

ESAD: Endoscopic Surgeon Action Detection Dataset

Vivek Singh Bawa^{c,1}, Gurkirt Singh^{d,1}, Francis Kaping^A^c, Inna Skarga-Bandurova^c, Alice Leporini^b, Carmela Landolfo^b, Armando Stabile^b, Francesco Setti^a, Riccardo Muradore^a, Elettra Oleari^b, Fabio Cuzzolin^c

^a*University of Verona, Italy*

^b*San Raffaele Hospital, Milan, Italy*

^c*Oxford Brookes University, Oxford, UK*

^d*Computer Vision Lab, ETH Zurich, Switzerland*

Abstract

In this work, we take aim towards increasing the effectiveness of surgical assistant robots. We intended to make assistant robots safer by making them aware about the actions of surgeon, so it can take appropriate assisting actions. In other words, we aim to solve the problem of surgeon action detection in endoscopic videos. To this, we introduce a challenging dataset for surgeon action detection in real world endoscopic videos. Action classes are picked based on the feedback of surgeons and annotated by medical professional. Given a video frame, we draw bounding box around surgical tool which is performing action and label it with action label. Finally, we present a frame-level action detection baseline model based on recent advances in object detection. Results on our new dataset show that our presented dataset provides enough interesting challenges for future method and it can serve as strong benchmark corresponding research in surgeon action detection in endoscopic videos.

Keywords: Action detection, endoscopic video, surgeon action detection, prostatectomy, surgical robotics

Email address: vsingh@brookes.ac.uk (Vivek Singh Bawa)

¹Authors have equal contribution

1. Introduction

Minimally Invasive Surgery (MIS) is a very sensitive medical procedure. A general MIS surgical procedure involves two surgeons: main surgeon and assistant surgeon. Success of a MIS procedure depends upon multiple factors, such as, attentiveness of main surgeon and assistant surgeon, competence of surgeons, effective coordination between the main surgeon and assistant surgeon etc.

According to Lancet Commission, each year 4.2 million people die within 30 days of surgery [18]. Another study at John Hopkins University states that 10% of total deaths in USA are due to medical error [17]. There is no definite measure to compute and predict the risk factor involving surgeons. This make it very critical to monitor the set of action performed by surgeons in real time, so that, any unfortunate event can be avoided.

Artificial Intelligence is being used in a lot of applications where human error has to be mitigated. The proposed dataset is another step in same direction. To make the surgical procedure safer, we should be able to identify and track the actions of main as well as assistant surgeon. This dataset is developed with the assistance of medical professionals as well as expert surgeon. More details of the data set can be found in section 4.

Although there a lot of datasets for different action detection task computer vision. But there is no existing dataset for action detection in medical computer vision, specifically for MIS surgeries. Given the complexity of the scene and difficulty in the detection of surgeon action, this dataset will pave a path forward and set a benchmark for the medical computer vision research community. The task of action detection in medical scenario is a lot different from the general computer vision (more discussion in section 3), hence all the standard computer vision algorithm can not be directly deployed. The dataset will also lay the foundation for more robust algorithms which will be used in future surgical systems to accomplish tasks, such as, autonomous assistant surgeons, surgeon feedback systems, surgical anomaly detection, and so on.

2. Literature review

Action detection or activity analysis for medical images is an under explored field. Hence most of the literature in this field will be borrowed from general activity analysis literature. The earlier works like [20] use hand motion of the surgeon to recognize the action preformed. Voros *et al.* [26] uses

motion of tools to detect the point of interaction between tool and the organ. Kocev *et al.* [11] uses point cloud generated using Microsoft Kinect camera to build the augmented reality model of real time actions performed by surgeon. In [25], authors used weakly supervised approach based on Gaussian Mixture Models (GMM) to recognize the surgeon actions. The approach was not developed for real surgical images and only recognised actions with assumption of one action per frame. Azari *et al.* [1] use video of surgeon hand motion to predict the surgical maneuvers. Li *et al.* [13] uses sub-action categories for early stage prediction of main surgical actions.

Most of the literature we will be borrowing from the human activity detection problem. In general, there are two types of activity analysis methods: static and dynamic. Static methods only have spatial information (image data) without any temporal context to current frame [23, 3, 19]. The dynamic activity detection methods use video data which give temporal context to the motion or structure under the observation [22, 9, 6, 7].

Singh *et al.* [23] used Single Shot multi-box Detector (SSD) [16] to detect the activity in the frame. SSD is a very successful algorithm in object detection which predicts the object bounding boxes in a single shot, making it one of the fastest detection algorithms available. [3] used RCNN for the region proposal and these proposals are used to learn the context information to produce a more accurate activity class. Saha *et al.* [22] proposed activity detection module called as 3D-RPN (3 dimensional region proposal network) which uses spatial as well as temporal information from the same sequence. The model takes two different frames, from the same action sequence, separated by Δt time to learn the temporal context to the current frame.

Tian *et al.* [24] uses deformable part based model [2] to detect the activity in the action frame. Peng *et al.* [19] developed a motion region proposal network which was based on faster RCNN [21]. two streams (images and optical flow) were used in faster RCNN to generate the activity proposals. Jain *et al.* [8] use super-voxels to generate activity bounding boxes. The paper produces 2D+t bounding boxes with selective search sampling from the videos. Kalogeiton [9] and [7] develop action tube based methods. Both of the models predict the action tube which provides spatial bounding boxes in each of the frame from start to end of an action in a video. Li *et al.* [12] proposed Recurrent Tubelet Proposal and Recognition (RTPR) networks to predict action tubes from start to end of the action in video. The model has two networks, one for proposal and other for recognition. Combination of Con-

volutional Neural Network (CNN) and Long Short Term Memory (LSTM) Network to learn the recurrent nature of the action proposals.

3. Problem statement

The task of surgeon action detection is very novel and complex. The key factor that makes this task different from other activity detection task is the appearance of the surgical scene. The most dominant issues with medical images are:

- Most important contributor to the difficulty is deformable nature of the organs. As shown in figure 1, the organs do not hold a fixed shape in contrast to human activity detection problem, where body has fixed shape and shows a identifiable position. Additionally, boundaries, shapes and color variance between two different organs is minimal, making it very different from standard computer vision tasks.
- The scene captured in using endoscope camera is in very close proximity, hence it is unable to show complete organs or its surroundings. hence there is very little contextual information. General activity dataset like Kinetic [10] or AVA [4] have color, texture, shape and context information making it easier to learn the scene features.
- Motion and orientation of endoscope in near proximity makes organs appear very different from different angles.
- The set of action defined in this dataset provide very accurate description of surgeon actions (e.g., CuttingMesocolon, PullingProstate etc.) to make prediction more informative and useful. Hence in presented dataset, it becomes highly important to know the organ under operation to accurately detect and predict the action.

4. ESAD Dataset

4.0.1. Annotation protocol

A set of protocols is developed to guide the annotators in their work. This helped minimise the ambiguity in deciding the size of bounding boxes around each action instance, as well as their locations. All annotators were provided a set of instructions with examples in order to standardise the procedure as much as possible. The following guidelines were enforced for the annotation:

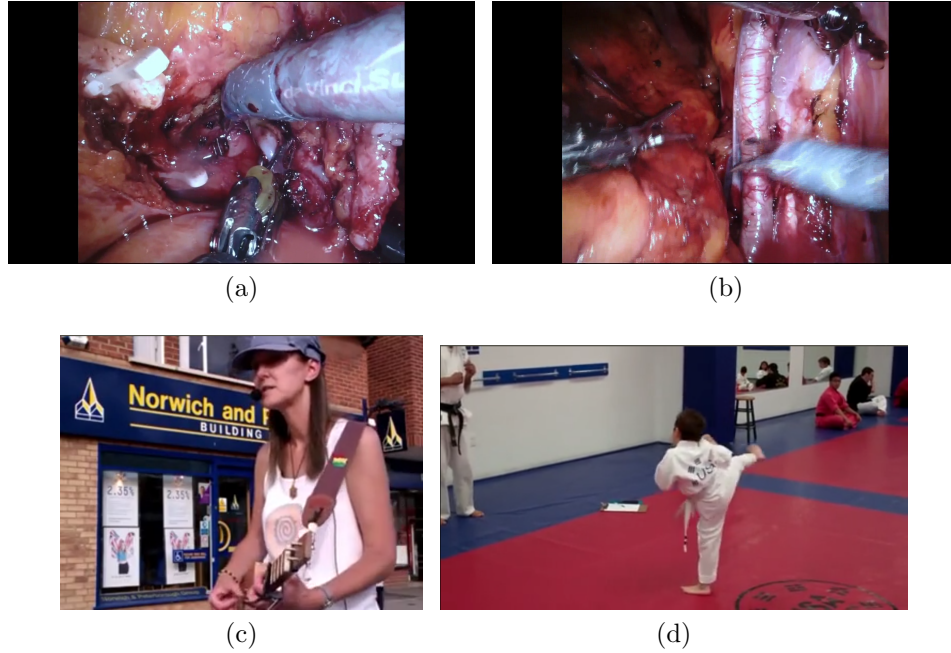


Figure 1: Image 1a and 1b are samples from ESAD dataset. Image 1c and 1d are samples from Kinetic-400 dataset [10], which is built from video on YouTube. The difference between the ESAD and Kinetic [10] dataset are evident. Endoscope video in ESAD dataset captures images from very close distance losing all contextual information unlike human activity videos at YouTube. Additionally, general activity dataset like Kinetic [10] or AVA [4] have color, texture, shape and context information making it easier to learn the scene features.

- Each bounding box should contain both the organ and tool performing the action under consideration, as each action class is highly dependent on the organ under operation.
- To balance the presence of tools and organs or tissue in a bounding box, bounding boxes are restricted to containing 30%-70% of either tools or organs.
- An action label is only assigned when a tool is close enough to the appropriate organ, as informed by the medical expert. Similarly, an action stops as soon as the tool starts to move away from the organ.
- Each video frame can have two actions, whose bounding boxes are

allowed to overlap.

4.0.2. Structure of dataset

After a rigorous analysis of the actions performed by surgeons during a typical prostatectomy procedure, we selected 21 action categories for ESAD dataset. Decision is made keeping in mind that action categories should not be too simple that they do not contribute any useful information. Similar, problem is faced by in previous medical action recognition datasets [1, 20]. Furthermore, the action classes should not be too complex, making it impossible to model the task. We concluded the action class list with the help of multiple surgeons and medical professionals. Detailed list of classes is shown in table 1 along with number of action instances for each category in the whole dataset.

For the creation of ESAD dataset, we collected four complete prostatectomy procedure with the consent of the patients and hospital. On an average, each video is 2 hours 20 minutes long. Each video is recorded videos at 30 FPS but annotation are performed at 1 FPS to maintain sufficient variation in the scene. Each frame can have multiple number of action instances. Each instance is annotated with a bounding box and its action label from classes. Surgeon and medical professionals were involved in making the decisions on the appearance of action classes as well as the area of the bounding boxes. As tools operate in close proximity, dataset have a lot of action instances with overlapping bounding boxes. Each annotation is verified by a medical professional.

Some sample from the ESAD dataset are shown in figure 2. In the images, it can be seen that bounding boxes are centred around the tool as laparoscopic tools represent the subject of action, but dataset also makes sure to include organ under operation. The reasons for that is that most of the surgical action have different names depending on the organ they are operating on despite of the same motion of tools.

4.0.3. Dataset split

The dataset is divided into three different sets: training, validation and test. The two surgeries with with the maximum number of action instances were selected as training set. The one video with the most balanced number of samples for each action class was used as the test set. The objective of dataset is to provide a fair evaluation to all types of algorithms. Hence,

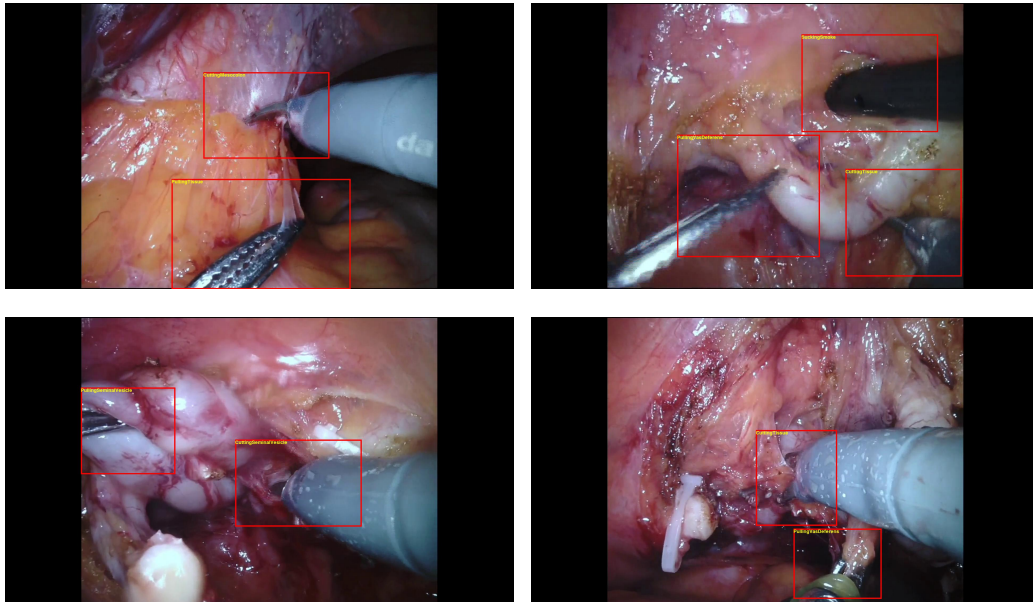


Figure 2: Samples from the ESAD dataset with annotations for surgeon actions. Red box in figure denotes the bounding box. The images are video output from the endoscope during prostatectomy procedure.

the last procedure is selected for validation set. The number of instances (labelled bounding boxes) for each action class can be seen in Table 1.

Distribution of samples for each action category for each of the three splits is shown in figure 3. It is clear from the bar chart the dataset is highly skewed in term of class imbalance. The reason for this is the nature of surgical procedures. As shown in figure 3, classes *PullingTissue* and *CuttingTissue* contain highest number of samples as this is the most common action performed by surgeon during prostatectomy. Whereas, classes like *BaggingProstate* and *CuttingThread* have lowest samples due to short duration of these activities per procedure.

5. Results and discussion

In this section, we start by presenting a baseline model 5.1 used to establish a baseline on our dataset. We show results of baseline mode in 5.3 section with evaluation metric described in 5.2. Finally, we will discuss the

Label	Train	Val	Test	Total instances
CuttingMesocolon	315	179	188	682
PullingVasDeferens	457	245	113	815
ClippingVasDeferens	33	25	48	106
CuttingVasDeferens	71	22	36	129
ClippingTissue	215	44	15	274
PullingSeminalVesicle	2712	342	436	3490
ClippingSeminalVesicle	118	35	33	186
CuttingSeminalVesicle	2509	196	307	3012
SuckingBlood	3753	575	1696	6024
SuckingSmoke	381	238	771	1390
PullingTissue	4877	2177	2024	9078
CuttingTissue	3715	1777	2055	7547
BaggingProstate	34	5	37	76
BladderNeckDissection	1621	283	519	2423
BladderAnastomosis	3585	298	1828	5711
PullingProstate	958	12	451	1421
ClippingBladderNeck	151	24	18	193
CuttingThread	108	22	40	170
UrethraDissection	351	56	439	846
CuttingProstate	1845	56	48	1949
PullingBladderNeck	189	509	105	803

Table 1: List of actions for ESAD dataset with number of samples for training, validation and test.

understanding gained from the results 5.3 in ??.

5.1. Baseline model

The baseline model is based on Feature Pyramidal Network (FPN) architecture. The concept was originally proposed by Lin *et al.* [14]. The paper uses convolutional CNN architecture with pooling layers. Residual networks (ResNet) [5] is used as a backbone network for the detection model. Output of each residual block is used to build the pyramid features. Residual feature maps on different level of pyramid are then feed to a sub-net made of 4 convolutional layers and finally a convolutional layer to predict the class scores and bounding box coordinated respectively. Similar to the original paper [14], we freeze the batch normalisation layers of ResNet based backbone networks.

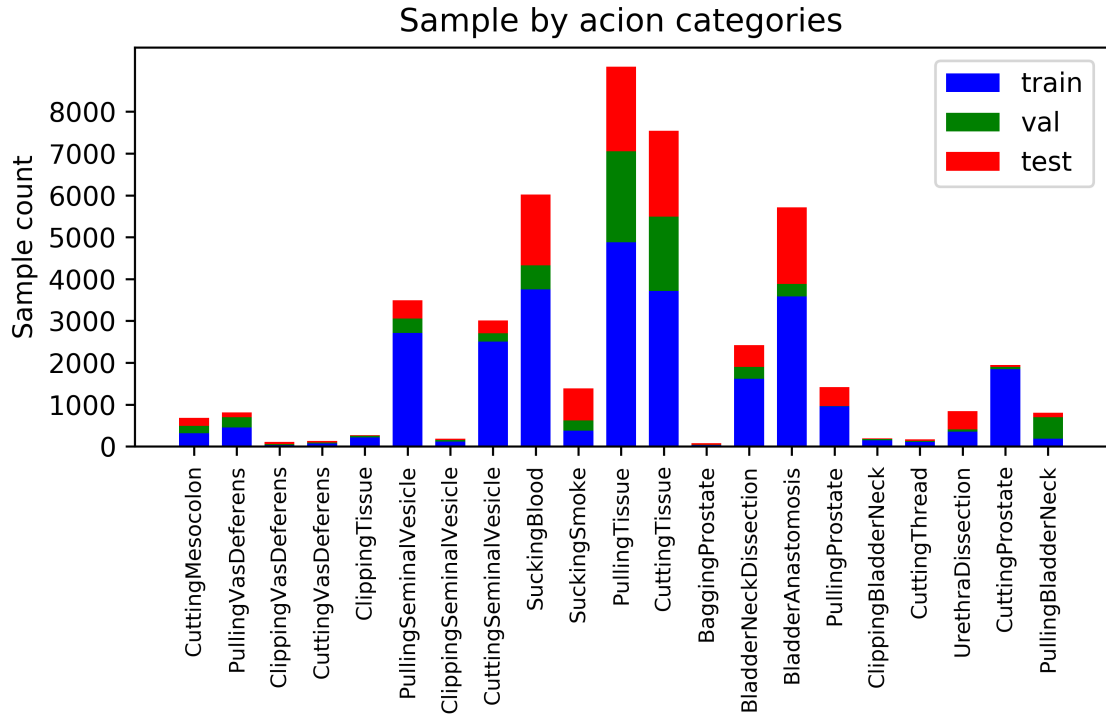


Figure 3: Distribution of samples for all action categories in training, validation and test splits. Blue, green and red colors of the bars represent training, validation and test set of ESAD.

Also, few initial layers are also frozen to avoid the overfitting. Finally, non maximum suppression (NMS) is used to discard the false positives in the model predictions at the test time.

We try two different loss functions to train classification sub-net of our baseline model. Since our baseline model is single stage model and based on [14], following [14], we train FPN with online hard example mining (OHEM)-loss [16] and Focal loss [15]. We use smooth-l1 loss [21] to train regression sub-net.

Implementation details: We train baseline model with various input image size, e.g. 200, 400 or 600. Short size of the input image is resized to input image size and longer side is resized with same scale. We set the learning rate 0.01 and batch size to 16. The networks are trained for 7K iteration with learning rate drop by the factor of 10 after 5K iterations. The

Loss	image size	AP_{10}	AP_{30}	AP_{50}	AP_{mean}	$Test - AP_{mean}$
Focal	200	33.8	17.7	6.6	19.4	15.7
Focal	400	35.9	19.4	8.0	21.1	16.1
Focal	600	29.2	17.6	8.7	18.5	14.0
Focal	800	31.9	20.1	8.7	20.2	12.4
OHEM	200	35.1	18.7	6.3	20.0	11.3
OHEM	400	33.9	19.2	7.4	20.2	13.6
OHEM	600	37.6	23.4	11.2	24.1	12.5
OHEM	800	36.8	24.3	12.2	24.4	12.3

Table 2: Results of the baseline models with different loss function and input image sizes, where backbone network fixed to ResNet50. AP_{10} , AP_{30} , AP_{50} , and AP_{mean} are presented on validation-set, while $Test - AP_{mean}$ is computed based on test-set similar to AP_{mean} .

complete model is implemented in pytorch and is provided as open access at <https://github.com/Viveksbawa/SARAS-ESAD-Baseline>. At the moment, the source code supports pytorch1.5 and Ubuntu with Anaconda distribution of python. It is tested on machines with 2/4/8 GPUs.

5.2. Evaluation metric

We also used three different IOU thresholds to compute the average precision (AP). We AP computed at 0.1, 0.3 and 0.5 AP are named as AP_{10} , AP_{30} and AP_{50} , respectively. Then, mean is computed at three thresholds to get a final evaluation score. the purpose of computing three different APs is to capture quality of detection as well as classification. As we know this is a new and complex task, it is very difficult to get good detection accuracy at higher threshold (can be seen in AP_{50} column). Hence, we want to identify both- how accurately model can detect the classes of actions present in the scene as well as the location of their bounding boxes.

5.3. Results

The results achieved by model with both of the losses are shown in table 2. We trained model on four different image sizes: 200, 400, 600, 800. Motive behind it is to observe the effect of tool sizes in the image on the detection accuracy. As we can observe in the table, with increase in image size, models with OHEM loss functions is able to achieve better detection accuracy on validation set. While the same can not said for test-set.

Loss	backbone	AP_{10}	AP_{30}	AP_{50}	AP_{mean}	$Test - AP_{mean}$
Focal	ResNet18	35.1	18.9	8.1	20.7	15.3
OHEM	ResNet18	36.0	20.7	7.7	21.5	13.8
Focal	ResNet34	34.6	18.9	6.4	19.9	14.3
OHEM	ResNet34	36.7	20.4	7.1	21.4	13.8
Focal	ResNet50	35.9	19.4	8.0	21.1	16.1
OHEM	ResNet50	33.9	19.2	7.4	20.2	13.6
Focal	ResNet101	32.5	17.2	6.1	18.6	14.0
OHEM	ResNet101	36.6	20.1	7.4	21.3	12.3

Table 3: Results of the baseline models with different loss function, backbone networks, where input image size is fixed to 400. AP_{10} , AP_{30} , AP_{50} , and AP_{mean} are presented on validation-set, while $Test - AP_{mean}$ is computed based on test-set similar to AP_{mean} .

The table 3 shows the results achieved by base model with different backbone networks while keeping the input image size fixed to 400 on both validation and test sets. It is clear from Tables 3 and 2 that OHEM loss performs better at validation set and focal loss performs better at test set.

From above results, it is clear that the presented baseline method is still far from achieve satisfactory performance. We hope that this will server as good benchmark for future methods which are specifically designed for endoscopic video.

6. Conclusion

This paper presents first of its kind dataset for action detection in surgical images. The dataset is developed on real videos collected from the prostatectomy procedures. This dataset aim to provide a benchmark for medical computer vision community to develop and test the state of the art algorithms for surgical robotics. We also released a baseline model along with the dataset which is developed using fully convolutional network architecture. Model is tested with two different type of loss functions: online hard example mining and focal loss. Focal loss based model is able to generalize much better for the test set. Additionally, we found out that bigger images size results in better model performance, generally. Complexity of dataset is also evaluated with different backbone architectures. Medium complexity/depth models like ResNet-34 perform better that higher depth models.

7. Acknowledgement

This work is conducted under the European Unions Horizon 2020 research and innovation programme under grant agreement No. 779813 and name of the project is Smart Autonomous Robotic Assistant Surgeon (SARAS project).

References

- [1] David P Azari, Yu Hen Hu, Brady L Miller, Brian V Le, and Robert G Radwin. Using surgeon hand motions to predict surgical maneuvers. *Human factors*, 61(8):1326–1339, 2019.
- [2] Pedro Felzenszwalb, David McAllester, and Deva Ramanan. A discriminatively trained, multiscale, deformable part model. In *2008 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2008.
- [3] Georgia Gkioxari, Ross Girshick, and Jitendra Malik. Contextual action recognition with r* cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1080–1088, 2015.
- [4] Chunhui Gu, Chen Sun, David A Ross, Carl Vondrick, Caroline Pantofaru, Yeqing Li, Sudheendra Vijayanarasimhan, George Toderici, Susanna Ricco, Rahul Sukthankar, et al. Ava: A video dataset of spatio-temporally localized atomic visual actions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6047–6056, 2018.
- [5] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [6] Rui Hou, Chen Chen, and Mubarak Shah. An end-to-end 3d convolutional neural network for action detection and segmentation in videos. *arXiv preprint arXiv:1712.01111*, 2017.
- [7] Rui Hou, Chen Chen, and Mubarak Shah. Tube convolutional neural network (t-cnn) for action detection in videos. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5822–5831, 2017.

- [8] Mihir Jain, Jan Van Gemert, Hervé Jégou, Patrick Bouthemy, and Cees GM Snoek. Action localization with tubelets from motion. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 740–747, 2014.
- [9] Vicky Kalogeiton, Philippe Weinzaepfel, Vittorio Ferrari, and Cordelia Schmid. Action tubelet detector for spatio-temporal action localization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4405–4413, 2017.
- [10] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017.
- [11] Bojan Kocev, Felix Ritter, and Lars Linsen. Projector-based surgeon–computer interaction on deformable surfaces. *International journal of computer assisted radiology and surgery*, 9(2):301–312, 2014.
- [12] Dong Li, Zhaofan Qiu, Qi Dai, Ting Yao, and Tao Mei. Recurrent tubelet proposal and recognition networks for action detection. In *Proceedings of the European conference on computer vision (ECCV)*, pages 303–318, 2018.
- [13] Ye LI, Jun OHYA, Toshio CHIBA, Rong XU, and Hiromasa YAMASHITA. Subaction based early recognition of surgeons’ hand actions from continuous surgery videos. *IEEEJ transactions on image electronics and visual computing*, 4(2):124–135, 2016.
- [14] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017.
- [15] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.
- [16] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multi-

- box detector. In *European conference on computer vision*, pages 21–37. Springer, 2016.
- [17] Martin Makary and Michael Daniel. Study suggests medical errors now third leading cause of death in the u.s. https://www.hopkinsmedicine.org/news/media/releases/study_suggests_medical_errors_now_third_leading_cause_of_death_in_the_us, 2016. Online; accessed 15-April-2020.
- [18] Dmitri Nepogodiev, Janet Martin, Bruce Biccard, Alex Makupe, Aneel Bhangu, Adesoji Ademuyiwa, et al. Global burden of postoperative death. *Lancet*, 393(401):33139–8, 2019.
- [19] Xiaojiang Peng and Cordelia Schmid. Multi-region two-stream r-cnn for action detection. In *European conference on computer vision*, pages 744–759. Springer, 2016.
- [20] Eduard Petlenkov, Sven Nomm, Juri Vain, and Fujio Miyawaki. Application of self organizing kohonen map to detection of surgeon motions during endoscopic surgery. In *2008 IEEE International Joint Conference on Neural Networks (IEEE World Congress on Computational Intelligence)*, pages 2806–2811. IEEE, 2008.
- [21] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015.
- [22] Suman Saha, Gurkirt Singh, and Fabio Cuzzolin. Amtnet: Action-micro-tube regression by end-to-end trainable deep architecture. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4414–4423, 2017.
- [23] Gurkirt Singh, Suman Saha, Michael Sapienza, Philip HS Torr, and Fabio Cuzzolin. Online real-time multiple spatiotemporal action localisation and prediction. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3637–3646, 2017.
- [24] Yicong Tian, Rahul Sukthankar, and Mubarak Shah. Spatiotemporal deformable part models for action detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2642–2649, 2013.

- [25] Beatrice van Amsterdam, Hirenkumar Nakawala, Elena De Momi, and Danail Stoyanov. Weakly supervised recognition of surgical gestures. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 9565–9571. IEEE, 2019.
- [26] Sandrine Voros and Gregory D Hager. Towards real-time tool-tissue interaction detection in robotically assisted laparoscopy. In *2008 2nd IEEE RAS & EMBS International Conference on Biomedical Robotics and Biomechatronics*, pages 562–567. IEEE, 2008.