# Prediction of laparoscopic procedure duration using unlabeled, multimodal sensor data

**Sebastian Bodenstedt · Martin Wagner ·
Lars Mündermann · Hannes Kenngott ·
Beat Müller-Stich · Michael Breucha ·
Sören Torge Mees · Jürgen Weitz ·
Stefanie Speidel**

**Abstract** *Purpose* The course of surgical procedures is often unpredictable, making it difficult to estimate the duration of procedures beforehand. This uncertainty makes scheduling surgical procedures a difficult task. A context-aware method that analyses the workflow of an intervention online and automatically predicts the remaining duration would alleviate these problems. As basis for such an estimate, information regarding the current state of the intervention is a requirement.

*Methods* Today, the operating room contains a diverse range of sensors. During laparoscopic interventions, the endoscopic video stream is an ideal source of such information. Extracting quantitative information from the video is challenging though, due to its high dimensionality. Other surgical devices (e.g. insufflator, lights, etc.) provide data streams which are, in contrast to the video stream, more compact and easier to quantify. Though whether such streams offer sufficient information for estimating the duration of surgery is uncertain. In this paper, we propose and compare methods, based on convolutional neural networks, for continuously predicting the duration of laparoscopic interventions based on unlabeled data, such as from endoscopic image and surgical device streams.

S. Bodenstedt · S. Speidel
Department for Translational Surgical Oncology, National Center for Tumor Diseases (NCT), Partner Site Dresden, Dresden, Germany
E-mail: Firstname.Lastname@nct-dresden.de

M. Wagner · H. Kenngott · B. Müller-Stich
Department of General, Visceral and Transplant Surgery, University of Heidelberg, Heidelberg

L. Mündermann
KARL STORZ SE & Co. KG, Tuttlingen, Germany

M. Breucha · S.T. Mees · J. Weitz
Department of Visceral, Thoracic and Vascular Surgery, Faculty of Medicine and University Hospital Carl Gustav Carus, TU Dresden, Dresden, Germany

*Results* The methods are evaluated on 80 recorded laparoscopic interventions of various types, for which surgical device data and the endoscopic video streams are available. Here the combined method performs best with an overall average error of 37% and an average halftime error of approximately 28%.

*Conclusion* In this paper, we present, to our knowledge, the first approach for online procedure duration prediction using unlabeled endoscopic video data and surgical device data in a laparoscopic setting. Furthermore, we show that a method incorporating both vision and device data performs better than methods based only on vision, while methods only based on tool usage and surgical device data perform poorly, showing the importance of the visual channel.

**Keywords** Surgical Workflow Analyses · Duration Prediction · SensorOR

## 1 Introduction

The time in the operating room (OR) and the time of the operating staff are cost intensive hospital resources and have to be allocated precisely. Planning the usage of the OR cannot be static, as a procedure that takes longer than previously estimated can cause the following surgeries to be pushed back or even canceled, thereby inconveniencing both patient and the surgical team. On the other hand, if a procedure finishes early, the OR stays unused, incurring unnecessary idle time for the surgical personnel. To prevent this from occurring, OR schedulers need to dynamically update timetables. This is complicated by the unpredictability of the surgical workflow, which makes it difficult to estimate the duration of procedures beforehand. Therefore, the OR schedulers have to be constantly kept in the loop of the progress of ongoing interventions. For this, they have to periodically inquire into the status of interventions, resulting in highly subjective predictions and avoidable interruptions of procedures.

Integrated ORs are becoming more prevalent in hospitals, making it possible to access data streams from surgical devices such as cameras, insufflator, lights, etc. during interventions. Such data streams can provide information that enable context-aware assistance, such as automatically and continuously predicting the progress of an ongoing intervention. Especially the endoscopic video stream, via which laparoscopic interventions are performed, contains a large supply of information. Workflow analysis methods can be used to segment interventions into surgical phases. Often tool usage is employed to determine the current surgical phase [2,6,10,14], which provides an indicator for progress surgery. Convolutional neural networks (CNNs) have also been used to determine the surgical phase directly from the endoscopic video stream [4, 13,15].

Surgical phase detection methods can be used to approximate the duration of surgical procedures, but these methods generally require a sufficient amount of labeled examples as training input. Furthermore, seeing that phase models are generally specified to a certain type of intervention, multiple detectors would need to be trained. Therefore, using a phase based method as a general

solution to determine the remaining duration of surgeries would require an unfeasible large amount of labeled training data. In [7], the authors propose a system that determines the remaining time of surgery during laparoscopic cholecystectomies without surgical phases, but directly from the usage of the electrosurgical device. Two recurrent CNNs for predicting the remaining time of surgery directly from endoscopic video of cholecystectomies are presented in [1] and in [16]. While [16] also does not rely on the annotation of surgical phases, only visual features are used as input. In [3], we also propose a method for computing the duration of laparoscopic surgeries of varying types using a combination of a CNN and a gated recurrent unit [5] (GRU) with unlabeled video sequences.

In this paper, we propose and evaluate three methods, based on recurrent CNNs, for directly predicting and refining the duration of laparoscopic interventions. These methods do not require labeled training data and function for different types of laparoscopic procedures. The first method uses surgical device data collected from the OR as input, the second endoscopic image data. The third method is the combination of the previous methods. The evaluation of the methods is performed on a dataset containing 80 laparoscopic surgeries of varying types and on the Cholec80 dataset [15]. To our knowledge, our approach is the first method to predict the duration of laparoscopic interventions based on a combination of unlabeled vision and surgical device data.

## 2 Methods

A requirement for predicting the remaining duration of laparoscopic procedures is information regarding the current state of the surgery. As the endoscopic video stream serves as basis for the actions of the surgeon, we assume that it contains sufficient information on the state of the procedure, though extracting quantitative information is challenging due to the high dimensionality of the data stream. Here, we propose a recurrent CNN that allows predicting the duration of surgery from an endoscopic video stream. Furthermore, we propose a variation of this recurrent CNN that explicitly performs tool presence detection and utilizes this information to further enhance its predictions.

On the other hand, integrated ORs are starting to become more prevalent in hospitals. These ORs make it possible to access data streams, in the form of time series, from other surgical devices. These time series are, in contrast to the video stream, more compact and easier to quantify, but contain a smaller amount of information. We hypothesize that both streams, video and device data, contain complementary information, meaning that combining the two should increase prediction accuracy. To evaluate this claim, we propose a fusion of the two streams into one recurrent CNN (figure 1).
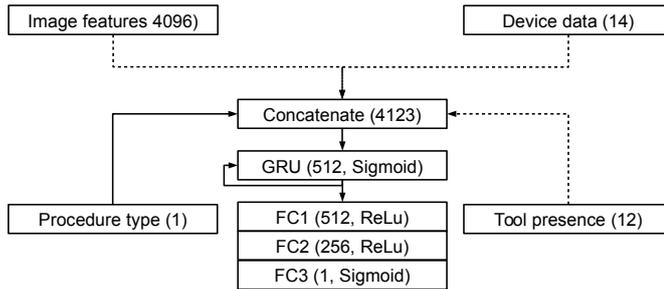
Fig. 1: The recurrent CNN topology for predicting surgery duration from images, device data and tool presence. The dotted lines indicate optional inputs. *FC* stands for a fully connected layer and *GRU* for a gated recurrent unit. The values in parenthesis indicate the number of hidden units and the nonlinearity used.

## 2.1 V-Net: Vision based estimation

Building upon our work in [3], the topology of our recurrent CNN for predicting duration of surgery from the endoscopic video, V-Net, is identical to the network for phase segmentation proposed in [4], which consists of an Alexnet [12] style network (image features in figure 1) combined with a GRU to incorporate temporal information. Only the final layers differs, we use a fully-connected layer with a single output and a sigmoid activation function. We pretrain the layers preceding the GRU with the method proposed in [4]. For this, we extract 2.2 million frames from 324 unlabeled laparoscopic procedures. The Alexnet is then trained in a Siamese fashion to extract features from images that allow sorting two random, given frames from the same procedure into the correct temporal order.

As input for V-Net, we sampled videos at a rate of one frame per second and downsampled each image to a resolution of $320 \times 240$ pixels. This is performed to reduce data size and training time.

Furthermore, we provide V-Net with information on the type of laparoscopic surgery being performed (procedure type in figure 1). For this, we categorize our dataset into 5 general types of laparoscopic surgery (table 1). We assign each frame a number between 0 and 1 as label, i.e. the label of frame $i$ from a video consisting of $N$ frames is $l_i = \frac{i}{N}$. During inference, we can then directly compute the duration prediction $\tilde{N}_i$: $\tilde{N}_i = \frac{i}{y_i}$, where $y_i \in [0, 1]$ is the predicted progress of the procedure.

Since the layers preceding the GRU are pretrained, we only optimize the weights of the newly added layers. For this, we use Adam [11] with a learning rate of $10^{-6}$ and a sigmoid cross-entropy loss. The network is trained for 50 epochs.

| ID | Surgery type | Samples in dataset | Average length (min.) |
|----|-------------|--------------------|-----------------------|
| 1 | Colorectal | 39 | 156 |
| 2 | Upper Gastrointestinal and Bariatric | 11 | 107 |
| 3 | Hepato-Pancreatico-Biliary | 4 | 102 |
| 4 | General Laparoscopic | 21 | 41 |
| 5 | Singular case | 5 | 91 |

Table 1: The categorization of laparoscopic surgery used as input. Also shown, the number of occurrences and average length in the **MultiType** dataset.

## 2.2 T-Net: Tool presence based estimation

Seeing that surgical tool presence in the endoscopic video is an important cue for progress of surgery [15], we hypothesize that information on the types of tools currently visible could be used to predict procedure duration. Tool usage could in principle be determined through external sensors such as RFID tags or a monitoring system. To simulate such a system and to identify the tools currently in use in the endoscopic video stream, we decided to apply a CNN trained on weakly labeled data were only tool usage is annotated in form of binary labels. Such data is easily annotated and is generally not specific to one type of surgical procedure.

For this, we modified a pretrained ResNet-152 [8] for tool presence detection. Here, we replace the pretrained fully connected output layer with another fully connected layer, consisting of 12 nodes (one per tool type) with sigmoid nonlinearities. The weights of the $conv5\_x$ layers and of the new fully connected layer are fine-tuned using a dataset of 24 colorectal laparoscopies in which the presence of 12 possible surgical tools has been labeled in one frame per minute. We used every labeled frame in the dataset to train the ResNet-152, overall 3800 frames. Adam [11] with a learning rate of $10^{-5}$ and binary cross-entropy are used. The network was then tested on a separate set of 9 further colorectal laparoscopies, consisting of 1543 annotated video frames. An average F1-score over all class of 81.2% was achieved.

To use tool presence information for predicting procedure duration, we use the output of the ResNet152 directly as input for the GRU described in the previous section, replacing the image features (tool presence in figure 1). T-Net is trained with the same hyper-parameters as V-Net.

## 2.3 D-Net: Surgical device data based estimation

The input consists of 14 values, each representing a different signal from a surgical device (a list of the devices and signals used can be found in table 2). Analog to the video stream, we select one second as the size of each time step. For signals with a data rate higher than 1 Hz, we discard all values, except the most recent one. For signals with a lower data rate, we use the most recent value, even if it was older than one second.

To predict the remaining duration of surgery, we first normalize the data, given the value ranges provided in table 2. The data is then fed into the previously described GRU (device data in figure 1). D-Net is trained with the same hyper-parameters as the previous networks.

## 2.4 TD-Net: Tool presence and surgical device data based estimation

To utilize tool presence information and surgical data for predicting the duration of surgery, we combine D-Net and T-Net into one architecture. This ID-Net is trained with identical hyper-parameters as the previous networks.

## 2.5 VT-Net: Vision and tool presence based estimation

To combine the tool presence information for predicting duration of surgery with the image features, we extend the architecture of V-Net into VT-Net, allowing it to additionally accept information on tool presence. VT-Net is trained with the same hyper-parameters as the previous networks.

## 2.6 VTD-Net: incorporating surgical device data

To incorporate the data stream provided by surgical devices, we extend VT-Net with a further input. VTD-Net is trained with the same hyper-parameters as the previous networks.

| Device | Signal | Data type | Value range |
|---|---|---|---|
| Insufflator | Current gas flow rate | Continuous | 0-215 |
| | Target gas flow rate | Continuous | 10-300 |
| | Current gas pressure | Continuous | 0-255 |
| | Target gas pressure | Continuous | 9-23 |
| | Used gas volume | Continuous | 0-9501 |
| | Gas supply pressure | Continuous | 0-760 |
| | Device on? | Binary | 0,1 |
| OR lights | All lights off? | Binary | 0,1 |
| | Intensity light 1 | Continuous | 0-100 |
| | Intensity light 2 | Continuous | 0-100 |
| Endoscopic light source | Intensity | Continuous | 0-100 |
| Endoscopic camera | White balance | Binary | 0,1 |
| | Gains | Continuous | 0-3298 |
| | Exposure index | Continuous | 0-834 |

Table 2: Surgical devices and signals used

## 3 Experiments and results

The basis of our evaluation is a dataset, containing recordings of 80 laparoscopic surgeries of different procedure types (**MultiType**). For each surgery, the dataset contains the endoscopic video stream and data collected from different surgical devices as listed in table 2. The procedures were all recorded in the same OR using the integrated operating room system OR1$^{\text{TM}}$ (Karl Storz GmbH & Co KG, Tuttlingen, Germany). The average procedure length in the dataset is 106 minutes. The datasets used for pretraining V-Net and training the ResNet-152 for tool presence detection do not overlap with this dataset.

To evaluate the proposed methods, we divide the dataset into four sets of equal size and perform a leave-one-set-out evaluation for each method. While dividing the dataset, we ascertain that the distribution of the different types of surgery into each set is balanced.

During testing, we compute the duration prediction $\tilde{N}_i$ at each frame $i$: $\tilde{N}_i = \frac{i}{y_i}$, where $y_i \in [0, 1]$ is the predicted progress of the procedure. With $\tilde{N}_i$, we can compute the absolute duration prediction error in seconds, $|\tilde{N}_i - N|$, and the duration prediction error relative to the length $N$ of of each procedure, $\frac{|\tilde{N}_i - N|}{N}$. The relative error gives a more appropriate impression on how well each method can predict procedure duration.

To put the performance of the proposed methods into perspective, we also propose different baselines for computing the remaining duration of a given surgery.

*Naive* As a **naive** baseline, we provide the duration prediction error that would occur if the average procedure duration over the training data were used as value for the predicted progress of the procedure.

*Type* We also provide a **type**-based baseline, where we instead compute the average procedure duration separately for each procedure category in table 1.

*Twinanda et al.* As a further baseline, we implemented the method proposed by Twinanda et al. [16] for predicting the remaining surgery duration. The method uses a ResNet152 to extract features from the images, which are then processed by a long short-term memory unit (LSTM) [9] to predict the remaining surgery time in minutes. As the operations in **MultiType** are significantly longer than the cholecystectomies in the dataset in [16], we set the parameter $s_{norm}$ to 20. Furthermore, the paper did not specify how many units the LSTM contained, we opted to use 512 units. Also the image size was not mentioned, we therefore used the default input resolution of ResNet152, which is $224 \times 224$ pixels. For training and testing **MultiType**, we performed a leave-one-set-out evaluation, training the ResNet152 and the LSTM in two steps as proposed in [16], though all the training data was used for both CNN and LSTM. We used 60000 iterations to train the LSTM. All other parameter values were set to equal values as described in [16].

For each method, we provide both the absolute and the relative average duration prediction error (see tables 3a and 3b). To measure how the error progresses during the course of a procedure, we compute the average error during each quarter of a given procedure (Q1-Q4).

| Method | Q1 | Q2 | Q3 | Q4 | Mean |
|---|---|---|---|---|---|
| V-Net | $3818 \pm 486$ | $2496 \pm 258$ | $1611 \pm 230$ | $1353 \pm 200$ | $2320 \pm 846$ |
| T-Net | $5078 \pm 1004$ | $1611 \pm 991$ | $1857 \pm 1000$ | $4316 \pm 1006$ | $3215 \pm 2005$ |
| D-Net | $4406 \pm 1315$ | $1928 \pm 788$ | $2301 \pm 910$ | $4482 \pm 1667$ | $3280 \pm 2373$ |
| ID-Net | $4473 \pm 1330$ | $1835 \pm 815$ | $2181 \pm 912$ | $4191 \pm 1613$ | $3170 \pm 2288$ |
| VT-Net | $4271 \pm 433$ | $2421 \pm 318$ | $1052 \pm 285$ | $1457 \pm 322$ | $2300 \pm 886$ |
| VTD-Net | $4143 \pm 449$ | $2289 \pm 312$ | $1071 \pm 279$ | $1313 \pm 312$ | $2204 \pm 875$ |
| Baseline (Naive) | $3208 \pm 2085$ | $3208 \pm 2085$ | $3208 \pm 2085$ | $3208 \pm 2085$ | $3208 \pm 2085$ |
| Baseline (Type) | $2093 \pm 1787$ | $2093 \pm 1787$ | $2093 \pm 1787$ | $2093 \pm 1787$ | $2093 \pm 1787$ |
| Twinanda et al. | $3862 \pm 531$ | $2427 \pm 508$ | $1405 \pm 463$ | $1828 \pm 507$ | $2380 \pm 1480$ |

(a) Absolute error (in seconds)

| Method | Q1 | Q2 | Q3 | Q4 | Mean |
|---|---|---|---|---|---|
| V-Net | $53\% \pm 9\%$ | $42\% \pm 6\%$ | $42\% \pm 5\%$ | $44\% \pm 4\%$ | $45\% \pm 16\%$ |
| T-Net | $74\% \pm 15\%$ | $23\% \pm 15\%$ | $28\% \pm 15\%$ | $79\% \pm 15\%$ | $51\% \pm 30\%$ |
| D-Net | $65\% \pm 20\%$ | $28\% \pm 13\%$ | $45\% \pm 15\%$ | $80\% \pm 22\%$ | $54\% \pm 37\%$ |
| TD-Net | $66\% \pm 20\%$ | $27\% \pm 13\%$ | $42\% \pm 16\%$ | $75\% \pm 22\%$ | $53\% \pm 36\%$ |
| VT-Net | $61\% \pm 8\%$ | $35\% \pm 6\%$ | $25\% \pm 5\%$ | $37\% \pm 6\%$ | $39\% \pm 15\%$ |
| VTD-Net | $59\% \pm 8\%$ | $32\% \pm 5\%$ | $24\% \pm 5\%$ | $35\% \pm 6\%$ | $37\% \pm 14\%$ |
| Baseline (Naive) | $124\% \pm 245\%$ | $124\% \pm 245\%$ | $124\% \pm 245\%$ | $124\% \pm 245\%$ | $124\% \pm 245\%$ |
| Baseline (Type) | $57\% \pm 99\%$ | $57\% \pm 99\%$ | $57\% \pm 99\%$ | $57\% \pm 99\%$ | $57\% \pm 99\%$ |
| Twinanda et al. | $65\% \pm 9\%$ | $54\% \pm 8\%$ | $51\% \pm 8\%$ | $66\% \pm 9\%$ | $59\% \pm 23\%$ |

(b) Relative error

Table 3: The average duration prediction errors on the **MultiType** dataset. (a) shows the mean absolute error of all the methods and (b) the mean relative error for all the methods.

For further evaluation, V-Net, T-Net and VT-Net are applied to the publicly available **Cholec80** dataset [15], which consists of 80 videos from laparoscopic cholecystectomies with an average length of 38 minutes. We divide the dataset into four sets of equal size and similar average procedure length and perform four leave-one-set-out evaluations for each method. The absolute and the relative average duration prediction errors can be found in tables 4a and 4b. The type-based baseline is not available, as the dataset consists only of a single procedure type. For the baseline based on Twinanda et al., all parameters were set to identical values as in [16], though we again used all the data in the training set to train both the ResNet152 and the LSTM. No surgical device data is contained in the dataset, meaning the networks relying on surgical device data cannot be used.

The results on **MultiType** show that all three of our methods that are based on the image features, outperform the baseline methods. All the image feature based methods also outperform the baseline based on Twinanda et al.,

| Method | Q1 | Q2 | Q3 | Q4 | Mean |
|---|---|---|---|---|---|
| V-Net | $934 \pm 238$ | $804 \pm 135$ | $581 \pm 119$ | $353 \pm 92$ | $668 \pm 233$ |
| T-Net | $1700 \pm 350$ | $500 \pm 333$ | $722 \pm 350$ | $735 \pm 351$ | $915 \pm 705$ |
| VT-Net | $920 \pm 147$ | $746 \pm 104$ | $434 \pm 93$ | $236 \pm 80$ | $584 \pm 214$ |
| Baseline (Naive) | $768 \pm 669$ | $768 \pm 669$ | $768 \pm 669$ | $768 \pm 669$ | $768 \pm 669$ |
| Twinanda et al. | $932 \pm 284$ | $619 \pm 208$ | $478 \pm 211$ | $324 \pm 217$ | $588 \pm 402$ |

(a) Absolute error (in seconds)

| Method | Q1 | Q2 | Q3 | Q4 | Mean |
|---|---|---|---|---|---|
| V-Net | $42\% \pm 12\%$ | $38\% \pm 7\%$ | $31\% \pm 6\%$ | $22\% \pm 5\%$ | $33\% \pm 11\%$ |
| T-Net | $74\% \pm 15\%$ | $22\% \pm 14\%$ | $31\% \pm 15\%$ | $35\% \pm 15\%$ | $41\% \pm 31\%$ |
| VT-Net | $34\% \pm 7\%$ | $28\% \pm 5\%$ | $17\% \pm 4\%$ | $12\% \pm 4\%$ | $23\% \pm 9\%$ |
| Baseline (Naive) | $66\% \pm 27\%$ | $66\% \pm 27\%$ | $66\% \pm 27\%$ | $66\% \pm 27\%$ | $66\% \pm 27\%$ |
| Twinanda et al. | $36\% \pm 13\%$ | $27\% \pm 9\%$ | $23\% \pm 11\%$ | $17\% \pm 11\%$ | $25\% \pm 17\%$ |

(b) Relative error

Table 4: The average duration prediction errors on the **Cholec80** dataset. (a) shows the mean absolute error of all the methods and (b) the mean relative error for all the methods.

also it can be noted that the methods presented here show a lower standard deviation of the error. One explanation why even V-Net performs better than Twinanda et al. could be due to the pretraining that was performed. Furthermore, VTD-Net provides more accurate results than the other two image feature based networks, demonstrating that surgical device data does indeed contain complementary information. Table 5 shows that VTD-Net performs consistently on the first three procedure categories. As can be seen in table 1, the surgeries in the first category are significantly longer than those in the second and third categories, though the relative error seems to stay consistent. Part of the general laparoscopic category are diagnostic laparoscopies, which are significantly shorter ($< 15$min) than the average procedure in the dataset and are less standardized and therefore difficult to predict, explaining this drop in performance. The last category contains only singular cases, which differ from the other categories, making predictions difficult.

The non-image based methods generally perform worse than the image based methods, leading us to conclude that while these data sources contain hints relevant for the progress of surgery that are useful for augmenting the image features, by themselves they don't contain sufficient information to accurately predict the progress of surgery.

On **Cholec80**, both image feature based methods outperform the naive baseline, while VT-Net also outperformed V-Net. The baseline based on Twinanda et al. outperformed V-Net and performed similarly to VT-Net, thought the standard deviation of the error of both V-Net and VT-Net was smaller. The method of Twinanda et al. seems to perform better at the beginning of surgery, while VT-Net performs better in the second half of surgery. Similarly as for **MultiType**, the method solemnly based on tool usage performs poorly. A

direct comparison with the results in [1] is not possible, as the authors use a private dataset for testing.

| Type | Q1 | Q2 | Q3 | Q4 | Mean |
|------|----|----|----|----|------|
| 1 | $6020 \pm 961$ | $3452 \pm 733$ | $1350 \pm 652$ | $1381 \pm 710$ | $3050 \pm 2136$ |
| 2 | $4102 \pm 863$ | $2096 \pm 541$ | $852 \pm 474$ | $1226 \pm 527$ | $2069 \pm 1460$ |
| 3 | $4070 \pm 822$ | $2090 \pm 510$ | $764 \pm 376$ | $1053 \pm 475$ | $1994 \pm 1483$ |
| 4 | $1098 \pm 457$ | $496 \pm 241$ | $706 \pm 246$ | $1270 \pm 292$ | $893 \pm 572$ |
| 5 | $1419 \pm 522$ | $707 \pm 308$ | $927 \pm 256$ | $1518 \pm 382$ | $1143 \pm 746$ |

(a) Absolute error (in seconds)

| Type | Q1 | Q2 | Q3 | Q4 | Mean |
|------|----|----|----|----|------|
| 1 | $67\% \pm 12\%$ | $37\% \pm 8\%$ | $15\% \pm 7\%$ | $18\% \pm 8\%$ | $34\% \pm 24\%$ |
| 2 | $62\% \pm 14\%$ | $30\% \pm 9\%$ | $12\% \pm 7\%$ | $22\% \pm 9\%$ | $32\% \pm 22\%$ |
| 3 | $65\% \pm 14\%$ | $32\% \pm 8\%$ | $12\% \pm 6\%$ | $18\% \pm 8\%$ | $32\% \pm 24\%$ |
| 4 | $43\% \pm 20\%$ | $25\% \pm 12\%$ | $45\% \pm 13\%$ | $72\% \pm 14\%$ | $46\% \pm 28\%$ |
| 5 | $45\% \pm 21\%$ | $26\% \pm 13\%$ | $46\% \pm 11\%$ | $71\% \pm 15\%$ | $47\% \pm 30\%$ |

(b) Relative error

Table 5: The average duration prediction errors of VTD-Net on the **Multi-Type** dataset, broken down according to procedure type. (a) shows the mean absolute error and (b) the mean relative error.

## 4 Discussion

In this paper, we presented, to our knowledge, the first approach for online procedure duration prediction using unlabeled endoscopic video data and surgical device data in a laparoscopic setting. On **MultiType**, VTD-Net showed an overall average prediction error of 37% and a halftime error of about 28%, which are lower than the results from the baselines. Furthermore, we showed that a method incorporating both vision and device data performs better than methods based only on vision, while methods only based on tool usage and surgical device data performed poorly, showing the importance of the visual channel. The evaluation showed that the presented methods currently produce larger than average errors on irregular procedures with short length (shorter than 15 min) and on singular cases. As the methods performed better on **Cholec80**, which contains mostly shorter and more standardized operations, we assume this is due to a lack of training data of the irregular cases. When looking at complex procedures of medium length (second and third category) and long length (first category), the relative error stays consistent.

On **Cholec80**, VT-Net achieved an average prediction error as well as a halftime error of 23%, outperforming the baseline. This difference of performance on the two datasets can be contributed to the fact that **MultiType** contained multiple procedure types and also had a higher variance in procedure duration. Seeing that the results of the proposed methods perform significantly

better on a simple operation like cholecystectomies than on a more complex and diverse operations as in the **MultiType** dataset leads us to conclude that the proposed methods themselves are sound, but that more training data is required for more complex cases.

As our results indicate that combing vision and device data provides more information on the progress of surgery, we assume that data available to the anesthetist, such as heart rate, blood pressure and drug doses, would provide even more valuable insights.

**Conflict of interest** S. Bodenstedt, M. Wagner, L. Mündermann, H. Kenngott, B. Müller-Stich, M. Breucha, S. Mees, J. Weitz and S. Speidel declare that they have no conflict of interest.

**Ethical approval** All procedures performed in studies involving human participants were in accordance with the ethical standards of the institutional and/or national research committee and with the 1964 Declaration of Helsinki and its later amendments or comparable ethical standards

**Informed consent** Informed consent was obtained from the study participants

# References

1. Aksamentov, I., Twinanda, A.P., Mutter, D., Marescaux, J., Padoy, N.: Deep neural networks predict remaining surgery duration from cholecystectomy videos. In: International Conference on Medical Image Computing and Computer-Assisted Intervention, pp. 586–593. Springer (2017)
2. Blum, T., Feußner, H., Navab, N.: Modeling and segmentation of surgical workflow from laparoscopic video. In: International Conference on Medical Image Computing and Computer-Assisted Intervention, pp. 400–407. Springer (2010)
3. Bodenstedt, S.: Image-based scene analysis for computer-assisted laparoscopic surgery. Ph.D. thesis, Karlsruhe Institute of Technology (2018). DOI 10.5445/IR/1000084137
4. Bodenstedt, S., Wagner, M., Katić, D., Mietkowski, P., Mayer, B., Kenngott, H., Müller-Stich, B., Dillmann, R., Speidel, S.: Unsupervised temporal context learning using convolutional neural networks for laparoscopic workflow analysis. ArXiv e-prints (2017)
5. Cho, K., van Merrienboer, B., Bahdanau, D., Bengio, Y.: On the properties of neural machine translation: Encoder–decoder approaches. In: Proceedings of SSST-8, Eighth Workshop on Syntax, Semantics and Structure in Statistical Translation, pp. 103–111. Association for Computational Linguistics (2014). DOI 10.3115/v1/W14-4012. URL http://aclweb.org/anthology/W14-4012
6. Dergachyova, O., Bouget, D., Huaulmé, A., Morandi, X., Jannin, P.: Automatic data-driven real-time segmentation and recognition of surgical workflow. International journal of computer assisted radiology and surgery **11**(6), 1081–1089 (2016)
7. Guédon, A.C.P., Paalvast, M., Meeuwsen, F.C., Tax, D.M.J., van Dijke, A.P., Wauben, L.S.G.L., van der Elst, M., Dankelman, J., van den Dobbelsteen, J.J.: 'it is time to prepare the next patient' real-time prediction of procedure duration in laparoscopic cholecystectomies. Journal of Medical Systems **40**(12), 271 (2016). DOI 10.1007/s10916-016-0631-1
8. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 770–778 (2016)
9. Hochreiter, S., Schmidhuber, J.: Long short-term memory. Neural computation **9**(8), 1735–1780 (1997)
10. Katić, D., Wekerle, A.L., Gärtner, F., Kenngott, H., Müller-Stich, B.P., Dillmann, R., Speidel, S.: Knowledge-driven formalization of laparoscopic surgeries for rule-based intraoperative context-aware assistance. In: International Conference on Information Processing in Computer-Assisted Interventions, pp. 158–167. Springer (2014)
11. Kingma, D., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)

12. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: Advances in neural information processing systems, pp. 1097–1105 (2012)
13. Lea, C., Choi, J.H., Reiter, A., Hager, G.: Surgical phase recognition: From instrumented ors to hospitals around the world. M2CAI 2016 (2016)
14. Padoy, N., Blum, T., Ahmadi, S.A., Feussner, H., Berger, M.O., Navab, N.: Statistical modeling and recognition of surgical workflow. Medical image analysis **16**(3), 632–641 (2012)
15. Twinanda, A.P., Shehata, S., Mutter, D., Marescaux, J., de Mathelin, M., Padoy, N.: Endonet: A deep architecture for recognition tasks on laparoscopic videos. IEEE transactions on medical imaging **36**(1), 86–97 (2017)
16. Twinanda, A.P., Yengera, G., Mutter, D., Marescaux, J., Padoy, N.: Rsdnet: Learning to predict remaining surgery duration from laparoscopic videos without manual annotations. IEEE Transactions on Medical Imaging pp. 1–1 (2018). DOI 10.1109/TMI.2018.2878055