2020　2021

Removing the Burden of Data Labeling: Automatic Surgical Video Understanding with Unsupervised Learning

LI, LIANGZHI

1,100,000

1.

2.

The establishment of techniques capable of understanding what, where, and when is happening with no requirements on human efforts will ease the task of surgery indexing, clinical training, etc.

As most of the existing automatic surgical video analysis models require a large number of manually labeled data for training, this project aims to design a learning method to perform spatial and temporal segmentations with smaller requirements of humans' input. During this project, I mainly studied the following sub-topics towards the goal.
1. Surgical images/frames semantic segmentation in a weakly-supervised way. I developed a new training strategy for video semantic segmentation models to utilized unlabeled data to improve their segmentation performance.
2. Surgical videos temporal analysis using no labels. I developed a retrieval-based method to automatically predict surgical duration.

Computer Vision

Few-shot Learning  Semantic Segmentation  Video Understanding  Medical Images  Computer Vision  Surgical Analysis  Deep Learning

１．研究開始当初の背景

　　Well-analyzed surgical videos are of great value for not only training clinical students but also helping professionals. For that purpose, there have been lots of works on the **automatic analysis of surgical videos**. Encouraging results have been achieved in the understanding of the spatial information (what and where are the objects) as well as the temporal information (what and when is happening). However, it remains **a challenging unsolved problem** when no enough ground-truth labels are available, as most of the existing models belong to *supervised learning*, which needs labels to teach the deep learning models, as shown in Figure 1. This fact has extremely limited the broad application of the automatic surgical video analysis system. Designing a new
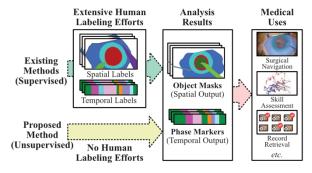


Figure 1. Supervised and unsupervised methods

computer vision (CV) method which can remove the labeling burden is **the first step to enable low-cost automatic surgery analysis**. Therefore, the key scientific question this project aims to answer is, how to extract spatial and temporal information from surgical videos with no labeled data available.

２．研究の目的

　　The objective of this research is to make some progress towards enabling the unsupervised spatial understanding (USU) and unsupervised temporal understanding (UTU) for surgical videos. As shown in Figure 2, we will first design a new USU method using hidden temporal information (*local consistency* for generating object proposals, *i.e.*, finding the pixels which always move together and thus may be an independent object, and *motion clues* for
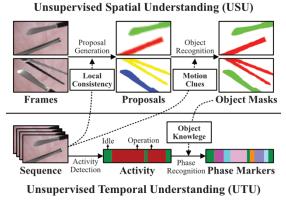


Figure 2. The proposed unsupervised method

differentiating and merging proposals into object masks, *i.e.*, which proposals should be same object). Then we will address the existing problems in UTU methods with the obtained spatial information (*object knowledge* for introducing more temporal differences, *i.e.*, which instruments or tissues are involved). The research objectives (RO) include:

- **RO1**: Enable the USU of surgery scenes using temporal clues.
- **RO2**: Utilize spatial knowledge to improve UTU performance.

３．研究の方法

There are some attempts in CV aiming at performing video analysis without ground-truth labels. These *unsupervised learning* methods **need no human-made labels**. However, few researchers have adopted them in the medical area, as the objects in medical images are either too similar to differentiate, which is <u>a key obstacle to USU</u>, or too small to cause significant changes which are <u>required for UTU</u>. In fact, medical images do have some inherent information to make *unsupervised learning* possible. For example, the objects in medical images usually have <u>unique movements which are much easier to recognize</u>, *e.g.*, some forceps and scissors may have similar appearances but with different movements (grasping/cutting); the exists of some objects may also serve as <u>a sign to infer the possible surgical phases</u> because using different combinations of instruments usually means different surgical phases. This kind of spatial-temporal information is pervasive among various surgery types and can help to enable automatic surgery understanding and **has not yet been studied**.

Therefore, to make progress towards enabling automatic surgery video analysis, this project is divided into three modules.

- **M1**: We will design a new USU model for object segmentation.
- **M2**: We will propose a new UTU model for phase segmentation.
- **M3**: We plan to develop an automatic segmentation software to output spatial and temporal results.

*The first 12 months* will be equally and sequentially dedicated to M1 and M2, respectively for the object and phase segmentation. We will make prototype implementations and test their performance. For *the last 6 months* in M3, we will develop software for surgery analysis using the proposed method. Also, we are going to submit two journal/conference papers to report our work in M1 and M2.

４．研究成果

(1) Weakly supervised semantic segmentation

Semantic video segmentation is a key challenge for various applications. As shown in Figure 3, I present a new model named Noisy-LSTM, which is trainable in an end-to-end manner, with convolutional LSTMs (ConvLSTMs) to leverage the temporal coherence in video frames, together with a simple yet effective training strategy that replaces a frame in a given video sequence with noises. The training strategy spoils the temporal coherence in video frames and thus makes the temporal
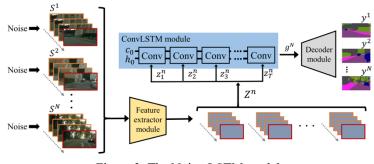


Figure 3. The Noisy-LSTM model.

links in ConvLSTMs unreliable; this may consequently improve the ability of the model to extract features from video frames and serve as a regularizer to avoid overfitting, without

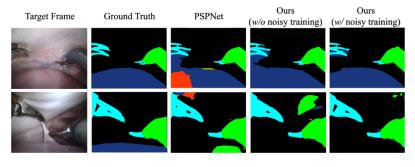requiring extra data annotations or computational costs. Experimental results demonstrate that the proposed model can achieve state-of-the-art performances on both the CityScapes and EndoVis2018 (as shown in Figure 4) datasets.



Figure 4. Segmentation performance on surgical videos.

(2) Video temporal segmentation and duration estimation

Estimating the surgery length has the potential to be utilized as skill assessment, surgical training, or efficient surgical facility utilization especially if it is done in real-time as a remaining surgery duration (RSD). Surgical length reflects a certain level of efficiency and mastery of the surgeon in a well-standardized surgery such as cataract surgery. Real-time RSD estimation can enable optimization of operating room scheduling as well. In this paper, we design and develop a real-time RSD estimation method for cataract surgery that does not require manual labeling.
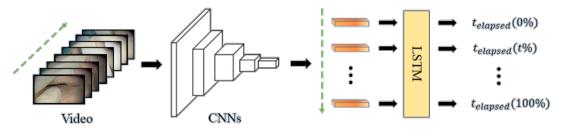


Figure 5. Surgery temporal segmentation and duration estimation model.

As shown in Figure 5, a regression method consisting of convolutional neural networks (CNNs) and long short-term memory (LSTM) is designed for RSD estimation. The model is firstly trained and evaluated for the single main target surgeon. Then, the fine-tuning strategy is used to transfer the model to the data of the other two surgeons. Mean average error (MAE in seconds) was used to evaluate the performance of the RSD estimation. The proposed method is compared with the naïve method which is based on the statistic of the historical data. A transferability experiment is also set to demonstrate the generalizability of the method.
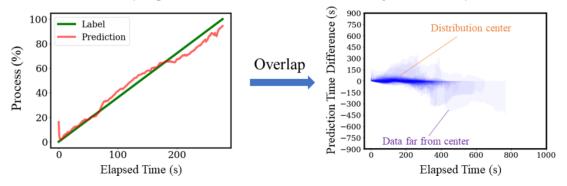


Figure 6. Surgery duration estimation performance.

As shown in Figure 6, the model shows low prediction error. The mean surgical time for the sample videos was 318.7 seconds (s) (standard deviation 83.4 s) for the main surgeon for the initial training. In our experiments, the lowest MAE of 19.4s (equal to about 6.4% of the mean surgical time) is achieved by our best-trained model for the independent test data of the main target surgeon. It reduces the MAE by 35.5s (-10.2%) compared to the naïve method. The fine-tuning strategy transfers the model trained for the main target to the data of other surgeons with only a small number of training data (20% of the pre-training). The MAEs for the other two surgeons are 28.3s and 30.6s with the fine-tuning model, which decreased by -8.1s and -7.5s than the Per-surgeon model (average declining of -7.8 s and 1.3% of video duration).

(3) Few-shot Learning

Few-shot learning (FSL) approaches, mostly neural network-based, are assuming that the pre-trained knowledge can be obtained from base (seen) categories and transferred to novel (unseen) categories. However, the black-box nature of neural networks makes it difficult to understand what is actually transferred, which may hamper its application in some
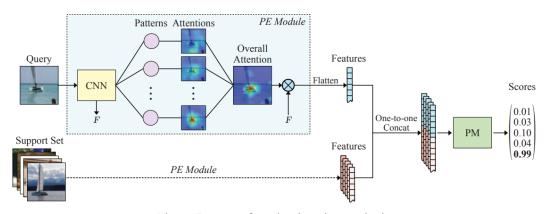


Figure 7. A new few-shot learning method.

risk-sensitive areas. As shown in Figure 7, I reveal a new way (MTUNet) to perform explainable FSL for image classification, using discriminative patterns and pairwise matching.

MTUNet learns discriminative patterns for a given set of images of the base categories as shown in Figure 7 and uses all these patterns to represent both support and query images. With this representation, pairwise matching scores are computed among the support and query images, based on which the prediction for the query image is done. Both the patterns and the overall representation can be easily visualized to reveal the reason for the matching scores. Experimental results prove that the proposed method can achieve satisfactory explainability (as shown in Figure 8).



Figure 8. FSL performance.

| | | | | | |
|---|---|---|---|---|---|
| 1 | 1 | 1 | 1 | | |
| Wang Bowen  Li Liangzhi  Nakashima Yuta  Kawasaki Ryo  Nagahara Hajime  Yagi Yasushi | | | | | 9 |
| Noisy-LSTM Improving Temporal Awareness for Video Semantic Segmentation | | | | | 2021 |
| IEEE Access | | | | | 46810  46820 |
| DOI<br>10.1109/ACCESS.2021.3067928 | | | | | |
| | | | | | |

| | | |
|---|---|---|
| 2 | 0 | 2 |
| Wang Bowen  Li Liangzhi  Manisha Verma  Nakashima Yuta  Kawasaki Ryo  Nagahara Hajime | | |
| MTUNet: Few-shot Image Classification with Visual Explanations | | |
| IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Responsible Computer Vision Workshop | | |
| 2021 | | |

| |
|---|
| Liangzhi Li, Bowen Wang, Manisha Verma, Yuta Nakashima, Ryo Kawasaki, Hajime Nagahara |
| SCOUTER: Slot Attention-based Classifier for Explainable Image Recognition |
| IEEE/CVF International Conference on Computer Vision (ICCV) |
| 2021 |

0

Github Code for the Semantic Segmentation
https://github.com/wbw520/NoisyLSTM
Github Code for the Few-shot Learning
https://github.com/wbw520/MTUNet

O