

Better than humans? Machine learning-based anatomy recognition in minimally-invasive abdominal surgery

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

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

✉ 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

30 **Abstract**

31 **Background:** Lack of anatomy recognition represents a clinically relevant risk factor in abdominal
32 surgery. While machine learning methods have the potential to aid in recognition of visible patterns
33 and structures, limited availability and diversity of (annotated) laparoscopic image data restrict the
34 clinical potential of such applications in practice. This study explores the potential of machine
35 learning algorithms to identify and delineate abdominal organs and anatomical structures using a
36 robust and comprehensive dataset, and compares algorithm performance to that of humans.

37 **Methods:** Based on the Dresden Surgical Anatomy Dataset providing 13195 laparoscopic images
38 with pixel-wise segmentations of eleven anatomical structures, two machine learning algorithms
39 were developed: individual segmentation algorithms for each structure, and a combined algorithm
40 with a common encoder and structure-specific decoders. Performance was assessed using F1
41 score, Intersection-over-Union (IoU), precision, recall, and specificity. Using the example of
42 pancreas segmentation on a sample dataset of 35 images, algorithm performance was compared
43 to that of a cohort of 28 physicians, medical students, and medical laypersons.

44 **Results:** Mean IoU for segmentation of intraabdominal structures ranged from 0.28 to 0.83 and
45 from 0.32 to 0.81 for the structure-specific and the combined semantic segmentation model,
46 respectively. Average inference for the structure-specific (one anatomical structure) and the
47 combined model (eleven anatomical structures) took 20 ms and 54 ms, respectively. The
48 structure-specific model performed equal to or better than 27 out of 28 human participants in
49 pancreas segmentation.

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

54

55 Introduction

56 Computer vision describes the computerized analysis of digital images aiming at the automation
57 of human visual capabilities, most commonly using machine learning methods, in particular deep
58 learning. This approach has transformed medicine in recent years, with successful applications
59 including computer-aided diagnosis of colonic polyp dignity in endoscopy^{1,2}, detection of clinically
60 actionable genetic alterations in histopathology³, and melanoma detection in dermatology⁴.
61 Availability of large amounts of training data is the defining prerequisite for successful application
62 of deep learning methods. With the establishment of laparoscopy as the gold standard for a variety
63 of surgical procedures⁵⁻⁸ and the increasing availability of computing resources, these concepts
64 have gradually been applied to abdominal surgery. The overwhelming majority of research efforts
65 in the field of Artificial Intelligence (AI)-based analysis of intraoperative surgical imaging data (i.e.
66 video data from laparoscopic or open surgeries) has focused on classifying images with respect
67 to the presence and/or location of previously annotated surgical instruments or anatomical
68 structures⁹⁻¹³ or on analysis of surgical proficiency¹⁴⁻¹⁶ based on recorded procedures. However,
69 almost all research endeavors in the field of computer vision in laparoscopic surgery have
70 concentrated on preclinical stages and to date, no AI model based on intraoperative surgical
71 imaging data could demonstrate a palpable clinical benefit.¹⁷ Among the studies closest to clinical
72 application are recent works on identification of instruments and hepatobiliary anatomy during
73 cholecystectomy for automated assessment of the critical view of safety¹³, and on the automated
74 segmentation of safe and unsafe preparation zones during cholecystectomy¹⁸.

75 In surgery, patient outcome heavily depends on experience and performance of the surgical
76 team.^{19,20} In a recent analysis of Human Performance Deficiencies in major cardiothoracic,
77 vascular, abdominal transplant, surgical oncology, acute care, and general surgical operations,
78 more than half of the cases with postoperative complications were associated with identifiable
79 human error. Among these errors, lack of recognition (including misidentified anatomy) accounted
80 for 18.8%, making it the most common Human Performance Deficiency overall.²¹ While AI-based
81 systems identifying anatomical risk and target structures would theoretically have the potential to
82 alleviate this risk, limited availability and diversity of (annotated) laparoscopic image data
83 drastically restrict the clinical potential of such applications in practice.

84 To advance and diversify the applications of computer vision in laparoscopic surgery, we have
85 recently published the Dresden Surgical Anatomy Dataset²², providing 13195 laparoscopic images
86 with high-quality annotations of the presence and exact location of eleven intraabdominal
87 anatomical structures: abdominal wall, colon, intestinal vessels (inferior mesenteric artery and
88 inferior mesenteric vein with their subsidiary vessels), liver, pancreas, small intestine, spleen,

89 stomach, ureter and vesicular glands. Here, we present the first study evaluating automated
90 detection and localization of organs and anatomical structures in laparoscopic view based on this
91 dataset, and, using the example of delineation of the pancreas, compare algorithm performance
92 to that of humans.

93

94 **Methods**

95 **Patient cohort**

96 Video data from 32 robot-assisted anterior rectal resections or rectal extirpations were gathered
97 at the University Hospital Carl Gustav Carus Dresden between February 2019 and February 2021.
98 All included patients had a clinical indication for the surgical procedure, recommended by an
99 interdisciplinary tumor board. The procedures were performed using the da Vinci® Xi system
100 (Intuitive Surgical, Sunnyvale, CA, USA) with a standard Da Vinci® Xi/X Endoscope with Camera
101 (8 mm diameter, 30° angle, Intuitive Surgical, Sunnyvale, CA, USA, Item code 470057). Surgeries
102 were recorded using the CAST system (Orpheus Medical GmbH, Frankfurt a.M., Germany). Each
103 record was saved at a resolution of 1920 x 1080 pixels in MPEG-4 format.

104 All experiments were performed in accordance with the ethical standards of the Declaration of
105 Helsinki and its later amendments. The local Institutional Review Board (ethics committee at the
106 Technical University Dresden) reviewed and approved this study (approval number: BO-EK-
107 137042018). The trial was registered on clinicaltrials.gov (trial registration ID: NCT05268432).
108 Written informed consent to laparoscopic image data acquisition, data annotation, data analysis,
109 and anonymized data publication was obtained from all participants. Before publication, all data
110 was anonymized according to the general data protection regulation of the European Union.

111

112 **Dataset**

113 Based on the full-length surgery recordings and respective temporal annotations of organ visibility,
114 individual image frames were extracted and annotated as described previously. The resulting
115 Dresden Surgical Anatomy Dataset comprises 13195 distinct images with pixel-wise
116 segmentations of eleven anatomical structures: abdominal wall, colon, intestinal vessels (inferior
117 mesenteric artery and inferior mesenteric vein with their subsidiary vessels), liver, pancreas, small
118 intestine, spleen, stomach, ureter and vesicular glands. Moreover, the dataset comprises binary
119 annotations of the presence of each of these organs for each image. The dataset is publicly
120 available via the following link: <https://figshare.com/s/d7a60b74989a9cab2f7f>.

121

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

- 124 — Training set (at least 12 surgeries per anatomical structure): surgeries 1, 4, 5, 6, 8, 9, 10,
125 12, 15, 16, 17, 19, 22, 23, 24, 25, 27, 28, 29, 30, 31.
- 126 — Validation set (3 surgeries per anatomical structure): surgeries 3, 21, 26.
- 127 — Test set (5 surgeries per anatomical structure): surgeries 2, 7, 11, 13, 14, 18, 20, 32.

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

134

135 **Structure-specific semantic segmentation model**

136 To segment each anatomical structure, a separate convolutional neural network for segmentation
137 a DeeplabV3²³ model with a ResNet50 backbone with default PyTorch pretraining on the COCO
138 dataset²⁴, was used. The networks were trained using cross-entropy loss and the AdamW
139 optimizer²⁵ for 100 epochs with a starting learning rate of 10^{-4} and a linear learning rate scheduler
140 decreasing the learning rate by 0.9 every 10 epochs. For data augmentation, we applied random
141 scaling and rotation, as well as brightness adjustments. The final model for each organ was
142 selected via the Intersection-over-Union (IoU) on the validation dataset and evaluated using the
143 Dresden Surgical Anatomy Dataset with the abovementioned training-validation-test split (Figure
144 1).

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

149

150 **Combined semantic segmentation model**

151 A convolutional neural network with a common encoder and eleven decoders for combined
152 segmentation of the eleven anatomical structures was trained. The used architecture is an
153 extension of DeepLabV3²³. A shared ResNet50 backbone with default PyTorch pretraining on the
154 COCO dataset²⁴, was used. For each class, a DeepLabV3 decoder was then run on the features

155 extracted from a given image by the backbone. As the images are only annotated for binary
156 classes, the loss is only calculated for the decoder associated with the class annotated in a given
157 image. The remaining training procedure was identical to the structure-specific model. The model
158 was trained and evaluated using the Dresden Surgical Anatomy Dataset with the abovementioned
159 training-validation-test split (Figure 1).

160 Segmentation performance was assessed using F1 score, IoU, precision, recall, and specificity
161 on the test folds.²⁶

162

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

164 To determine the clinical potential of automated segmentation of anatomical risk structures, the
165 segmentation performance of 28 humans was compared to that of the structure-specific semantic
166 segmentation model using the example of the pancreas. The local Institutional Review Board
167 (ethics committee at the Technical University Dresden) reviewed and approved this study
168 (approval number: BO-EK-566122021). All participants provided written informed consent to
169 anonymous study participation, data acquisition and analysis, and publication. In total, 28
170 participants (physician and non-physician medical staff, medical students, and medical
171 laypersons) marked the pancreas in 35 images from the Dresden Surgical Anatomy Dataset with
172 bounding boxes. These images originated from 26 different surgeries, and the pancreas was
173 visible in 16 of the 35 images. Each of the previously selected 35 images was shown once, the
174 order being arbitrarily chosen but identical for all participants. The open-source annotation
175 software Computer Vision Annotation Tool (CVAT) was used for annotations. In cases where the
176 pancreas was seen in multiple, non-connected locations in the image, participants were asked to
177 create separate bounding boxes for each area.

178 Based on the structure-specific semantic segmentation model, axis-aligned bounding boxes
179 marking the pancreas were generated in the 35 images from the pixel-wise segmentation. To
180 guarantee that the respective images were not part of the training data, four-fold cross validation
181 was used, i.e. the origin surgeries were split into four equal-sized batches, and algorithms were
182 trained on three batches that did not contain the respective origin image before being applied to
183 segmentation.

184 To compare human and algorithm performance, the bounding boxes created by each participant
185 and the structure-specific semantic segmentation model were compared to bounding boxes
186 derived from the Dresden Surgical Anatomy Dataset, which were defined as ground truth. IoU
187 between the manual or algorithmical bounding box and the ground truth was used to compare
188 segmentation accuracy.

189

190 **Data Availability**

191 The Dresden Surgical Anatomy Dataset is publicly available via the following link:
192 <https://figshare.com/s/d7a60b74989a9cab2f7f>. All other data generated and analyzed during the
193 current study are available from the corresponding authors on reasonable request. To gain
194 access, data requestors will need to sign a data access agreement.

195

196 **Code Availability**

197 The most relevant scripts used for dataset compilation are publicly available via the following link:
198 <https://zenodo.org/record/6958337#.YzsBdnZBzOg>.

199

200 **Results**

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

202 Structure-specific multi-layer convolutional neural networks (Figure 1) were trained to segment
203 the abdominal wall, the colon, intestinal vessels (inferior mesenteric artery and inferior mesenteric
204 vein with their subsidiary vessels), the liver, the pancreas, the small intestine, the spleen, the
205 stomach, the ureter, and vesicular glands. Table 1 displays mean F1 score, IoU, precision, recall,
206 and specificity for individual anatomical structures as predicted by the structure-specific
207 algorithms. Out of the analyzed segmentation models, performance was lowest for the vesicular
208 glands (mean IoU: 0.28 ± 0.21) and the pancreas (mean IoU: 0.28 ± 0.27), while excellent
209 predictions were achieved for the abdominal wall (mean IoU: 0.83 ± 0.14) and the small intestine
210 (mean IoU: 0.80 ± 0.18). In segmentation of the pancreas, the ureter, the vesicular glands and the
211 intestinal veins, there was a proportion of images with no detection or no overlap between ground
212 truth, while for all remaining anatomical structures, this proportion was minimal (Figure 2 a). While
213 the images, in which the highest IoUs were observed, mostly displayed large organ segments that
214 were clearly visible (Figure 2 b), the images with the lowest IoU were of variable quality with
215 confounding factors such as blood, smoke, soiling of the endoscope lens, extreme zoom, or
216 pictures blurred by camera shake (Figure 2 c).

217 Inference on a single image with a resolution of 640 x 512 pixels required, on average, 20 ms on
218 an Nvidia A5000, resulting in a frame rate of 50 frames per second. This runtime includes one
219 decoder, meaning that only the segmentation for one anatomical class is included.

220

221 **Table 1: Summary of performance metrics for anatomical structure segmentation using**
222 **structure-specific models based on the DeepLabv3 architecture.** For each metric, mean and
223 standard deviation are displayed.
224

Anatomical structure	F1 score	IoU	Precision	Recall	Specificity
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

225

226 **Machine Learning-based anatomical structure segmentation in a combined model**
227 Using all annotated images from the Dresden Surgical Anatomy Dataset, a combined model with
228 a mutual encoder and organ-specific decoders was trained (Figure 1). Table 2 displays mean F1
229 score, IoU, precision, recall, and specificity for anatomical structure segmentation in the combined
230 model. The performance of the combined model was overall similar to that of structure-specific
231 models (Table 1), with highest segmentation performance for the abdominal wall (IoU: 0.81 ± 0.16)
232 and the small intestine (IoU: 0.77 ± 0.19), and the lowest performance for the vesicular glands
233 (IoU: 0.32 ± 0.24) and the pancreas (IoU: 0.33 ± 0.25). As for the structure-specific models, few
234 images were not or entirely mispredicted; the proportion of such images was largest in the
235 pancreas and the vesicular glands (Figure 3 a). Similar trends towards an impact of segment size,
236 zoom, endoscope lens soiling, blurry images, and presence of blood or smoke were seen as for
237 structure-specific models when comparing image quality of well-predicted images (Figure 3 b) and
238 images with poor or no prediction (Figure 3 c).

239 Inference on a single image with a resolution of 640 x 512 pixels required, on average, 54 ms on
240 an Nvidia A5000, resulting in a frame rate of about 18.5 frames per second. This runtime includes
241 all 11 decoders, meaning that segmentations for all organ classes are included.
242

243 **Table 2: Summary of performance metrics for anatomical structure segmentation using the**
244 **combined model (common encoder with structure-specific decoders).** For each metric,
245 mean and standard deviation are displayed.
246

Anatomical structure	F1 score	IoU	Precision	Recall	Specificity
Abdominal wall	0.89 ± 0.12	0.81 ± 0.16	0.85 ± 0.16	0.95 ± 0.07	0.96 ± 0.05
Colon	0.79 ± 0.17	0.69 ± 0.20	0.77 ± 0.18	0.86 ± 0.19	0.96 ± 0.05
Inferior mesenteric artery	0.63 ± 0.21	0.49 ± 0.20	0.57 ± 0.22	0.78 ± 0.23	0.99 ± 0.01
Intestinal veins	0.60 ± 0.25	0.47 ± 0.22	0.60 ± 0.24	0.72 ± 0.31	0.99 ± 0.01
Liver	0.81 ± 0.19	0.72 ± 0.22	0.80 ± 0.22	0.89 ± 0.15	0.98 ± 0.05
Pancreas	0.44 ± 0.30	0.33 ± 0.25	0.54 ± 0.34	0.49 ± 0.36	0.98 ± 0.10
Small intestine	0.85 ± 0.15	0.77 ± 0.19	0.80 ± 0.19	0.94 ± 0.12	0.95 ± 0.06
Spleen	0.80 ± 0.19	0.70 ± 0.22	0.72 ± 0.20	0.95 ± 0.15	0.99 ± 0.01
Stomach	0.70 ± 0.27	0.60 ± 0.28	0.70 ± 0.26	0.81 ± 0.29	0.99 ± 0.01
Ureter	0.52 ± 0.26	0.39 ± 0.23	0.45 ± 0.25	0.71 ± 0.34	0.99 ± 0.01
Vesicular glands	0.44 ± 0.29	0.32 ± 0.24	0.47 ± 0.31	0.53 ± 0.35	0.98 ± 0.02

247

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

249 To approximate the clinical value of the previously described algorithms for anatomical structure
250 segmentation, the performance of the structure-specific model was compared to that of a cohort
251 of 28 physicians, medical students, and persons with no medical background (Figure 4 a), and
252 different degrees of experience in laparoscopic surgery (Figure 4 b). A vulnerable anatomical
253 structure with – measured by classical metrics of overlap (Tables 1 and 2) – comparably weak
254 segmentation performance of the trained algorithms, the pancreas was selected as an example.
255 Comparing bounding box segmentations of the pancreas of participants and the machine learning
256 model, the medical and laparoscopy-specific experience of human participants was mirrored by
257 the respective IoUs describing the overlap between the pancreas annotation and the ground truth.
258 The pancreas-specific segmentation model (IoU: 0.29) performed equal to or better than 27 out
259 of the 28 human participants (Figures 4 c and d). Overall, these results demonstrate that the
260 developed models have clinical potential to improve the recognition of vulnerable anatomical
261 structures.
262

263 Discussion

264 In surgery, misinterpretation of visual cues can result in objectifiable errors with serious
265 consequences.²¹ Based on a robust public dataset providing 13195 laparoscopic images with
266 segmentations of eleven intra-abdominal anatomical structures, this study explores the potential
267 of machine learning for automated segmentation of these organs, and compares algorithmic

268 segmentation quality to that of humans with varying experience in minimally-invasive abdominal
269 surgery.

270 In summary, the presented findings suggest that machine learning-based segmentation of
271 intraabdominal organs and anatomical structures is possible and has the potential to provide
272 clinically valuable information. At an average runtime of 54 ms per image, corresponding to a
273 frame rate of 18.5 frames per second, the combined model would facilitate near-real-time
274 identification of eleven anatomical structures. These runtimes mirror the performance of a non-
275 optimized version of the model, which can be significantly improved using methods such as
276 TensorRT from Nvidia. Measured by classical metrics of overlap between segmentation and
277 ground truth, predictions were, overall, better for large and similar-appearing organs such as the
278 abdominal wall, the liver, the stomach, and the spleen as compared to smaller and more diverse-
279 appearing organs such as the pancreas, the ureter, or vesicular glands. Furthermore, poor image
280 quality (i.e. images blurred by camera movements, presence of blood or smoke in images) was
281 linked to lower accuracy of machine learning-based segmentations. These findings imply that
282 computer vision studies in laparoscopy should be carefully interpreted taking representativity and
283 potential selection of underlying training and validation data into consideration.

284 Measured by classical metrics of overlap (e.g. IoU, F1 score, precision, recall, specificity) that are
285 commonly used to evaluate segmentation performance, the structure-specific models and the
286 combined model performed similarly with average IoUs ranging from 0.28 to 0.83 and from 0.32
287 to 0.81, respectively. Interpretation of these metrics, however, represents a major challenge in
288 computer vision applications in medical domains such as dermatology and endoscopy²⁷⁻²⁹ as well
289 as non-medical domains such as autonomous driving³⁰. In the specific use case of laparoscopic
290 surgery, evidence suggests that such technical metrics alone are not sufficient to characterize the
291 clinical potential and utility of segmentation algorithms.^{31,32} In this context, the subjective clinical
292 utility of a bounding box-based detection system recognizing the common bile duct and the cystic
293 duct at average precisions of 0.32 and 0.07, respectively, demonstrated by Tokuyasu *et al.*,
294 supports this hypothesis.¹² In the presented analysis, the trained structure-specific machine
295 learning algorithm performed equal to or better than all human participants in the specific task of
296 bounding box segmentation of the pancreas except for one expert with over 10 years of
297 experience. This suggests that even for structures such as the pancreas with seemingly poor
298 segmentation quality (segmentation IoU of the best-performing model: 0.33 ± 0.25 in the test set)
299 have the potential to provide clinically valuable help in anatomy recognition. Notably, the best
300 average IoUs achieved in this comparative study were 0.29 (for the structure-specific model) and
301 0.36 (for the best human participant), which would both be considered less reliable segmentation

302 quality measures on paper. This encourages further discussion about metrics for segmentation
303 quality assessment in clinical AI. In the future, the potential of the described dataset²² and organ
304 segmentation algorithms could be exploited for educational purposes^{33,34}, for guidance systems
305 facilitating real-time detection of risk and target structures^{18,32,35}, or as an auxiliary function
306 integrated in more complex surgical assistance systems, such as guidance systems relying on
307 automated liver registration³⁶.

308 The limitations of this work are mostly related to the dataset and general limitations of machine
309 learning-based segmentation: First, the Dresden Surgical Anatomy Dataset is a monocentric
310 dataset based on 32 robot-assisted rectal surgeries. Therefore, the images used for algorithm
311 training and validation display organs from specific angles, which could limit generalizability and
312 transferability of the presented findings to other minimally-invasive abdominal surgeries,
313 particularly non-robotic procedures. Second, annotations were required for training of machine
314 learning algorithms, potentially inducing some bias towards the way that organs were annotated
315 in the algorithms, which may differ from individual healthcare professionals' way of recognizing an
316 organ. This is particularly relevant for organs such as the ureters or the pancreas, which often
317 appear covered by layers of tissue. Here, computer vision-based algorithms that solely consider
318 the laparoscopic images provided by the Dresden Surgical Anatomy Dataset for identification of
319 risk structures will only be able to identify an organ once it is visible. For an earlier recognition of
320 such hidden risk structures, more training data with meaningful annotations would be necessary.
321 Importantly, the presented comparison to human performance focused on segmentation of visible
322 anatomy as well, neglecting that humans (and possibly computers, too) could already identify a
323 risk structure hidden underneath tissue layers. The existing limitations notwithstanding, the
324 presented study represents an important addition to the growing body of research on medical
325 image analysis in laparoscopic surgery, particularly by linking technical metrics to human
326 performance.

327 In conclusion, this study demonstrates that machine learning methods have the potential to
328 provide clinically relevant near-real-time assistance in anatomy recognition in minimally-invasive
329 surgery. This study is the first to use the recently published Dresden Surgical Anatomy Dataset,
330 providing baseline algorithms for organ segmentation and evaluating the clinical relevance of such
331 algorithms. Future research should investigate other segmentation methods, the transferability of
332 these results to other surgical procedures, and the clinical impact of real-time surgical assistance
333 systems and didactic applications based on automated segmentation algorithms.

334
335

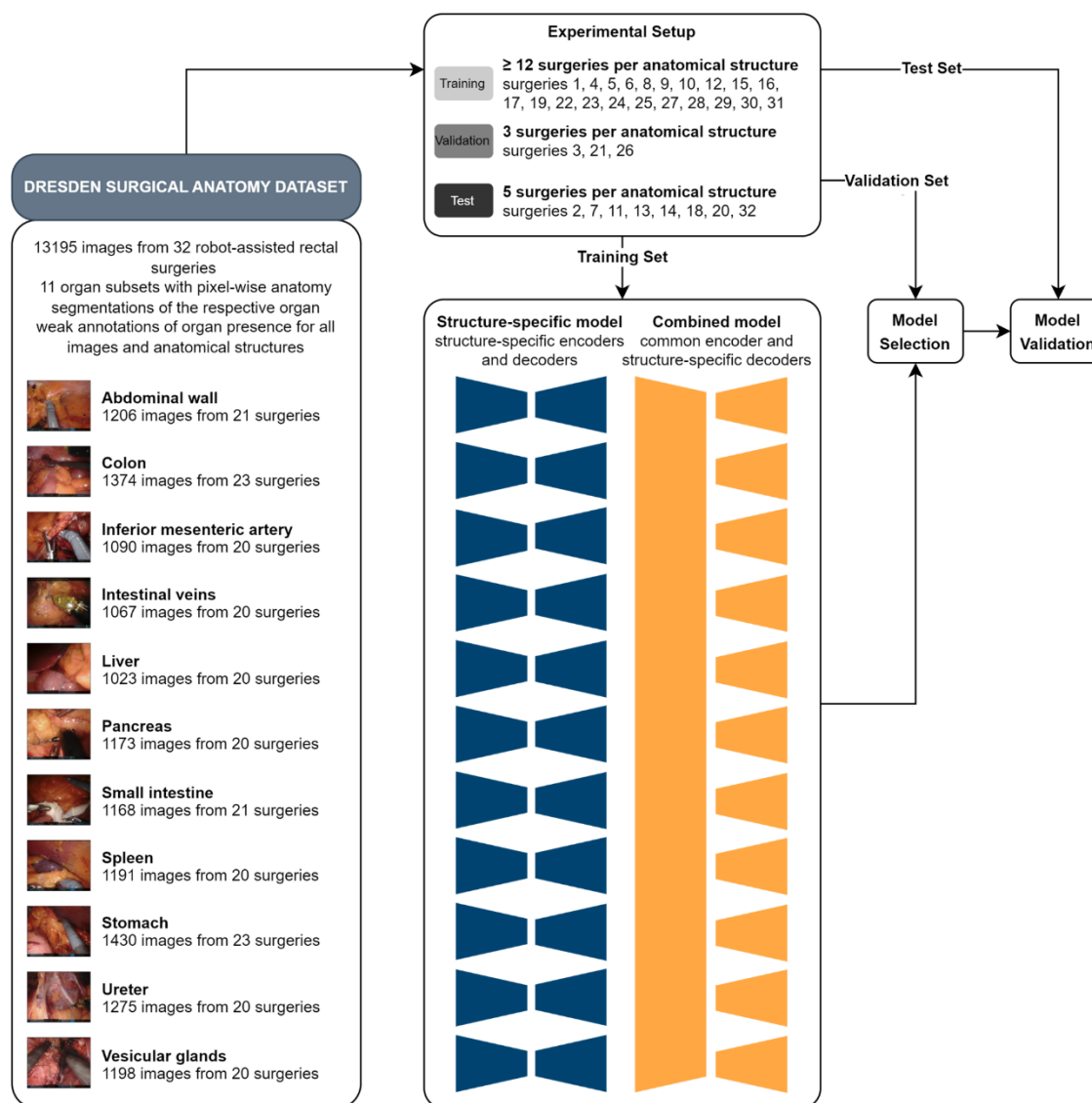
336 References

- 337 1. Wang, P. *et al.* Effect of a deep-learning computer-aided detection system on adenoma
338 detection during colonoscopy (CADE-DB trial): a double-blind randomised study. *Lancet*
339 *Gastroenterol. Hepatol.* **5**, 343–351 (2020).
- 340 2. Wang, P. *et al.* Real-time automatic detection system increases colonoscopic polyp and
341 adenoma detection rates: a prospective randomised controlled study. *Gut* **68**, 1813–1819
342 (2019).
- 343 3. Kather, J. N. *et al.* Pan-cancer image-based detection of clinically actionable genetic
344 alterations. *Nat. Cancer* **2020 18 1**, 789–799 (2020).
- 345 4. Esteva, A. *et al.* Dermatologist-level classification of skin cancer with deep neural networks.
346 *Nature* **542**, 115–118 (2017).
- 347 5. Simillis, C. *et al.* Open Versus Laparoscopic Versus Robotic Versus Transanal Mesorectal
348 Excision for Rectal Cancer: A Systematic Review and Network Meta-analysis. *Ann. Surg.*
349 **270**, 59–68 (2019).
- 350 6. Zhao, J. J. *et al.* Comparative outcomes of needlescopic, single-incision laparoscopic,
351 standard laparoscopic, mini-laparotomy, and open cholecystectomy: A systematic review
352 and network meta-analysis of 96 randomized controlled trials with 11,083 patients. *Surgery*
353 **170**, 994–1003 (2021).
- 354 7. Luketich, J. D. *et al.* Outcomes after minimally invasive esophagectomy: review of over
355 1000 patients. *Ann. Surg.* **256**, 95–103 (2012).
- 356 8. Thomson, J. E. *et al.* Laparoscopic versus open surgery for complicated appendicitis: a
357 randomized controlled trial to prove safety. *Surg. Endosc.* **29**, 2027–2032 (2015).
- 358 9. Islam, M., Atputharuban, D. A., Ramesh, R. & Ren, H. Real-time instrument segmentation
359 in robotic surgery using auxiliary supervised deep adversarial learning. *IEEE Robot. Autom.*
360 *Lett.* **4**, 2188–2195 (2019).
- 361 10. Roß, T. *et al.* Comparative validation of multi-instance instrument segmentation in
362 endoscopy: results of the ROBUST-MIS 2019 challenge. *Med. Image Anal.* **70**, 101920
363 (2020).
- 364 11. Shvets, A. A., Rakhlin, A., Kalinin, A. A. & Iglovikov, V. I. Automatic Instrument
365 Segmentation in Robot-Assisted Surgery using Deep Learning. in *Proceedings - 17th IEEE*
366 *International Conference on Machine Learning and Applications, ICMLA 2018* 624–628
367 (Institute of Electrical and Electronics Engineers Inc., 2019).
368 doi:10.1109/ICMLA.2018.00100.
- 369 12. Tokuyasu, T. *et al.* Development of an artificial intelligence system using deep learning to

- 370 indicate anatomical landmarks during laparoscopic cholecystectomy. *Surg. Endosc.* 2020
371 354 **35**, 1651–1658 (2020).
- 372 13. Mascagni, P. *et al.* Artificial Intelligence for Surgical Safety: Automatic Assessment of the
373 Critical View of Safety in Laparoscopic Cholecystectomy Using Deep Learning. *Ann. Surg.*
374 (2020) doi:10.1097/SLA.0000000000004351.
- 375 14. Jin, A. *et al.* Tool Detection and Operative Skill Assessment in Surgical Videos Using
376 Region-Based Convolutional Neural Networks. *Proc. - 2018 IEEE Winter Conf. Appl.*
377 *Comput. Vision, WACV 2018* **2018-January**, 691–699 (2018).
- 378 15. Funke, I. *et al.* Using 3D Convolutional Neural Networks to Learn Spatiotemporal Features
379 for Automatic Surgical Gesture Recognition in Video. *Med. Image Comput. Comput. Assist.*
380 *Interv. – MICCAI 2019. Lect. Notes Comput. Sci.* **11768**, 467–475 (2019).
- 381 16. Lavanchy, J. L. *et al.* Automation of surgical skill assessment using a three-stage machine
382 learning algorithm. *Sci. Reports 2021* **11** **11**, 1–9 (2021).
- 383 17. Maier-Hein, L. *et al.* Surgical data science – from concepts toward clinical translation. *Med.*
384 *Image Anal.* **76**, 102306 (2022).
- 385 18. Madani, A. *et al.* Artificial Intelligence for Intraoperative Guidance. *Ann. Surg.* (2020)
386 doi:10.1097/sla.0000000000004594.
- 387 19. Fecso, A. B., Szasz, P., Kerezov, G. & Grantcharov, T. P. The effect of technical
388 performance on patient outcomes in surgery. *Ann. Surg.* **265**, 492–501 (2017).
- 389 20. Mazzocco, K. *et al.* Surgical team behaviors and patient outcomes. *Am. J. Surg.* **197**, 678–
390 685 (2009).
- 391 21. Suliburk, J. W. *et al.* Analysis of Human Performance Deficiencies Associated With Surgical
392 Adverse Events. *JAMA Netw. Open* **2**, e198067–e198067 (2019).
- 393 22. Carstens, M. *et al.* The Dresden Surgical Anatomy Dataset for abdominal organ
394 segmentation in surgical data science. *Figshare* (2022).
- 395 23. Chen, L.-C., Papandreou, G., Schroff, F. & Adam, H. Rethinking Atrous Convolution for
396 Semantic Image Segmentation. *arXiv* (2017) doi:10.48550/arxiv.1706.05587.
- 397 24. Lin, T. Y. *et al.* Microsoft COCO: Common Objects in Context. *Lect. Notes Comput. Sci.*
398 *(including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics)* **8693 LNCS**, 740–
399 755 (2014).
- 400 25. Loshchilov, I. & Hutter, F. Decoupled Weight Decay Regularization. *7th Int. Conf. Learn.*
401 *Represent. ICLR 2019* (2017) doi:10.48550/arxiv.1711.05101.
- 402 26. Leger, S. *et al.* A comparative study of machine learning methods for time-to-event survival
403 data for radiomics risk modelling. *Sci. Rep.* **7**, 11 (2017).

- 404 27. Renard, F., Guedria, S., Palma, N. De & Vuillerme, N. Variability and reproducibility in deep
405 learning for medical image segmentation. *Sci. Rep.* **10**, 1–16 (2020).
- 406 28. Powers, D. M. W. & Ailab. Evaluation: from precision, recall and F-measure to ROC,
407 informedness, markedness and correlation. *arXiv* (2020) doi:10.48550/arxiv.2010.16061.
- 408 29. Parikh, R. B., Teeple, S. & Navathe, A. S. Addressing Bias in Artificial Intelligence in Health
409 Care. *JAMA* **322**, 2377–2378 (2019).
- 410 30. Zhang, Y., Mehta, S. & Caspi, A. Rethinking Semantic Segmentation Evaluation for
411 Explainability and Model Selection. (2021).
- 412 31. Reinke, A. *et al.* Common Limitations of Image Processing Metrics: A Picture Story. *arXiv*
413 (2021) doi:10.48550/arxiv.2104.05642.
- 414 32. Hashimoto, D. A. *et al.* Computer Vision Analysis of Intraoperative Video: Automated
415 Recognition of Operative Steps in Laparoscopic Sleeve Gastrectomy. *Ann. Surg.* **270**, 414–
416 421 (2019).
- 417 33. Hu, Y. Y. *et al.* Complementing Operating Room Teaching With Video-Based Coaching.
418 *JAMA Surg.* **152**, 318–325 (2017).
- 419 34. Mizota, T., Anton, N. E. & Stefanidis, D. Surgeons see anatomical structures faster and
420 more accurately compared to novices: Development of a pattern recognition skill
421 assessment platform. *Am. J. Surg.* **217**, 222–227 (2019).
- 422 35. Ward, T. M. *et al.* Computer vision in surgery. *Surgery* **169**, 1253–1256 (2021).
- 423 36. Docea, R. *et al.* Simultaneous localisation and mapping for laparoscopic liver navigation : a
424 comparative evaluation study. in *Medical Imaging 2021: Image-Guided Procedures,*
425 *Robotic Interventions, and Modeling* (eds. Linte, C. A. & Siewerdsen, J. H.) vol. 11598 8
426 (SPIE, 2021).
- 427

428 Figures and Figure Captions



429

430 **Figure 1**

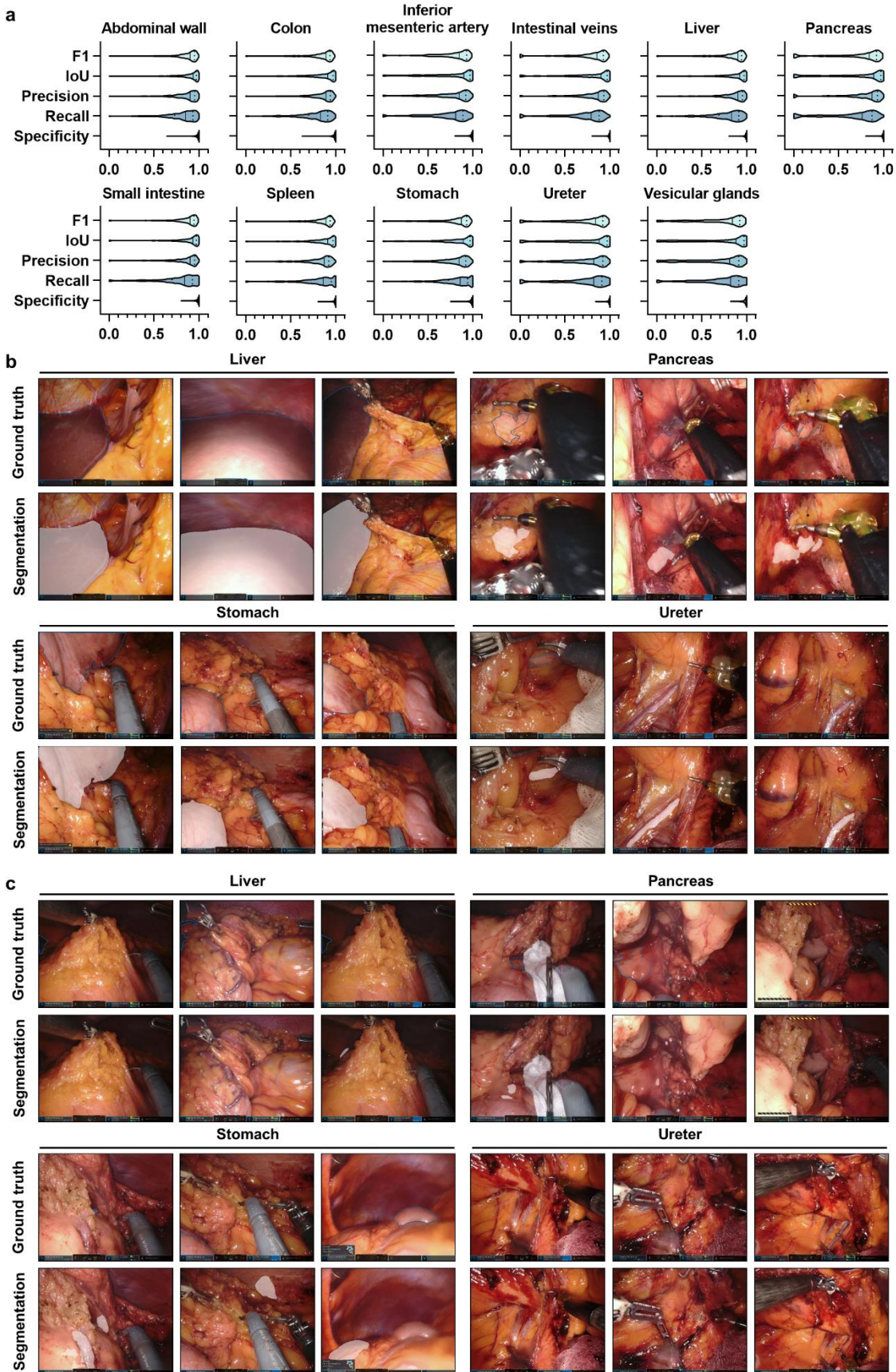
431 **Fig. 1 Schematic illustration of the structure-specific and combined machine learning**

432 **models used for semantic segmentation.** The Dresden Surgical Anatomy Dataset was split into

433 a training, a validation, and a test set. For spatial segmentation, two machine learning models

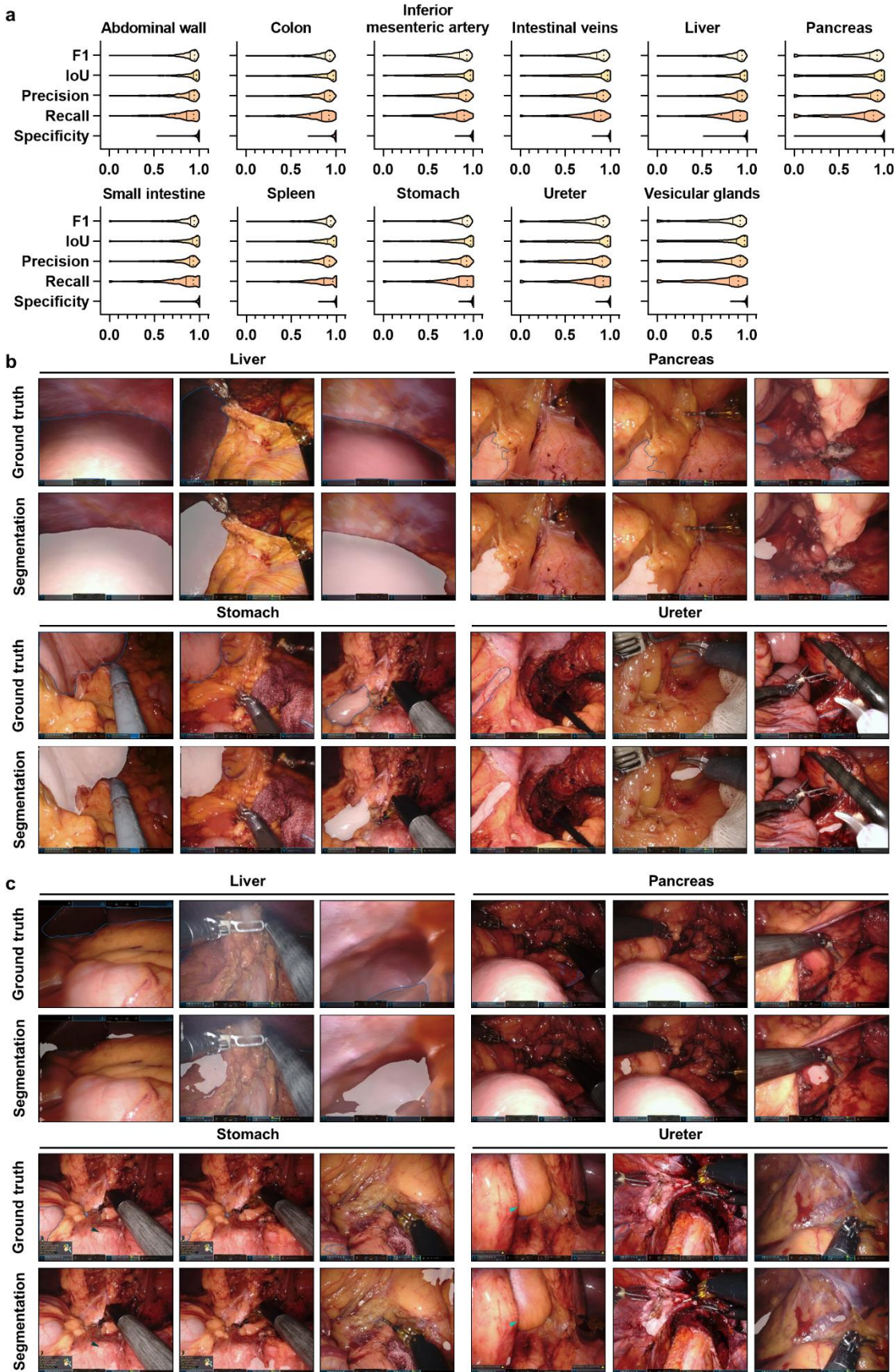
434 were trained: A structure-specific model with individual encoders and decoders, and a combined

435 model with a common encoder and structure-specific decoders.



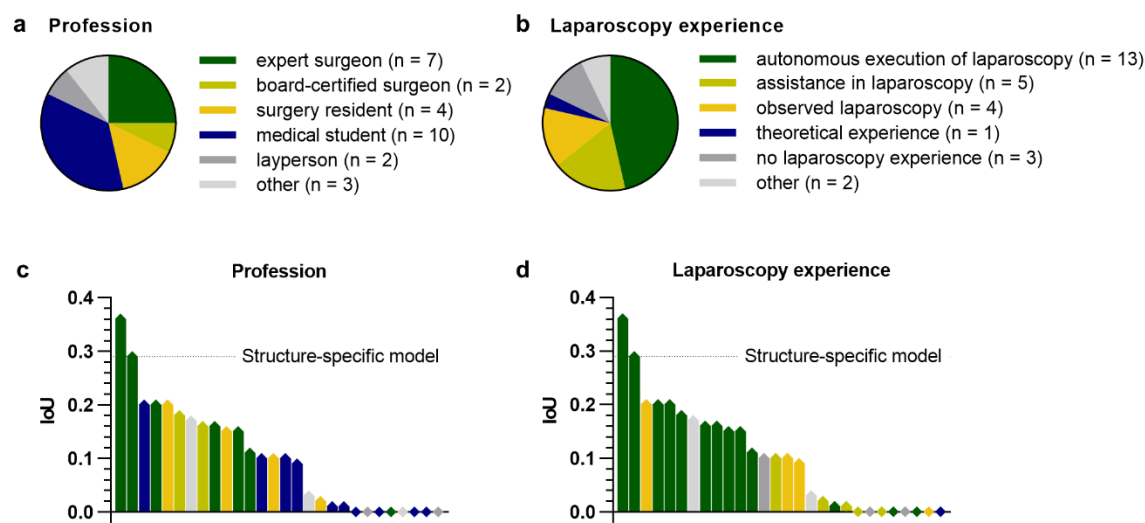
437 **Figure 2**

438 **Fig. 2 Pixel-wise organ segmentation with structure-specific models trained on the**
439 **respective organ subsets of the Dresden Surgical Anatomy Dataset. (a)** Violin plot
440 illustrations of performance metrics for structure-specific segmentation models. The median and
441 quartiles are illustrated as solid and dashed lines, respectively. **(b)** Example images with the
442 highest IoUs for liver, pancreas, stomach, and ureter segmentation with structure-specific
443 segmentation models. Ground truth is displayed in blue, and proposed segmentations are
444 displayed as white overlay. **(c)** Example images with the lowest IoUs for liver, pancreas, stomach,
445 and ureter segmentation with structure-specific segmentation models. Ground truth is displayed
446 in blue, and proposed segmentations are displayed as white overlay.
447



449 **Figure 3**

450 **Fig. 3 Pixel-wise organ segmentation with the combined model trained on the entire**
451 **Dresden Surgical Anatomy Dataset with a common encoder and structure-specific**
452 **decoders. (a)** Violin plot illustrations of performance metrics for the combined segmentation
453 model. The median and quartiles are illustrated as solid and dashed lines, respectively. **(b)**
454 Example images with the highest IoUs for liver, pancreas, stomach, and ureter segmentation with
455 the combined segmentation model. Ground truth is displayed in blue, and proposed
456 segmentations are displayed as white overlay. **(c)** Example images with the lowest IoUs for liver,
457 pancreas, stomach, and ureter segmentation with the combined segmentation model. Ground
458 truth is displayed in blue, and proposed segmentations are displayed as white overlay.
459



460

461 **Figure 4**

462 **Fig. 4 Comparison of pancreas segmentation performance of the structure-specific model**
463 **with a cohort of 28 human participants. (a)** Distribution of medical and non-medical professions
464 among human participants. **(b)** Distribution of laparoscopy experience among human participants.
465 **(c)** Waterfall chart displaying the average pancreas segmentation IoUs of participants with
466 different professions as compared to the IoU generated by the structure-specific model. **(d)**
467 Waterfall chart displaying the average pancreas segmentation IoUs of participants with varying
468 laparoscopy experience as compared to the IoU generated by the structure-specific model.
469

470 **Abbreviations**

471	AI	Artificial Intelligence
472	IoU	Intersection-over-Union
473	SD	Standard deviation

474

475 **Acknowledgements and Funding**

476 FRK, SL, JW, and SS were supported through project funding within the Else Kröner Fresenius
477 Center for Digital Health (EKFZ), Dresden, Germany (project “CoBot”). FRK received funding from
478 the Medical Faculty of the Technical University Dresden within the MedDrive Start program (grant
479 number 60487) and from the Joachim Herz Foundation (Add-On Fellowship for Interdisciplinary
480 Life Science). FMR received a doctoral student scholarship from the *Carus Promotionskolleg*
481 Dresden. The authors gratefully acknowledge excellent project coordination by Dr. Elisabeth
482 Fischermeier and Dr. Grit Krause-Jüttler.

483

484 **Author contributions**

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

490

491 **Competing interests**

492 The authors declare no conflicts of interest.

493