–16) based
90 on recorded procedures. However, almost all research endeavors in the field of computer vision
91 in laparoscopic surgery have concentrated on preclinical stages and to date, no AI model based
92 on intraoperative surgical imaging data could demonstrate a palpable clinical benefit (17,18).
93 Among the studies closest to clinical application are recent works on identification of instruments
94 and hepatobiliary anatomy during cholecystectomy for automated assessment of the critical view
95 of safety (13), and on the automated segmentation of safe and unsafe preparation zones during
96 cholecystectomy (19).

97 In surgery, patient outcome heavily depends on experience and performance of the surgical team
98 (20,21). In a recent analysis of Human Performance Deficiencies in major cardiothoracic, vascular,
99 abdominal transplant, surgical oncology, acute care, and general surgical operations, more than
100 half of the cases with postoperative complications were associated with identifiable human error.
101 Among these errors, lack of recognition (including misidentified anatomy) accounted for 18.8%,
102 making it the most common Human Performance Deficiency overall (22). Examples of
103 complications directly related to anatomical misperception are iatrogenic lesions to the ureter in
104 gynecologic procedures (23) and pancreatic injuries during splenic flexure mobilization in
105 colorectal surgery (24). While AI-based systems identifying anatomical risk and target structures
106 would theoretically have the potential to alleviate this risk, limited availability and diversity of
107 (annotated) laparoscopic image data drastically restrict the clinical potential of such applications
108 in practice.

109 To advance and diversify the applications of computer vision in laparoscopic surgery, we have
110 recently published the Dresden Surgical Anatomy Dataset (25), providing 13195 laparoscopic
111 images with high-quality (26), expert-reviewed annotations of the presence and exact location of
112 eleven intraabdominal anatomical structures: abdominal wall, colon, intestinal vessels (inferior
113 mesenteric artery and inferior mesenteric vein with their subsidiary vessels), liver, pancreas, small
114 intestine, spleen, stomach, ureter and vesicular glands. Here, we present the first study based on
115 this dataset and present machine learning models to assist in precisely delineating anatomical
116 structures, aiming to reduce surgical risks. Specifically, we evaluate automated detection and
117 localization of organs and anatomical structures in laparoscopic view using two state-of-the-art
118 model architectures: DeepLabv3 and SegFormer. To assess the clinical value of the presented
119 machine learning models, we compare algorithm segmentation performance to that of humans
120 using the example of delineation of the pancreas.

121

122 **METHODS**

123 ***Patient cohort***

124 Video data from 32 robot-assisted anterior rectal resections or rectal extirpations were gathered
125 at the University Hospital Carl Gustav Carus Dresden between February 2019 and February 2021.
126 All included patients had a clinical indication for the surgical procedure, recommended by an
127 interdisciplinary tumor board. The procedures were performed using the da Vinci® Xi system
128 (Intuitive Surgical, Sunnyvale, CA, USA) with a standard Da Vinci® Xi/X Endoscope with Camera
129 (8 mm diameter, 30° angle, Intuitive Surgical, Sunnyvale, CA, USA, Item code 470057). Surgeries
130 were recorded using the CAST system (Orpheus Medical GmbH, Frankfurt a.M., Germany). Each
131 record was saved at a resolution of 1920 x 1080 pixels in MPEG-4 format.

132 All experiments were performed in accordance with the ethical standards of the Declaration of
133 Helsinki and its later amendments. The local Institutional Review Board (ethics committee at the
134 Technical University Dresden) reviewed and approved this study (approval number: BO-EK-
135 140032021). The trial was registered on clinicaltrials.gov (trial registration ID: NCT05268432).
136 Written informed consent to laparoscopic image data acquisition, data annotation, data analysis,
137 and anonymized data publication was obtained from all participants. Before publication, all data
138 was anonymized according to the general data protection regulation of the European Union.

139 ***Patient cohort***

140 Video data from 32 robot-assisted anterior rectal resections or rectal extirpations were gathered
141 at the University Hospital Carl Gustav Carus Dresden between February 2019 and February 2021.

142 All included patients had a clinical indication for the surgical procedure, recommended by an
143 interdisciplinary tumor board. Patients were not specifically selected with respect to demographic
144 or physical parameters (i.e. age, sex, body-mass index, comorbidities, previous surgical
145 procedures) or disease-specific criteria (i.e. indication, disease stage). Respective details of the
146 underlying patient cohort have been published previously (25). The procedures were performed
147 using the da Vinci® Xi system (Intuitive Surgical, Sunnyvale, CA, USA) with a standard Da Vinci®
148 Xi/X Endoscope with Camera (8 mm diameter, 30° angle, Intuitive Surgical, Sunnyvale, CA, USA,
149 Item code 470057). Surgeries were recorded using the CAST system (Orpheus Medical GmbH,
150 Frankfurt a.M., Germany). Each record was saved at a resolution of 1920 x 1080 pixels in MPEG-
151 4 format.

152 All experiments were performed in accordance with the ethical standards of the Declaration of
153 Helsinki and its later amendments. The local Institutional Review Board (ethics committee at the
154 Technical University Dresden) reviewed and approved this study (approval number: BO-EK-
155 140032021). The trial was registered on clinicaltrials.gov (trial registration ID: NCT05268432).
156 Written informed consent to laparoscopic image data acquisition, data annotation, data analysis,
157 and anonymized data publication was obtained from all participants. Before publication, all data
158 was anonymized according to the general data protection regulation of the European Union.

159 ***Dataset***

160 Based on the full-length surgery recordings and respective temporal annotations of organ visibility,
161 individual image frames were extracted and annotated as described previously (25). In brief, three
162 independent annotators with substantial experience in robot-assisted rectal surgery created pixel-
163 wise annotations, which were subsequently reviewed by a surgeon with 4 years of experience in
164 robot-assisted rectal surgery. A detailed description of the annotation process including underlying
165 annotation protocols as well as analyses of annotator agreement and technical parameters has
166 been published previously (25). To guarantee real-world applicability of machine learning models
167 trained on the dataset, images with perturbations such as blurring due to camera movements,
168 soiling of the lens, and presence of blood or smoke were not specifically excluded. However, the
169 annotation protocols advised annotators to only annotate structures in soiled and blurry images if
170 the respective structures were clearly delineable. The resulting Dresden Surgical Anatomy
171 Dataset comprises 13195 distinct images with pixel-wise segmentations of eleven anatomical
172 structures: abdominal wall, colon, intestinal vessels (inferior mesenteric artery and inferior
173 mesenteric vein with their subsidiary vessels), liver, pancreas, small intestine, spleen, stomach,
174 ureter and vesicular glands. Moreover, the dataset comprises binary annotations of the presence

175 of each of these organs for each image. The dataset is publicly available via the following link:
176 <https://doi.org/10.6084/m9.figshare.21702600>.

177
178 For machine learning purposes, the Dresden Surgical Anatomy Dataset was split into training,
179 validation, and test data as follows (Figure 1):

- 180 — Training set (at least 12 surgeries per anatomical structure): surgeries 1, 4, 5, 6, 8, 9, 10,
181 12, 15, 16, 17, 19, 22, 23, 24, 25, 27, 28, 29, 30, 31.
- 182 — Validation set (3 surgeries per anatomical structure): surgeries 3, 21, 26.
- 183 — Test set (5 surgeries per anatomical structure): surgeries 2, 7, 11, 13, 14, 18, 20, 32.

184 This split is proposed for future works using the Dresden Surgical Anatomy Dataset to reproduce
185 the variance of the entire dataset within each subset, and to ensure comparability regarding clinical
186 variables between the training, the validation, and the test set. Surgeries for the test set were
187 selected to minimize variance regarding the number of frames over the segmented classes. Out
188 of the remaining surgeries, the validation set was separated from the training set using the same
189 criterion.

190

191 ***Structure-specific semantic segmentation models***

192 To segment each anatomical structure, a separate convolutional neural network for segmentation
193 of individual structures was trained. Specifically, we trained and compared two different
194 architectures: a Deeplabv3 (27) model with a ResNet50 backbone with default PyTorch pretraining
195 on the COCO dataset (28), and a SegFormer (29) model pretrained on the Cityscapes dataset
196 (30). The networks were trained using cross-entropy loss and the AdamW optimizer (31) for 100
197 epochs with a starting learning rate of 10^{-4} and a linear learning rate scheduler decreasing the
198 learning rate by 0.9 every 10 epochs. For data augmentation, we applied random scaling and
199 rotation, as well as brightness and contrast adjustments. The final model for each organ was
200 selected via the Intersection-over-Union (IoU, Supplementary Figure 1) on the validation dataset
201 and evaluated using the Dresden Surgical Anatomy Dataset with the abovementioned training-
202 validation-test split (Figure 1).

203 Segmentation performance was assessed using F1 score, IoU, precision, recall, and specificity
204 on the test folds. These parameters are commonly used technical measures of prediction
205 exactness, ranging from 0 (least exact prediction) to 1 (entirely correct prediction without any
206 misprediction, Supplementary Figure 1).

207

208 ***Combined semantic segmentation models***

209 A convolutional neural network with a common encoder and eleven decoders for combined
210 segmentation of the eleven anatomical structures was trained. As for the structure-specific
211 models, DeepLabv3-based (27) and SegFormer-based (29) models were used. For DeepLabv3,
212 a shared ResNet50 backbone with default PyTorch pretraining on the COCO dataset (28) was
213 used. For each class, a DeepLabv3 decoder was then run on the features extracted from a given
214 image by the backbone. Similarly, for SegFormer, an encoder, pretrained on the Cityscapes
215 dataset (30), was combined with eleven decoders.

216 As the images are only annotated for binary classes, the loss is only calculated for every pixel in
217 images, in which the structure associated with the current decoder is annotated. For images, in
218 which the associated class is not annotated, only the pixels that are annotated as belonging to
219 another class are included in the loss, e.g., pixels that were annotated as belonging to the class
220 "liver" can be used as negative examples for the class "pancreas". The remaining training
221 procedure was identical to the structure-specific model. The models were trained and evaluated
222 using the Dresden Surgical Anatomy Dataset with the abovementioned training-validation-test
223 split (Figure 1).

224 Segmentation performance was assessed using F1 score, IoU, precision, recall, and specificity
225 on the test folds.

226

227 ***Evaluation of the semantic segmentation models on an external dataset***

228 To explore generalizability, structure-specific and combined models based on both architectures
229 (DeepLabv3 and SegFormer) were deployed to laparoscopic image data from the publicly
230 available LapGyn4 dataset (32). Models were separately deployed for full-scene segmentations
231 and their performance was visually compared.

232

233 ***Comparative evaluation of algorithmic and human performance***

234 To determine the clinical potential of automated segmentation of anatomical risk structures, the
235 segmentation performance of 28 humans was compared to that of the structure-specific and the
236 combined semantic segmentation models using the example of the pancreas. The local
237 Institutional Review Board (ethics committee at the Technical University Dresden) reviewed and
238 approved this study (approval number: BO-EK-566122021). All participants provided written
239 informed consent to anonymous study participation, data acquisition and analysis, and publication.
240 In total, 28 participants (physician and non-physician medical staff, medical students, and medical

241 laypersons) marked the pancreas in 35 images from the Dresden Surgical Anatomy Dataset (25)
242 with bounding boxes. These images originated from 26 different surgeries, and the pancreas was
243 visible in 16 of the 35 images. Each of the previously selected 35 images was shown once, the
244 order being arbitrarily chosen but identical for all participants. The open-source annotation
245 software Computer Vision Annotation Tool (CVAT) was used for annotations. In cases where the
246 pancreas was seen in multiple, non-connected locations in the image, participants were asked to
247 create separate bounding boxes for each area.

248 Based on the structure-specific and the combined semantic segmentation models, axis-aligned
249 bounding boxes marking the pancreas were generated in the 35 images from the pixel-wise
250 segmentation. To guarantee that the respective images were not part of the training data, four-
251 fold cross validation was used, i.e., the origin surgeries were split into four equal-sized batches,
252 and algorithms were trained on three batches that did not contain the respective origin image
253 before being applied to segmentation.

254 To compare human and algorithm performance, the bounding boxes created by each participant
255 and the structure-specific as well as the combined semantic segmentation models were compared
256 to bounding boxes derived from the Dresden Surgical Anatomy Dataset, which were defined as
257 ground truth. IoU between the manual or automatic bounding box and the ground truth was used
258 to compare segmentation accuracy.

259

260 **DATA AND CODE AVAILABILITY**

261 ***Data Availability***

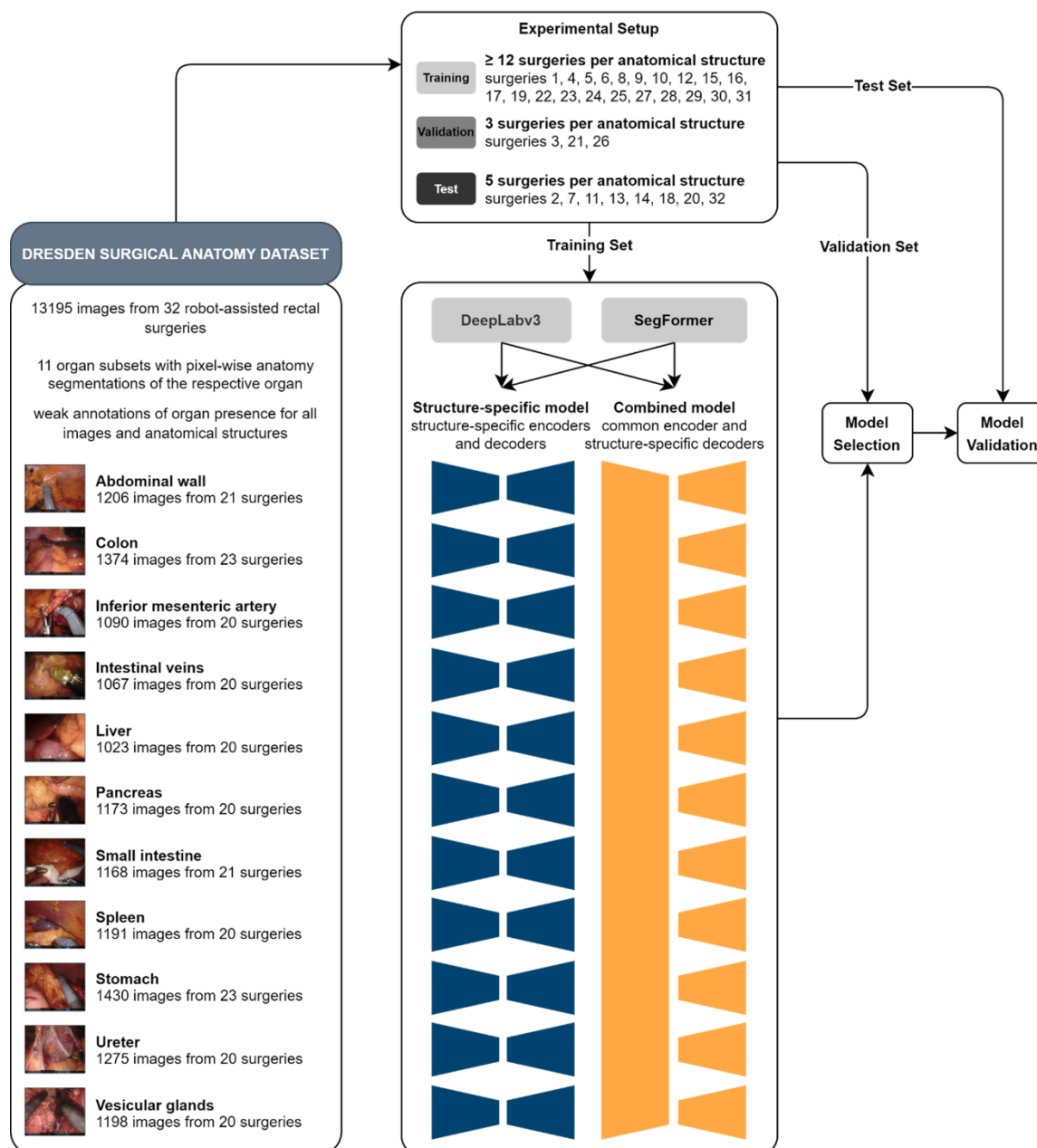
262 The Dresden Surgical Anatomy Dataset is publicly available via the following link:
263 <https://doi.org/10.6084/m9.figshare.21702600>. All other data generated and analyzed during the
264 current study are available from the corresponding authors on reasonable request. To gain
265 access, data requestors will need to sign a data access agreement.

266

267 ***Code Availability***

268 The most relevant scripts used for dataset compilation are publicly available via the following link:
269 <https://zenodo.org/record/6958337#.YzsBdnZBzOg>. The code used for segmentation algorithms
270 is available at https://gitlab.com/nct_tso_public/anatomy-recognition-dsad.

271



272
273
274
275
276
277
278
279

Figure 1: Schematic illustration of the structure-specific and combined machine learning models used for semantic segmentation. The Dresden Surgical Anatomy Dataset was split into a training, a validation, and a test set. For spatial segmentation, two sets of machine learning models – a structure-specific model with individual encoders and decoders, and a combined model with a common encoder and structure-specific decoders – were trained for DeepLabv3-based and SegFormer-based model architectures.

280 RESULTS

281 ***Machine Learning-based anatomical structure segmentation in structure-specific models***

282 Structure-specific multi-layer convolutional neural networks (Figure 1) based on two different
283 semantic segmentation architectures termed DeepLabv3 and SegFormer, were trained to
284 segment the abdominal wall, the colon, intestinal vessels (inferior mesenteric artery and inferior
285 mesenteric vein with their subsidiary vessels), the liver, the pancreas, the small intestine, the
286 spleen, the stomach, the ureter, and vesicular glands (Supplementary Table 1). Table 1 displays
287 technical metrics of overlap between the annotated ground truth and the model predictions (mean
288 F1 score, IoU, precision, recall, and specificity) for individual anatomical structures as predicted
289 by the structure-specific algorithms on the test data.

290 Out of the analyzed segmentation models based on DeepLabv3, performance was lowest for
291 vesicular glands (mean IoU: 0.28 ± 0.21), the pancreas (mean IoU: 0.28 ± 0.27), and the ureter
292 (mean IoU: 0.36 ± 0.25), while excellent predictions were achieved for the abdominal wall (mean
293 IoU: 0.83 ± 0.14) and the small intestine (mean IoU: 0.80 ± 0.18) (Supplementary Figure 1). In
294 segmentation of the pancreas, the ureter, vesicular glands and intestinal vessel structures, there
295 was a relevant proportion of images with no detection or no overlap between prediction and ground
296 truth, while for all remaining anatomical structures, this proportion was minimal (Figure 2 a). While
297 the images, in which the highest IoUs were observed, mostly displayed large organ segments that
298 were clearly visible (Figure 2 b), the images with the lowest IoU were of variable quality with
299 confounding factors such as blood, smoke, soiling of the endoscope lens, or pictures blurred by
300 camera shake (Figure 2 c). While overall segmentation performance of both architectures was
301 similar for structure-specific models, SegFormer-based models showed a trend towards better
302 performance than DeepLabv3-based models in segmentation of the pancreas, the spleen, and
303 the ureter (Table 1, Figure 2, Supplementary Figure 2).

304 To determine the models' capabilities to operate in real-time (frame rates of > 20 frames per
305 second), we determined their inference times per image. For the DeepLabv3-based structure-
306 specific models, inference on a single image with a resolution of 640×512 pixels required, on
307 average, 28 ms on an Nvidia A5000, resulting in a frame rate of almost 36 frames per second. In
308 contrast, the SegFormer-based structure-specific semantic segmentation models operated
309 considerably slower at an inference time of 53 ms per image, resulting in a frame rate of 18 frames
310 per second. This runtime includes one decoder, meaning that only the segmentation for one
311 anatomical class (organ or structure) is included.

312

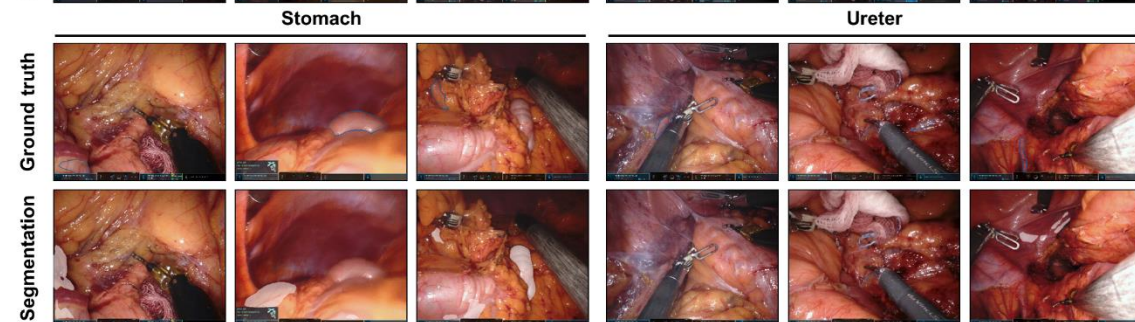
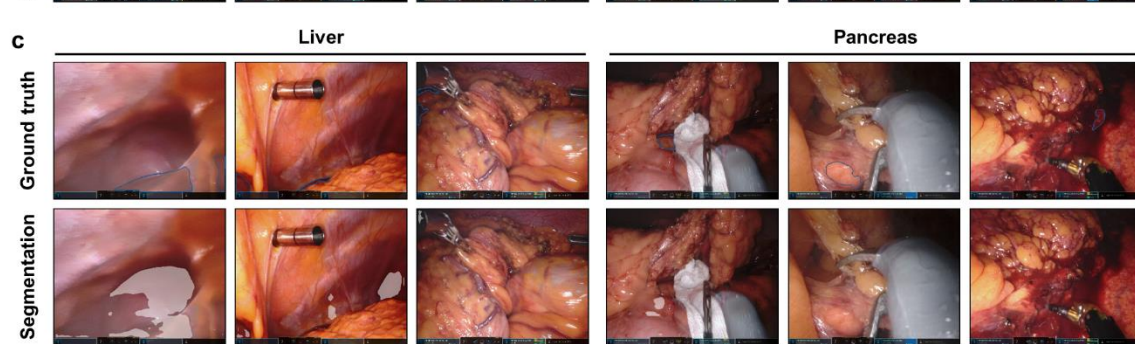
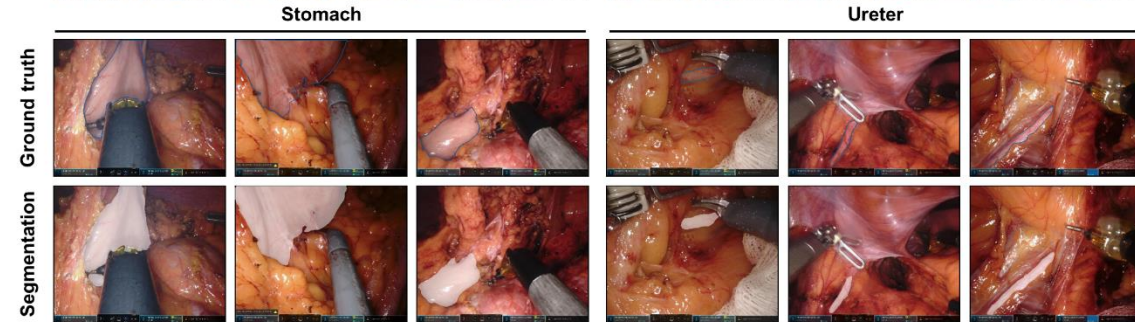
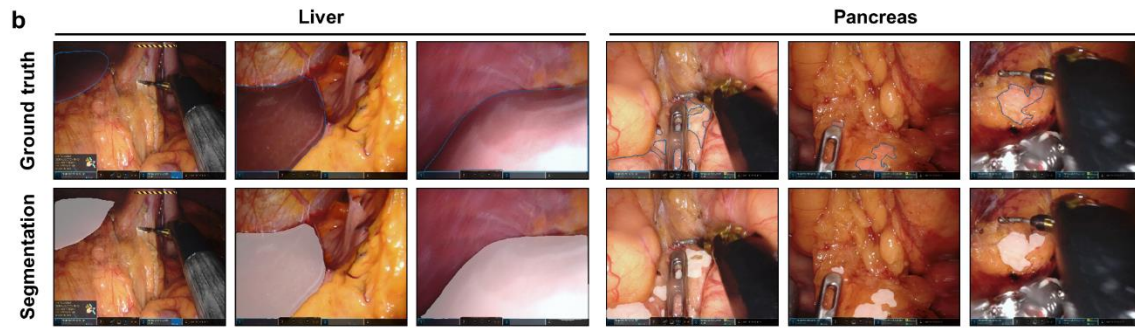
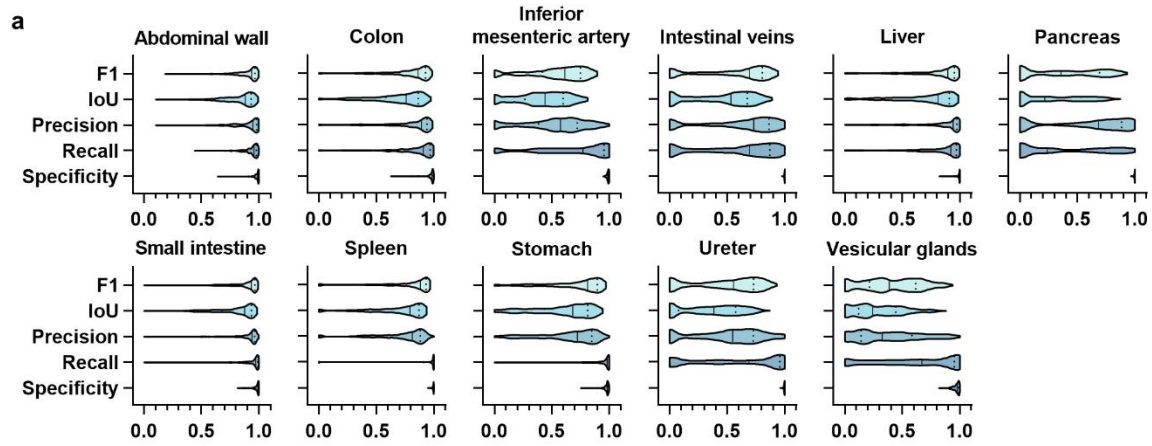
313 **Table 1: Summary of performance metrics for anatomical structure segmentation using DeepLabv3-based (a)**
 314 **and SegFormer-based (b) structure-specific models on the test dataset.** For each metric, mean and standard
 315 deviation are displayed.
 316

a	Anatomical structure	F1 score	IoU	Precision	Recall	Specificity
DeepLabv3 (Test)	Abdominal wall	0.90 ± 0.10	0.83 ± 0.14	0.89 ± 0.14	0.93 ± 0.07	0.97 ± 0.04
	Colon	0.79 ± 0.20	0.69 ± 0.22	0.80 ± 0.21	0.82 ± 0.21	0.97 ± 0.05
	Inferior mesenteric artery	0.54 ± 0.26	0.41 ± 0.22	0.55 ± 0.25	0.67 ± 0.33	0.99 ± 0.01
	Intestinal veins	0.54 ± 0.33	0.44 ± 0.29	0.70 ± 0.26	0.56 ± 0.36	1.00 ± 0.00
	Liver	0.80 ± 0.23	0.71 ± 0.25	0.85 ± 0.21	0.81 ± 0.24	0.98 ± 0.03
	Pancreas	0.37 ± 0.32	0.28 ± 0.27	0.59 ± 0.37	0.37 ± 0.36	1.00 ± 0.01
	Small intestine	0.87 ± 0.14	0.80 ± 0.18	0.87 ± 0.16	0.91 ± 0.15	0.97 ± 0.04
	Spleen	0.79 ± 0.23	0.69 ± 0.24	0.74 ± 0.22	0.90 ± 0.24	0.99 ± 0.01
	Stomach	0.71 ± 0.24	0.60 ± 0.25	0.65 ± 0.25	0.89 ± 0.21	0.98 ± 0.02
	Ureter	0.47 ± 0.30	0.36 ± 0.25	0.53 ± 0.28	0.57 ± 0.39	1.00 ± 0.00
Vesicular glands	0.40 ± 0.25	0.28 ± 0.21	0.37 ± 0.28	0.62 ± 0.35	0.97 ± 0.03	

b	Anatomical structure	F1 score	IoU	Precision	Recall	Specificity
SegFormer (Test)	Abdominal wall	0.91 ± 0.11	0.85 ± 0.15	0.90 ± 0.14	0.94 ± 0.09	0.98 ± 0.03
	Colon	0.77 ± 0.21	0.66 ± 0.22	0.73 ± 0.22	0.87 ± 0.22	0.95 ± 0.07
	Inferior mesenteric artery	0.60 ± 0.23	0.46 ± 0.21	0.58 ± 0.25	0.73 ± 0.29	0.99 ± 0.01
	Intestinal veins	0.65 ± 0.25	0.52 ± 0.24	0.62 ± 0.27	0.76 ± 0.27	1.00 ± 0.00
	Liver	0.83 ± 0.21	0.75 ± 0.24	0.82 ± 0.23	0.88 ± 0.18	0.98 ± 0.03
	Pancreas	0.47 ± 0.32	0.37 ± 0.28	0.61 ± 0.36	0.48 ± 0.36	0.99 ± 0.01
	Small intestine	0.89 ± 0.13	0.83 ± 0.17	0.87 ± 0.16	0.95 ± 0.10	0.97 ± 0.04
	Spleen	0.85 ± 0.19	0.78 ± 0.21	0.80 ± 0.19	0.95 ± 0.16	1.00 ± 0.01
	Stomach	0.75 ± 0.27	0.66 ± 0.28	0.76 ± 0.25	0.82 ± 0.29	0.99 ± 0.01
	Ureter	0.58 ± 0.27	0.46 ± 0.24	0.53 ± 0.26	0.74 ± 0.32	0.99 ± 0.01
Vesicular glands	0.43 ± 0.26	0.31 ± 0.22	0.40 ± 0.28	0.63 ± 0.35	0.97 ± 0.03	

317

318



320 **Figure 2: Pixel-wise organ segmentation with DeepLabv3-based structure-specific models trained on the**
321 **respective organ subsets of the Dresden Surgical Anatomy Dataset. (a)** Violin plot illustrations of performance
322 metrics for DeepLabv3-based structure-specific segmentation models on the test dataset. The median and quartiles are
323 illustrated as solid and dashed lines, respectively. **(b)** Example images from the test dataset with the highest IoUs for
324 liver, pancreas, stomach, and ureter segmentation with DeepLabv3-based structure-specific segmentation models.
325 Ground truth is displayed as blue line (upper panel), model segmentations are displayed as white overlay (lower panel).
326 **(c)** Example images from the test dataset with the lowest IoUs for liver, pancreas, stomach, and ureter segmentation
327 with DeepLabv3-based structure-specific segmentation models. Ground truth is displayed as blue line (upper panel),
328 model segmentations are displayed as white overlay (lower panel).

329

330 ***Machine Learning-based anatomical structure segmentation in combined models***

331 In contrast to structure-specific models, models with a mutual encoder and organ-specific
332 decoders could facilitate the identification of multiple organs at once, with the potential benefit of
333 faster operation for multiple classes instead of sequential operation of several class-specific
334 models. Therefore, combined models for both semantic segmentation architectures – DeepLabv3
335 and Segformer – were trained using annotated images from the Dresden Surgical Anatomy
336 Dataset across anatomical structure classes (Figure 1, Supplementary Table 2). Table 2 displays
337 mean F1 score, IoU, precision, recall, and specificity for anatomical structure segmentation in the
338 combined model.

339 The performance of the combined model based on DeepLabv3 was overall similar to that of
340 structure-specific models (Table 1), with highest segmentation performance for the abdominal wall
341 (mean IoU: 0.77 ± 0.15) and the small intestine (mean IoU: 0.72 ± 0.21), and the lowest
342 performance for the pancreas (mean IoU: 0.23 ± 0.29), the ureter (IoU: 0.29 ± 0.22) and vesicular
343 glands (IoU: 0.30 ± 0.23) (Supplementary Figure 1). In comparison to the respective structure-
344 specific models, the combined DeepLabv3-based model performed notably weaker in liver
345 segmentation, while performance for the other anatomical structures was similar. The proportion
346 of images for which the combined DeepLabv3-based model could not create a prediction or for
347 which predictions showed no overlap with the ground truth at all was largest in the ureter, the
348 pancreas, the stomach, the abdominal vessel structures, and the vesicular glands (Figure 3 a).
349 Similar to the DeepLabv3-based structure-specific models, trends towards an impact of segment
350 size, uncommon angles of vision, endoscope lens soiling, blurry images, and presence of blood
351 or smoke were seen when comparing image quality of well-predicted images (Figure 3 b) to
352 images with poor or no prediction (Figure 3 c). Similar to the structure-specific models,
353 segmentation performance of the SegFormer-based combined segmentation model was, overall,
354 similar to that of DeepLabv3-based models. For segmentation of the spleen, there was a trend
355 towards weaker performance of SegFormer-based models than DeepLabv3-based combined
356 models (Table 2, Figure 3, Supplementary Figure 3).

357 For the DeepLabv3-based combined models, inference on a single image with a resolution of 640
 358 x 512 pixels required, on average, 71 ms on an Nvidia A5000, resulting in a frame rate of about
 359 14 frames per second. As for structure-specific models of both architectures, SegFormer-based
 360 combined semantic segmentation models operated considerably slower at an inference time of
 361 102 ms per image, resulting in a frame rate of about 10 frames per second. This runtime includes
 362 all 11 decoders, meaning that segmentations for all anatomical classes (organs or structures) are
 363 included.

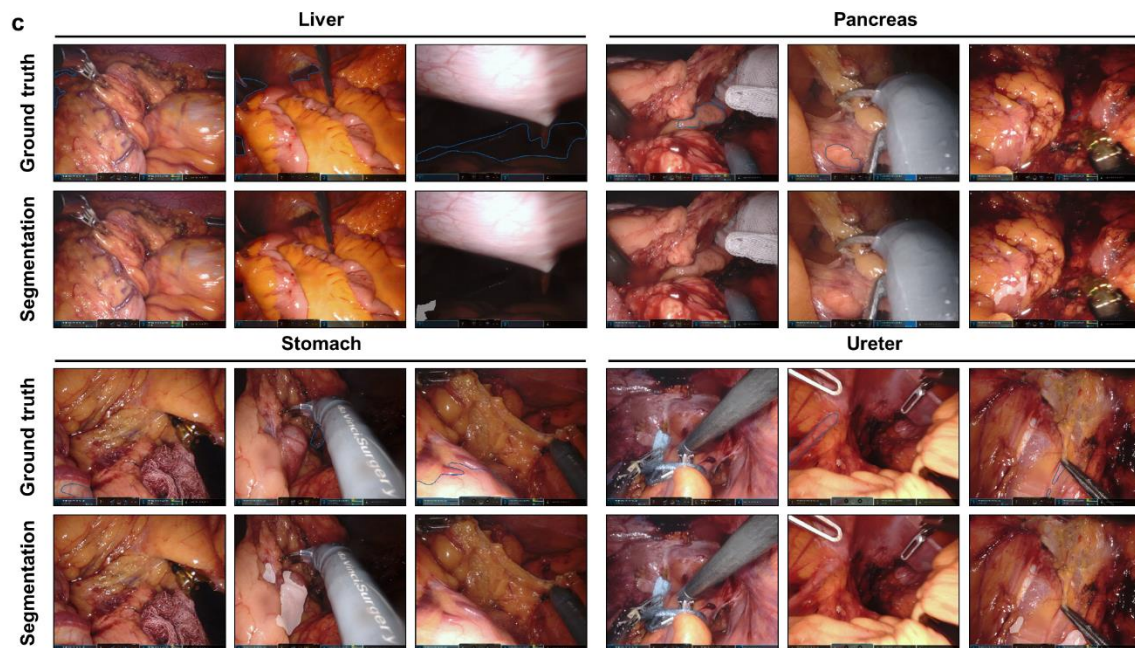
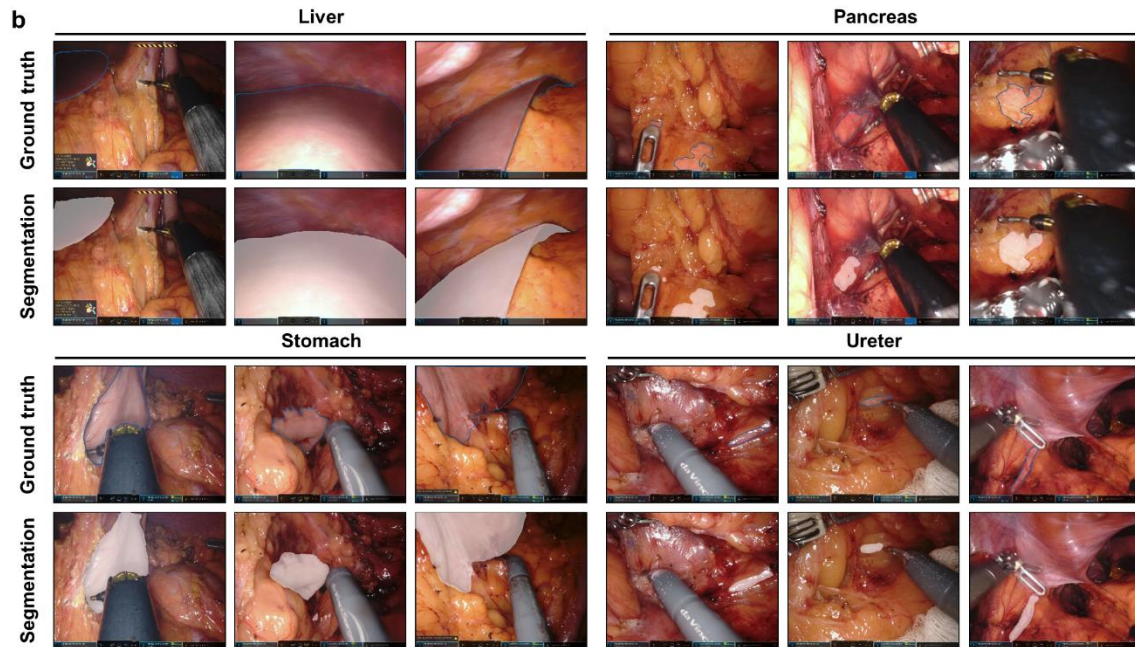
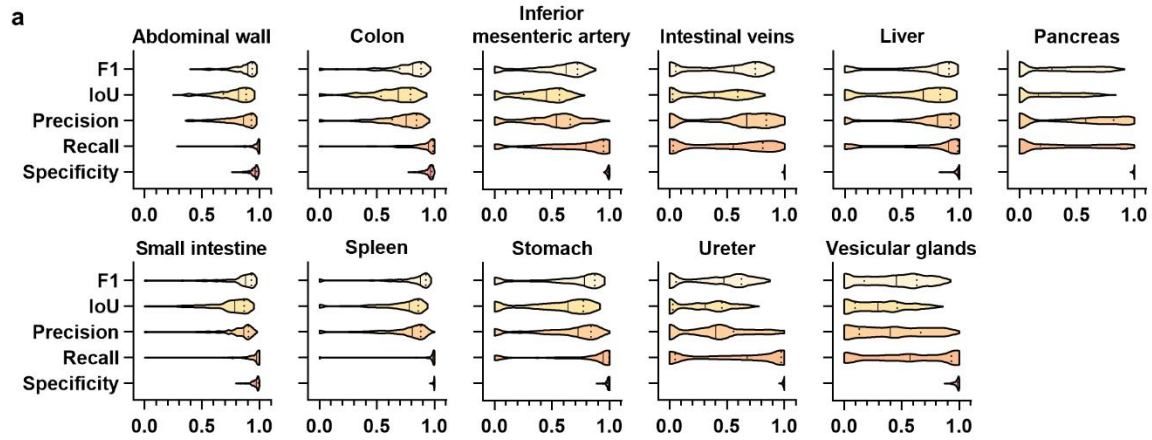
364

365 **Table 2: Summary of performance metrics for anatomical structure segmentation using the DeepLabv3-based**
 366 **(a) and SegFormer-based (b) combined models (common encoder with structure-specific decoders) on the test**
 367 **dataset.** For each metric, mean and standard deviation are displayed.
 368

a	Anatomical structure	F1 score	IoU	Precision	Recall	Specificity
DeepLabv3 (Test)	Abdominal wall	0.86 ± 0.11	0.77 ± 0.15	0.81 ± 0.15	0.95 ± 0.09	0.95 ± 0.04
	Colon	0.75 ± 0.19	0.63 ± 0.21	0.71 ± 0.18	0.84 ± 0.23	0.95 ± 0.04
	Inferior mesenteric artery	0.53 ± 0.25	0.40 ± 0.21	0.52 ± 0.22	0.68 ± 0.32	0.99 ± 0.01
	Intestinal veins	0.46 ± 0.32	0.35 ± 0.27	0.70 ± 0.23	0.48 ± 0.36	1.00 ± 0.00
	Liver	0.65 ± 0.34	0.57 ± 0.33	0.76 ± 0.23	0.69 ± 0.38	0.98 ± 0.03
	Pancreas	0.32 ± 0.30	0.23 ± 0.24	0.61 ± 0.33	0.32 ± 0.35	0.99 ± 0.01
	Small intestine	0.81 ± 0.19	0.72 ± 0.21	0.81 ± 0.17	0.87 ± 0.23	0.96 ± 0.03
	Spleen	0.78 ± 0.24	0.69 ± 0.24	0.76 ± 0.18	0.89 ± 0.26	0.99 ± 0.01
	Stomach	0.63 ± 0.32	0.53 ± 0.29	0.68 ± 0.23	0.74 ± 0.37	0.98 ± 0.02
	Ureter	0.40 ± 0.28	0.29 ± 0.22	0.44 ± 0.27	0.56 ± 0.40	0.99 ± 0.01
	Vesicular glands	0.42 ± 0.27	0.30 ± 0.23	0.41 ± 0.30	0.56 ± 0.36	0.98 ± 0.02
b	Anatomical structure	F1 score	IoU	Precision	Recall	Specificity
SegFormer (Test)	Abdominal wall	0.76 ± 0.24	0.66 ± 0.24	0.78 ± 0.15	0.86 ± 0.28	0.92 ± 0.07
	Colon	0.64 ± 0.27	0.52 ± 0.24	0.66 ± 0.19	0.78 ± 0.34	0.93 ± 0.06
	Inferior mesenteric artery	0.40 ± 0.24	0.28 ± 0.18	0.37 ± 0.21	0.68 ± 0.38	0.97 ± 0.02
	Intestinal veins	0.43 ± 0.33	0.33 ± 0.27	0.63 ± 0.24	0.52 ± 0.41	0.99 ± 0.01
	Liver	0.62 ± 0.35	0.53 ± 0.32	0.76 ± 0.19	0.71 ± 0.40	0.95 ± 0.07
	Pancreas	0.35 ± 0.32	0.26 ± 0.25	0.56 ± 0.28	0.43 ± 0.41	0.99 ± 0.01
	Small intestine	0.78 ± 0.19	0.67 ± 0.20	0.74 ± 0.14	0.90 ± 0.23	0.94 ± 0.06
	Spleen	0.71 ± 0.24	0.59 ± 0.23	0.65 ± 0.21	0.89 ± 0.25	0.99 ± 0.01
	Stomach	0.65 ± 0.32	0.55 ± 0.29	0.71 ± 0.22	0.75 ± 0.36	0.98 ± 0.02
	Ureter	0.38 ± 0.29	0.27 ± 0.23	0.43 ± 0.27	0.55 ± 0.40	0.99 ± 0.01
	Vesicular glands	0.38 ± 0.25	0.26 ± 0.20	0.32 ± 0.25	0.66 ± 0.35	0.96 ± 0.03

369

370



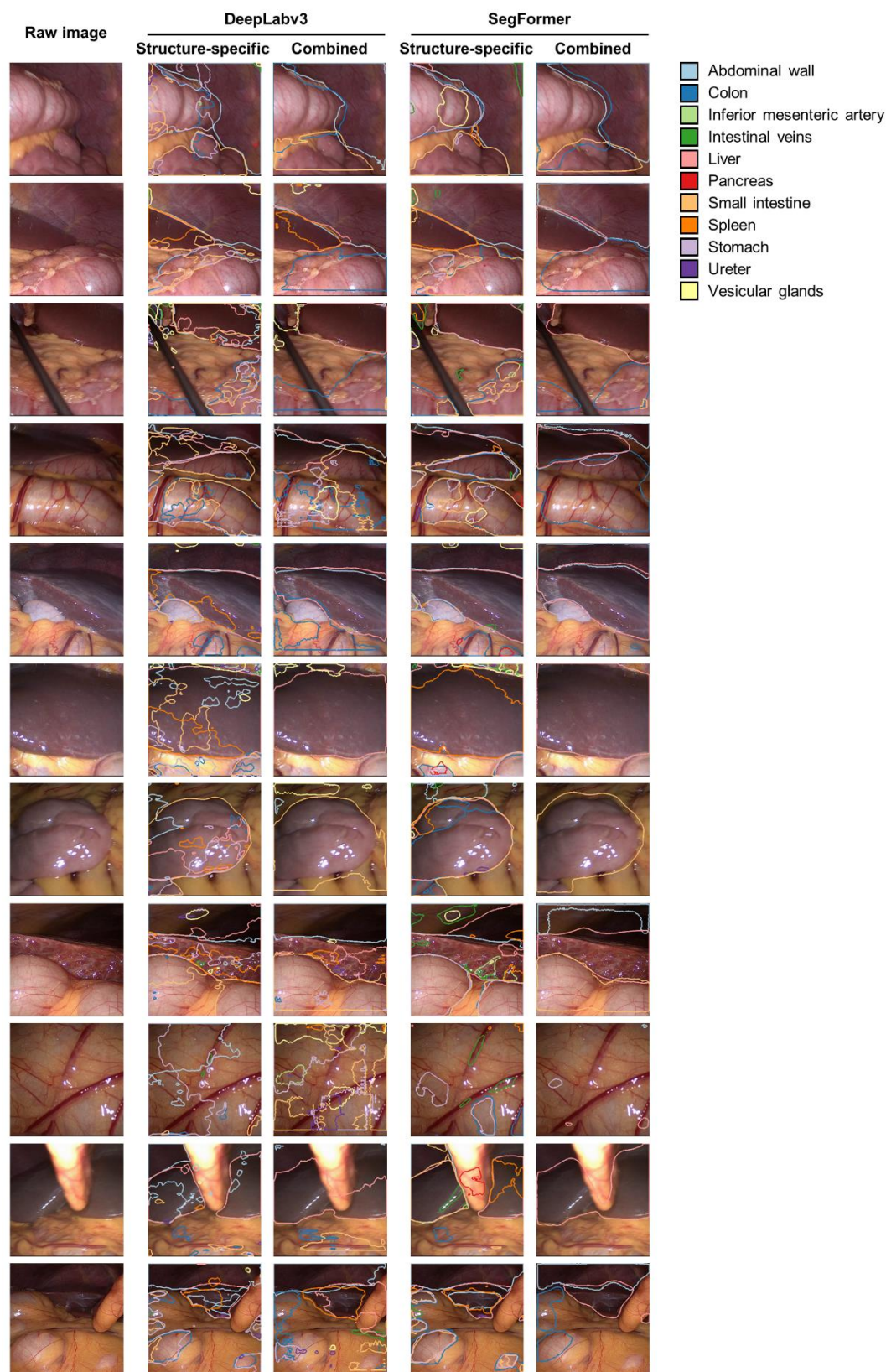
372 **Figure 3: Pixel-wise organ segmentation with the DeepLabv3-based combined model trained on the Dresden**
373 **Surgical Anatomy Dataset across anatomical structure classes with a common encoder and structure-specific**
374 **decoders. (a)** Violin plot illustrations of performance metrics for the DeepLabv3-based combined segmentation model
375 on the test dataset. The median and quartiles are illustrated as solid and dashed lines, respectively. **(b)** Example images
376 from the test dataset with the highest IoUs for liver, pancreas, stomach, and ureter segmentation with the DeepLabv3-
377 based combined segmentation model. Ground truth is displayed as blue line (upper panel), model segmentations are
378 displayed as white overlay (lower panel). **(c)** Example images from the test dataset with the lowest IoUs for liver,
379 pancreas, stomach, and ureter segmentation with the DeepLabv3-based combined segmentation model. Ground truth
380 is displayed as blue line (upper panel), model segmentations are displayed as white overlay (lower panel).
381

382 ***Performance of machine learning models on an external laparoscopic image dataset***

383 To evaluate model robustness on an external dataset, we deployed the different organ
384 segmentation models onto the publicly available LapGyn4 dataset (32) and qualitatively compared
385 their performance. Overall, the combined models better reflected true anatomical constellations
386 than the structure-specific models that generally lacked specificity. With respect to model
387 architecture, the SegFormer-based segmentations were considerably more robust than the
388 DeepLabv3-based models. Common mispredictions included confusion of liver and spleen,
389 misinterpretation of organs that were not part of the training dataset (i.e. the gallbladder), and poor
390 segmentation performance on less common images (i.e. extreme close-ups) (Figure 4).

391 In summary, the SegFormer-based combined semantic segmentation model resulted in robust
392 segmentations reproducing the true underlying anatomy. The remaining segmentation models
393 provided substantially less specific and less robust segmentation outputs on the external dataset.

394



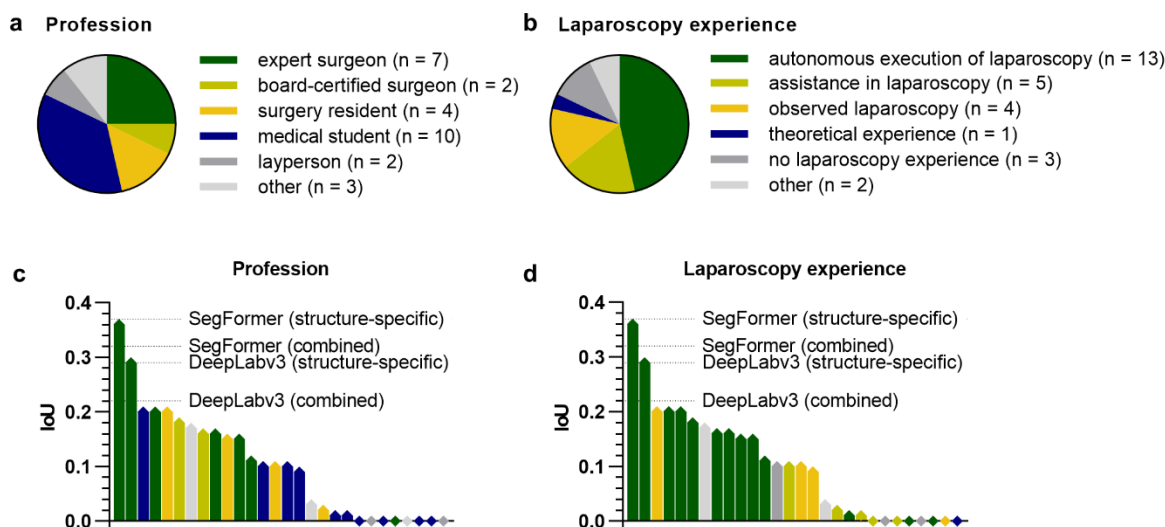
395
396 **Figure 4: Comparison of DeepLabv3-based and Segformer-based structure-specific and combined**
397 **segmentation model performance on an external laparoscopic image dataset (LapGyn4).** Models were deployed
398 to the publicly available LapGyn4 dataset of non-semantically segmented images from gynecological procedures in

399 conventional laparoscopic technique. Model segmentations for each organ are displayed. For the structure-specific
400 models, segmentations of the eleven individual segmentation models are overlaid in one image. Figure shows
401 representative images from the dataset.
402

403 **Performance of machine learning models in relation to human performance**

404 To approximate the clinical value of the previously described algorithms for anatomical structure
405 segmentation, the performances of the DeepLabv3-based and SegFormer-based structure-
406 specific and the combined models were compared to that of a cohort of 28 physicians, medical
407 students, and persons with no medical background (Figure 5 a), and different degrees of
408 experience in laparoscopic surgery (Figure 5 b). A vulnerable anatomical structure (24) with –
409 measured by classical metrics of overlap (Tables 1 and 2) – comparably weak segmentation
410 performance of the trained algorithms, the pancreas was selected as an example.

411 Comparing bounding box segmentations of the pancreas of human annotators, the medical and
412 laparoscopy-specific experience of participants was mirrored by the respective IoUs describing
413 the overlap between the pancreas annotation and the ground truth. The pancreas-specific
414 segmentation models based on DeepLabv3 (IoU: 0.29) and SegFormer (IoU: 0.37) as well as the
415 combined segmentation models based on DeepLabv3 (IoU: 0.21) and SegFormer (IoU: 0.32)
416 outperformed at least 26 out of the 28 human participants (Figures 5 c and d). Overall, these
417 results demonstrate that the developed models have clinical potential to improve the recognition
418 of vulnerable anatomical structures in laparoscopy.
419



420
421 **Figure 5: Comparison of pancreas segmentation performance of the structure-specific and the combined**
422 **semantic segmentation models with a cohort of 28 human participants. (a)** Distribution of medical and non-medical
423 **professions among human participants. (b)** Distribution of laparoscopy experience among human participants. **(c)**
424 **Waterfall chart displaying the average pancreas segmentation IoUs of participants with different professions as**
425 **compared to the IoU generated by the structure-specific and the combined semantic segmentation models. (d)**
426 **Waterfall chart displaying the average pancreas segmentation IoUs of participants with varying laparoscopy experience as**
427 **compared to the IoU generated by the structure-specific and the combined semantic segmentation models.**

428 **DISCUSSION**

429 In surgery, misinterpretation of visual cues can result in objectifiable errors with serious
430 consequences (22). Machine learning models could augment identification of anatomical
431 structures during minimally-invasive surgery and thereby contribute to a reduction of surgical risks.
432 However, data scarcity and suboptimal dataset quality, among other factors, drastically restrict the
433 clinical impact of applications in the field of surgical data science (17,33–37). Based on a robust
434 public dataset providing 13195 laparoscopic images with segmentations of eleven intra-abdominal
435 anatomical structures, this study explores the potential of machine learning for automated
436 segmentation of these organs, and compares algorithmic segmentation quality to that of humans
437 with varying experience in minimally-invasive abdominal surgery.

438 In summary, the presented findings suggest that machine learning-based segmentation of
439 intraabdominal organs and anatomical structures is possible and has the potential to provide
440 clinically valuable information. At an average runtime of 71 ms per image, corresponding to a
441 frame rate of 14 frames per second, the combined DeepLabv3-based model would facilitate near-
442 real-time identification of eleven anatomical structures. In contrast, the SegFormer-based model
443 is further from real-time performance at a runtime of 102 ms per image, resulting in a frame rate
444 of less than 10 frames per second. These runtimes mirror the performances of non-optimized
445 versions of the models, which can be significantly improved using methods such as TensorRT
446 from Nvidia. However, with respect to generalizability and robustness, we observed substantially
447 more accurate segmentation performance of the SegFormer-based models as compared to the
448 DeepLabv3-based models when deployed to an external conventional laparoscopic dataset.
449 Moreover, the structure-specific models exhibited a lack in accuracy and anatomical coherence,
450 which can be explained by their organ-specific training process.

451 Measured by classical metrics of overlap between segmentation and ground truth, predictions
452 were, overall, better for large and similar-appearing organs such as the abdominal wall, the liver,
453 the stomach, and the spleen as compared to smaller and more diverse-appearing organs such as
454 the pancreas, the ureter, or vesicular glands. Furthermore, poor image quality (i.e., images blurred
455 by camera movements, presence of blood or smoke in images) was linked to lower accuracy of
456 machine learning-based segmentations. Consequently, it is likely that a better nominal
457 performance of the machine learning models could be achieved through selection of images from
458 early phases of the surgery, in which such perturbations are not present. However, we purposely
459 did not exclude images with suboptimal image quality, as selection on image level would introduce
460 bias and thereby limit applicability. In this context, selection on patient level and on image level is
461 a common challenge in computer vision (38) that can lead to skewed reporting of outcomes and

462 poor performance on real-world data (34). Overall, our findings on the influence of image quality
463 on segmentation performance imply that computer vision studies in laparoscopy should be
464 carefully interpreted taking representativity and potential selection of underlying training and
465 validation data into consideration.

466 Measured by classical metrics of overlap (e.g., IoU, F1 score, precision, recall, specificity) that are
467 commonly used to evaluate segmentation performance, the structure-specific models and the
468 combined models provided comparable segmentation performances on the internal test dataset.
469 Interpretation of such metrics of overlap, however, represents a major challenge in computer
470 vision applications in medical domains such as dermatology and endoscopy (39–41) as well as
471 non-medical domains such as autonomous driving (42). In the specific use case of laparoscopic
472 surgery, evidence suggests that such technical metrics alone are not sufficient to characterize the
473 clinical potential and utility of segmentation algorithms (37,43). In this context, the subjective
474 clinical utility of a bounding box-based detection system recognizing the common bile duct and
475 the cystic duct at average precisions of 0.32 and 0.07, respectively, demonstrated by Tokuyasu
476 *et al.*, supports this hypothesis (12). In colorectal surgery, anatomical misinterpretation during
477 splenic flexure mobilization can result in iatrogenic lesions to the pancreas (24). In the presented
478 analysis, the trained structure-specific and combined machine learning algorithms outperformed
479 all human participants in the specific task of bounding box segmentation of the pancreas except
480 for two surgical specialists with over 10 years of experience. This suggests that even for structures
481 such as the pancreas with seemingly poor segmentation quality (segmentation IoU of the best-
482 performing model: 0.37 ± 0.28 in the test set) have the potential to provide clinically valuable help
483 in anatomy recognition. In this context, analysis of additional anatomical risk structures (i.e. ureters
484 and blood vessels) and inclusion of more advanced personnel in future comparison studies will
485 help better define the models' capabilities in comparison with (expert) surgeons. Notably, the best
486 average IoUs for pancreas segmentation achieved in this comparative study were 0.37 (for the
487 SegFormer-based structure-specific model) and 0.36 (for the best human participant), which
488 would both be considered less reliable segmentation quality measures on paper. This encourages
489 further discussion about metrics for segmentation quality assessment in clinical AI. In the future,
490 the potential of the described dataset (25) and organ segmentation algorithms could be exploited
491 for educational purposes (44,45), for guidance systems facilitating real-time detection of risk and
492 target structures (19,43,46,47), or as an auxiliary function integrated in more complex surgical
493 assistance systems, such as guidance systems relying on automated liver registration (48).

494 The limitations of this work are mostly related to the dataset and general limitations of machine
495 learning-based segmentation: First, the Dresden Surgical Anatomy Dataset is a monocentric

496 dataset based on 32 robot-assisted rectal surgeries. Therefore, the images used for algorithm
497 training and validation originate from one set of hardware and display organs from specific angles.
498 As a consequence, given the lack of a laparoscopic image dataset with similarly rigorous organ
499 annotations, generalizability and transferability of the presented findings to other centers and other
500 minimally-invasive abdominal surgeries, particularly non-robotic procedures, could only be
501 qualitatively investigated. Second, annotations were required for training of machine learning
502 algorithms, potentially inducing some bias towards the way that organs were annotated in the
503 resulting models. With respect to annotation quality, three individual annotations of each
504 anatomical structure were reviewed by a single surgical expert. This represents a major limitation
505 of the underlying dataset, which is reasoned in the time-consuming and effortful annotation
506 process making the inclusion of more expert surgeons unfeasible. Given that annotations were
507 based on specific annotation protocols including images (49) and all annotators had a medical
508 background with several years of experience in the field of human anatomy (25), the quality of
509 annotations can be considered high, despite the limited experience of the reviewing surgeon (4
510 years of experience in robot-assisted rectal surgery). This is particularly true when comparing the
511 underlying dataset with other datasets commonly used in surgical data science that are often
512 based on single annotations carried out by individuals without domain knowledge (17,26,50). Still,
513 the way that organs are annotated may differ from individual healthcare professionals' way of
514 recognizing an organ. This is particularly relevant for organs such as the ureters or the pancreas,
515 which often appear covered by layers of tissue. Here, computer vision-based algorithms that solely
516 consider the laparoscopic images provided by the Dresden Surgical Anatomy Dataset for
517 identification of risk structures will only be able to identify an organ once it is visible. For an earlier
518 recognition of such hidden risk structures, more training data with meaningful annotations would
519 be necessary. Importantly, the presented comparison to human performance focused on
520 segmentation of visible anatomy as well, neglecting that humans (and possibly computers, too)
521 could already identify a risk structure hidden underneath layers of tissue. Third, the dataset only
522 includes individual annotated images. In some structures such as the ureter, video data offers
523 considerably more information than still image data. In this context, it is conceivable that an
524 incorporation of temporal aspects could result in major improvements of both human and algorithm
525 recognition performance.

526 While the presented machine learning models show promise in improving the identification of
527 anatomical structures in laparoscopy, their clinical utility still needs to be explored. Successful
528 adoption of new technologies in surgery depends on factors beyond segmentation performance,
529 runtime and generalizability, such as visualization of intraoperative decision support (51), human-

530 machine interaction (52), and interface design. Therefore, interdisciplinary collaboration is critical
531 to better understand respective surgeon needs. Moreover, prospective trials are needed to
532 determine the impact of these factors on clinical outcomes. The existing limitations
533 notwithstanding, the presented study represents an important addition to the growing body of
534 research on medical image analysis in laparoscopic surgery, particularly by linking technical
535 metrics to human performance.

536 In conclusion, this study demonstrates that machine learning methods have the potential to
537 provide clinically relevant near-real-time assistance in anatomy recognition in minimally-invasive
538 surgery. This study is the first to use the recently published Dresden Surgical Anatomy Dataset,
539 providing baseline algorithms for organ segmentation and evaluating the clinical relevance of such
540 algorithms by introducing more clinically meaningful comparators beyond classical computer
541 vision metrics. Future research should investigate other segmentation methods, the potential to
542 integrate high-level anatomical knowledge into segmentation models (38), the transferability of
543 these results to other surgical procedures, and the clinical impact of real-time surgical assistance
544 systems and didactic applications based on automated segmentation algorithms. Furthermore,
545 seeing the DeepLabv3-based models outperform the SegFormer-based models in terms of run-
546 time, but lacking in accuracy and generalizability, future research could focus on combining the
547 two, in order to harness the best of both worlds.

548

549 **ACKNOWLEDGMENTS**

550 ***Assistance with the study***

551 The authors gratefully acknowledge excellent project coordination by Dr. Elisabeth Fischermeier
552 and Dr. Grit Krause-Jüttler.

553

554 ***Financial support and sponsorship***

555 This work has been funded by the Else Kröner Fresenius Center for Digital Health (EKFZ),
556 Dresden, Germany (project “CoBot”), by the German Research Foundation DFG within the Cluster
557 of Excellence EXC 2050: “Center for Tactile Internet with Human-in-the-Loop (CeTI)” (project
558 number 390696704) and by the German Federal Ministry of Health (BMG) within the “Surgomics”
559 project (grant number BMG 2520DAT82). Furthermore, FRK received funding from the Medical
560 Faculty of the Technical University Dresden within the MedDrive Start program (grant number
561 60487) and from the Joachim Herz Foundation (Add-On Fellowship for Interdisciplinary Life
562 Science). FMR received a doctoral student scholarship from the *Carus Promotionskolleg* Dresden.

563

564 ***Conflicts of interest***

565 The authors declare no conflicts of interest.

566

567 ***Presentation***

568 None.

569

570 REFERENCES

- 571 1. Wang P, Liu X, Berzin TM, Glissen Brown JR, Liu P, Zhou C, et al. Effect of a deep-learning
572 computer-aided detection system on adenoma detection during colonoscopy (CADE-DB
573 trial): a double-blind randomised study. *Lancet Gastroenterol Hepatol* [Internet]. 2020 Apr
574 1 [cited 2022 Oct 3];5(4):343–51. Available from:
575 <http://www.thelancet.com/article/S246812531930411X/fulltext>
- 576 2. Wang P, Berzin TM, Glissen Brown JR, Bharadwaj S, Becq A, Xiao X, et al. Real-time
577 automatic detection system increases colonoscopic polyp and adenoma detection rates: a
578 prospective randomised controlled study. *Gut* [Internet]. 2019 Oct 1 [cited 2022 Oct
579 3];68(10):1813–9. Available from: <https://gut.bmj.com/content/68/10/1813>
- 580 3. Kather JN, Heij LR, Grabsch HI, Loeffler C, Echle A, Muti HS, et al. Pan-cancer image-
581 based detection of clinically actionable genetic alterations. *Nat Cancer* 2020 18 [Internet].
582 2020 Jul 27 [cited 2022 Oct 3];1(8):789–99. Available from:
583 <https://www.nature.com/articles/s43018-020-0087-6>
- 584 4. Esteva A, Kuprel B, Novoa RA, Ko J, Swetter SM, Blau HM, et al. Dermatologist-level
585 classification of skin cancer with deep neural networks. *Nature* [Internet]. 2017 Feb 2 [cited
586 2021 Feb 15];542(7639):115–8. Available from:
587 <https://pubmed.ncbi.nlm.nih.gov/28117445/>
- 588 5. Simillis C, Lal N, Thoukididou SN, Kontovounisios C, Smith JJ, Hompes R, et al. Open
589 Versus Laparoscopic Versus Robotic Versus Transanal Mesorectal Excision for Rectal
590 Cancer: A Systematic Review and Network Meta-analysis. *Ann Surg* [Internet]. 2019 Jul 1
591 [cited 2022 Oct 6];270(1):59–68. Available from:
592 <https://pubmed.ncbi.nlm.nih.gov/30720507/>
- 593 6. Zhao JJ, Syn NL, Chong C, Tan HL, Ng JYX, Yap A, et al. Comparative outcomes of
594 needlescopic, single-incision laparoscopic, standard laparoscopic, mini-laparotomy, and
595 open cholecystectomy: A systematic review and network meta-analysis of 96 randomized
596 controlled trials with 11,083 patients. *Surgery* [Internet]. 2021 Oct 1 [cited 2022 Oct
597 6];170(4):994–1003. Available from: <https://pubmed.ncbi.nlm.nih.gov/34023139/>
- 598 7. Luketich JD, Pennathur A, Awais O, Levy RM, Keeley S, Shende M, et al. Outcomes after
599 minimally invasive esophagectomy: review of over 1000 patients. *Ann Surg* [Internet]. 2012
600 Jul [cited 2022 Oct 6];256(1):95–103. Available from:
601 <https://pubmed.ncbi.nlm.nih.gov/22668811/>
- 602 8. Thomson JE, Kruger D, Jann-Kruger C, Kiss A, Omoshoro-Jones JAO, Luvhengo T, et al.
603 Laparoscopic versus open surgery for complicated appendicitis: a randomized controlled

- 604 trial to prove safety. *Surg Endosc* [Internet]. 2015 Jul 19 [cited 2022 Oct 6];29(7):2027–32.
605 Available from: <https://pubmed.ncbi.nlm.nih.gov/25318368/>
- 606 9. Islam M, Atputharuban DA, Ramesh R, Ren H. Real-time instrument segmentation in
607 robotic surgery using auxiliary supervised deep adversarial learning. *IEEE Robot Autom*
608 *Lett.* 2019 Apr 1;4(2):2188–95.
- 609 10. Roß T, Reinke A, Full PM, Wagner M, Kenngott H, Apitz M, et al. Comparative validation
610 of multi-instance instrument segmentation in endoscopy: results of the ROBUST-MIS 2019
611 challenge. *Med Image Anal* [Internet]. 2020 Nov [cited 2021 Mar 21];70:101920. Available
612 from: <https://pubmed.ncbi.nlm.nih.gov/33676097/>
- 613 11. Shvets AA, Rakhlin A, Kalinin AA, Iglovikov VI. Automatic Instrument Segmentation in
614 Robot-Assisted Surgery using Deep Learning. In: *Proceedings - 17th IEEE International*
615 *Conference on Machine Learning and Applications, ICMLA 2018.* Institute of Electrical and
616 Electronics Engineers Inc.; 2019. p. 624–8.
- 617 12. Tokuyasu T, Iwashita Y, Matsunobu Y, Kamiyama T, Ishikake M, Sakaguchi S, et al.
618 Development of an artificial intelligence system using deep learning to indicate anatomical
619 landmarks during laparoscopic cholecystectomy. *Surg Endosc* 2020 354 [Internet]. 2020
620 Apr 18 [cited 2021 Jul 16];35(4):1651–8. Available from:
621 <https://link.springer.com/article/10.1007/s00464-020-07548-x>
- 622 13. Mascagni P, Vardazaryan A, Alapatt D, Urade T, Emre T, Fiorillo C, et al. Artificial
623 Intelligence for Surgical Safety: Automatic Assessment of the Critical View of Safety in
624 Laparoscopic Cholecystectomy Using Deep Learning. *Ann Surg.* 2020;
- 625 14. Jin A, Yeung S, Jopling J, Krause J, Azagury D, Milstein A, et al. Tool Detection and
626 Operative Skill Assessment in Surgical Videos Using Region-Based Convolutional Neural
627 Networks. *Proc - 2018 IEEE Winter Conf Appl Comput Vision, WACV 2018* [Internet]. 2018
628 Feb 24 [cited 2021 Jul 14];2018-January:691–9. Available from:
629 <https://arxiv.org/abs/1802.08774v2>
- 630 15. Funke I, Bodenstedt S, Oehme F, von Bechtolsheim F, Weitz J, Speidel S. Using 3D
631 Convolutional Neural Networks to Learn Spatiotemporal Features for Automatic Surgical
632 Gesture Recognition in Video. *Med Image Comput Comput Assist Interv – MICCAI 2019*
633 *Lect Notes Comput Sci* [Internet]. 2019 [cited 2022 Aug 19];11768:467–75. Available from:
634 https://link.springer.com/chapter/10.1007/978-3-030-32254-0_52
- 635 16. Lavanchy JL, Zindel J, Kirtac K, Twick I, Hosgor E, Candinas D, et al. Automation of surgical
636 skill assessment using a three-stage machine learning algorithm. *Sci Reports* 2021 111
637 [Internet]. 2021 Mar 4 [cited 2022 Oct 6];11(1):1–9. Available from:

- 638 <https://www.nature.com/articles/s41598-021-84295-6>
- 639 17. Maier-Hein L, Eisenmann M, Sarikaya D, März K, Collins T, Malpani A, et al. Surgical data
640 science – from concepts toward clinical translation. *Med Image Anal.* 2022 Feb
641 1;76:102306.
- 642 18. Kolbinger FR, Leger S, Carstens M, Rinner FM, Krell S, Chernykh A, et al. Artificial
643 Intelligence for context-aware surgical guidance in complex robot-assisted oncological
644 procedures: an exploratory feasibility study. *medRxiv [Internet].* 2022 May 3 [cited 2022
645 May 13];2022.05.02.22274561. Available from:
646 <https://www.medrxiv.org/content/10.1101/2022.05.02.22274561v1>
- 647 19. Madani A, Namazi B, Altieri MS, Hashimoto DA, Rivera AM, Pucher PH, et al. Artificial
648 Intelligence for Intraoperative Guidance. *Ann Surg.* 2020;
- 649 20. Fecso AB, Szasz P, Kerezov G, Grantcharov TP. The effect of technical performance on
650 patient outcomes in surgery. *Ann Surg [Internet].* 2017 Mar 1 [cited 2022 Oct 6];265(3):492–
651 501. Available from:
652 [https://journals.lww.com/annalsofsurgery/Fulltext/2017/03000/The_Effect_of_Technical_P](https://journals.lww.com/annalsofsurgery/Fulltext/2017/03000/The_Effect_of_Technical_Performance_on_Patient.11.aspx)
653 [erformance_on_Patient.11.aspx](https://journals.lww.com/annalsofsurgery/Fulltext/2017/03000/The_Effect_of_Technical_Performance_on_Patient.11.aspx)
- 654 21. Mazzocco K, Petitti DB, Fong KT, Bonacum D, Brookey J, Graham S, et al. Surgical team
655 behaviors and patient outcomes. *Am J Surg.* 2009 May 1;197(5):678–85.
- 656 22. Suliburk JW, Buck QM, Pirko CJ, Massarweh NN, Barshes NR, Singh H, et al. Analysis of
657 Human Performance Deficiencies Associated With Surgical Adverse Events. *JAMA Netw*
658 *Open [Internet].* 2019 Jul 3 [cited 2022 Oct 3];2(7):e198067–e198067. Available from:
659 <https://jamanetwork.com/journals/jamanetworkopen/fullarticle/2740065>
- 660 23. Adelman MR, Bardsley TR, Sharp HT. Urinary Tract Injuries in Laparoscopic Hysterectomy:
661 A Systematic Review. *J Minim Invasive Gynecol.* 2014 Jul 1;21(4):558–66.
- 662 24. Freund MR, Kent I, Horesh N, Smith T, Emile SH, Wexner SD. Pancreatic injuries following
663 laparoscopic splenic flexure mobilization. *Int J Colorectal Dis [Internet].* 2022 Apr 1 [cited
664 2023 Mar 20];37(4):967–71. Available from:
665 <https://link.springer.com/article/10.1007/s00384-022-04112-y>
- 666 25. Carstens M, Rinner FM, Bodenstedt S, Jenke AC, Weitz J, Distler M, et al. The Dresden
667 Surgical Anatomy Dataset for Abdominal Organ Segmentation in Surgical Data Science.
668 *Sci Data [Internet].* 2023 Jan 12 [cited 2023 Jan 12];10(1):1–8. Available from:
669 <https://www.nature.com/articles/s41597-022-01719-2>
- 670 26. Chen L-C, Papandreou G, Schroff F, Adam H. Rethinking Atrous Convolution for Semantic
671 Image Segmentation. *arXiv [Internet].* 2017 Jun 17 [cited 2022 Oct 10]; Available from:

- 672 <https://arxiv.org/abs/1706.05587v3>
- 673 27. Lin TY, Maire M, Belongie S, Hays J, Perona P, Ramanan D, et al. Microsoft COCO:
674 Common Objects in Context. Lect Notes Comput Sci (including Subser Lect Notes Artif
675 Intell Lect Notes Bioinformatics) [Internet]. 2014 May 1 [cited 2022 Nov 11];8693
676 LNCS(PART 5):740–55. Available from: <https://arxiv.org/abs/1405.0312v3>
- 677 28. Xie E, Wang W, Yu Z, Anandkumar A, Alvarez JM, Luo P. SegFormer: Simple and Efficient
678 Design for Semantic Segmentation with Transformers. Adv Neural Inf Process Syst. 2021
679 Dec 6;34:12077–90.
- 680 29. Cordts M, Omran M, Ramos S, Rehfeld T, Enzweiler M, Benenson R, et al. The Cityscapes
681 Dataset for Semantic Urban Scene Understanding. Proc IEEE Conf Comput Vis Pattern
682 Recognit. 2016;
- 683 30. Loshchilov I, Hutter F. Decoupled Weight Decay Regularization. 7th Int Conf Learn
684 Represent ICLR 2019 [Internet]. 2017 Nov 14 [cited 2022 Nov 9]; Available from:
685 <https://arxiv.org/abs/1711.05101v3>
- 686 31. Leibetseder A, Petscharnig S, Primus MJ, Kletz S, Münzer B, Schoeffmann K, et al.
687 LapGyn4: A Dataset for 4 Automatic Content Analysis Problems in the Domain of
688 Laparoscopic Gynecology. Proc 9th ACM Multimed Syst Conf [Internet]. 2018 [cited 2021
689 Jul 19];18. Available from: <https://doi.org/10.1145/3204949.3208127>
- 690 32. Reddy CL, Mitra S, Meara JG, Atun R, Afshar S. Artificial Intelligence and its role in surgical
691 care in low-income and middle-income countries. Lancet Digit Heal. 2019 Dec 1;1(8):e384–
692 6.
- 693 33. Moglia A, Georgiou K, Georgiou E, Satava RM, Cuschieri A. A systematic review on artificial
694 intelligence in robot-assisted surgery. Int J Surg. 2021 Nov 1;95:106151.
- 695 34. Anteby R, Horesh N, Soffer S, Zager Y, Barash Y, Amiel J, et al. Deep learning visual
696 analysis in laparoscopic surgery: a systematic review and diagnostic test accuracy meta-
697 analysis. Surg Endosc [Internet]. 2021 Apr 1 [cited 2021 Jul 16];35(4):1521–33. Available
698 from: <https://pubmed.ncbi.nlm.nih.gov/33398560/>
- 699 35. Kuo RYL, Harrison CJ, Jones BE, Geoghegan L, Furniss D. Perspectives: A surgeon's
700 guide to machine learning. Int J Surg. 2021 Oct 1;94:106133.
- 701 36. Reinke A, Tizabi MD, Sudre CH, Eisenmann M, Rädtsch T, Baumgartner M, et al. Common
702 Limitations of Image Processing Metrics: A Picture Story. arXiv [Internet]. 2021 Apr 12 [cited
703 2022 May 13]; Available from: <https://arxiv.org/abs/2104.05642v4>
- 704 37. Jin C, Udupa JK, Zhao L, Tong Y, Odhner D, Pednekar G, et al. Object recognition in
705 medical images via anatomy-guided deep learning. Med Image Anal. 2022 Oct

- 706 1;81:102527.
- 707 38. Renard F, Guedria S, Palma N De, Vuillerme N. Variability and reproducibility in deep
708 learning for medical image segmentation. *Sci Rep* [Internet]. 2020 Aug 13 [cited 2022 Oct
709 16];10(1):1–16. Available from: <https://www.nature.com/articles/s41598-020-69920-0>
- 710 39. Powers DMW, Ailab. Evaluation: from precision, recall and F-measure to ROC,
711 informedness, markedness and correlation. *arXiv* [Internet]. 2020 Oct 11 [cited 2022 Oct
712 16]; Available from: <https://arxiv.org/abs/2010.16061v1>
- 713 40. Parikh RB, Teeple S, Navathe AS. Addressing Bias in Artificial Intelligence in Health Care.
714 *JAMA* [Internet]. 2019 Dec 24 [cited 2022 Oct 16];322(24):2377–8. Available from:
715 <https://jamanetwork.com/journals/jama/fullarticle/2756196>
- 716 41. Zhang Y, Mehta S, Caspi A. Rethinking Semantic Segmentation Evaluation for
717 Explainability and Model Selection. 2021 Jan 21 [cited 2021 Jul 27]; Available from:
718 <https://arxiv.org/abs/2101.08418v1>
- 719 42. Hashimoto DA, Rosman G, Witkowski ER, Stafford C, Navarette-Welton AJ, Rattner DW,
720 et al. Computer Vision Analysis of Intraoperative Video: Automated Recognition of
721 Operative Steps in Laparoscopic Sleeve Gastrectomy. *Ann Surg* [Internet]. 2019 Sep 1
722 [cited 2021 Jul 16];270(3):414–21. Available from:
723 [https://journals.lww.com/annalsofsurgery/Fulltext/2019/09000/Computer_Vision_Analysis
724 _of_Intraoperative_Video_3.aspx](https://journals.lww.com/annalsofsurgery/Fulltext/2019/09000/Computer_Vision_Analysis_of_Intraoperative_Video_3.aspx)
- 725 43. Hu YY, Mazer LM, Yule SJ, Arriaga AF, Greenberg CC, Lipsitz SR, et al. Complementing
726 Operating Room Teaching With Video-Based Coaching. *JAMA Surg* [Internet]. 2017 Apr 1
727 [cited 2022 Oct 27];152(4):318–25. Available from:
728 <https://jamanetwork.com/journals/jamasurgery/fullarticle/2593311>
- 729 44. Mizota T, Anton NE, Stefanidis D. Surgeons see anatomical structures faster and more
730 accurately compared to novices: Development of a pattern recognition skill assessment
731 platform. *Am J Surg*. 2019 Feb 1;217(2):222–7.
- 732 45. Ward TM, Mascagni P, Ban Y, Rosman G, Padoy N, Meireles O, et al. Computer vision in
733 surgery. *Surgery* [Internet]. 2021 May 1 [cited 2022 Oct 27];169(5):1253–6. Available from:
734 <https://pubmed.ncbi.nlm.nih.gov/33272610/>
- 735 46. Chopra H, Baig AA, Arora S, Singh I, Kaur A, Emran T Bin. Artificial intelligence in surgery:
736 Modern trends. *Int J Surg*. 2022 Oct 1;106:106883.
- 737 47. Docea R, Pfeiffer M, Bodenstedt S, Kolbinger FR, Höller L, Wittig I, et al. Simultaneous
738 localisation and mapping for laparoscopic liver navigation : a comparative evaluation study.
739 In: Linte CA, Siewerdsen JH, editors. *Medical Imaging 2021: Image-Guided Procedures*,

- 740 Robotic Interventions, and Modeling [Internet]. SPIE; 2021 [cited 2021 Feb 25]. p. 8.
741 Available from: [https://www.spiedigitallibrary.org/conference-proceedings-of-](https://www.spiedigitallibrary.org/conference-proceedings-of-spie/11598/2582121/Simultaneous-localisation-and-mapping-for-laparoscopic-liver-navigation/10.1117/12.2582121.full)
742 [spie/11598/2582121/Simultaneous-localisation-and-mapping-for-laparoscopic-liver-](https://www.spiedigitallibrary.org/conference-proceedings-of-spie/11598/2582121/Simultaneous-localisation-and-mapping-for-laparoscopic-liver-navigation/10.1117/12.2582121.full)
743 [navigation/10.1117/12.2582121.full](https://www.spiedigitallibrary.org/conference-proceedings-of-spie/11598/2582121/Simultaneous-localisation-and-mapping-for-laparoscopic-liver-navigation/10.1117/12.2582121.full)
- 744 48. Shaalan D, Jusoh S. Visualization in Medical System Interfaces: UX Guidelines. Proc 12th
745 Int Conf Electron Comput Artif Intell ECAI 2020. 2020 Jun 1;
- 746 49. Henry KE, Kornfield R, Sridharan A, Linton RC, Groh C, Wang T, et al. Human-machine
747 teaming is key to AI adoption: clinicians' experiences with a deployed machine learning
748 system. npj Digit Med [Internet]. 2022 Jul 21 [cited 2023 Mar 30];5(1):1–6. Available from:
749 <https://www.nature.com/articles/s41746-022-00597-7>
750

751 **ABBREVIATIONS**

752	AI	Artificial Intelligence
753	IoU	Intersection-over-Union
754	SD	Standard deviation

755

756 **AUTHOR CONTRIBUTIONS**

757 FRK, JW, MD, SS, and SB conceptualized the study. FRK, FMR, and MC collected and annotated
758 clinical and video data and contributed to data analysis. ACJ, SK, SL, and SB implemented and
759 trained the neural networks and contributed to data analysis. JW, MD, and SS supervised the
760 project, provided infrastructure and gave important scientific input. FRK drafted the initial
761 manuscript text. All authors reviewed, edited, and approved the final manuscript.

762