

CholecTriplet2021: A benchmark challenge for surgical action triplet recognition

 Chinedu Innocent Nwoye ^{✉,a,*}, Deepak Alapatt^a, Tong Yu^a, Armine Vardazaryan^b, Fangfang Xia^c, Zixuan Zhao^c, Tong Xia^d, Fucang Jia^d,  Yuxuan Yang^e,  Hao Wang^e, Derong Yu^f,  Guoyan Zheng^f, Xiaotian Duan^g, Neil Getty^g, Ricardo Sanchez-Matilla^h, Maria Robu^h, Li Zhangⁱ, Huabin Chenⁱ, Jiacheng Wang^j, Liansheng Wang^j, Bokai Zhang^k,  Beerend Gerats^l, Sista Raviteja^m,  Rachana Sathish^m,  Rong Tao^f, Satoshi Kondoⁿ, Winnie Pang^o, Hongliang Ren^u,  Julian Ronald Abbing^l, Mohammad Hasan Sarhan^k, Sebastian Bodenstedt^p, Nithya Bhasker^p,  Bruno Oliveira^{q,r,s},  Helena R. Torres^{q,r,s}, Li Ling^e, Finn Gaida^t, Tobias Czempiel^t,  João L. Vilaça^q,  Pedro Morais^q,  Jaime Fonseca^s, Ruby Mae Egging^l, Inge Nicole Wijma^l, Chen Qianⁱ, Guibin Bianⁱ, Zhen Liⁱ, Velmurugan Balasubramanian^m,  Debdoot Sheet^m, Imanol Luengo^h,  Yuanbo Zhu^e,  Shuai Ding^e, Jakob-Anton Aschenbrenner^k, Nicolas Elini van der Kar^l, Mengya Xu^o, Mobarakol Islam^o, Lalithkumar Seenivasan^o, Alexander Jenke^p, Danail Stoyanov^{h,v}, Didier Mutter^{b,x},  Pietro Mascagni^{a,w},  Barbara Seeliger^{a,b,x}, Cristians Gonzalez^{b,x},  Nicolas Padoy^{a,b}

^aICube, University of Strasbourg, CNRS, France

^bIHU Strasbourg, France

^cDepartment of Computer Science, University of Chicago, US

^dLab for Medical Imaging and Digital Surgery, Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences, China

^eSchool of Management, Hefei University of Technology, Hefei, China

^fInstitute of Medical Robotics, Shanghai Jiao Tong University, Shanghai, China

^gArgonne National Laboratory, 9700 S Cass Ave, Lemont, IL 60439, US

^hDigital Surgery, a Medtronic Company, London, UK

ⁱInstitute of Automation, Chinese Academy of Sciences, China

^jDepartment of Computer Science at School of Informatics, Xiamen University, Xiamen, China

^kJohnson & Johnson

^lMeander Medical Centre, The Netherlands

^mIndian Institute of Technology Kharagpur, India

ⁿMuroran Institute of Technology, Japan

^oDepartment of Biomedical Engineering, National University of Singapore, Singapore

^pDepartment for Translational Surgical Oncology, National Center for Tumor Diseases Partner Site Dresden, Germany

^q2Ai School of Technology, IPCA, Barcelos, Portugal.

^rLife and Health Science Research Institute (ICVS), School of Medicine, University of Minho, Braga, Portugal.

^sAlgoritimi Center, School of Engineering, University of Minho, Guimeraes, Portugal.

^tTechnical University of Munich, Germany

^uDepartment of Electronic Engineering, The Chinese University of Hong Kong, Hong Kong

^vWellcome/EPSCRC Center for Interventional and Surgical Science, University College London, UK

^wFondazione Policlinico Universitario Agostino Gemelli IRCCS, Rome, Italy

^xUniversity Hospital of Strasbourg, France

ABSTRACT

Context-aware decision support in the operating room can foster surgical safety and efficiency by leveraging real-time feedback from surgical workflow analysis. Most existing works recognize surgical activities at a coarse-grained level, such as phases, steps or events, leaving out fine-grained interaction details about the surgical activity; yet those are needed for more helpful AI assistance in the operating room. Recognizing surgical actions as triplets of $\langle instrument, verb, target \rangle$ combination delivers comprehensive details about the activities taking place in surgical videos. This paper presents *CholecTriplet2021*: an endoscopic vision challenge organized at MICCAI 2021 for the recognition of surgical action triplets in laparoscopic videos. The challenge granted private access to the large-scale *CholecT50* dataset, which is annotated with action triplet information. In this paper, we present the challenge setup and assessment of the state-of-the-art deep learning methods proposed by the participants during the challenge. A total of 4 baseline methods from the challenge organizers and 19 new deep learning algorithms by competing teams are presented to recognize surgical action triplets directly from surgical videos, achieving mean average precision (mAP) ranging from 4.2% to 38.1%. This study also analyzes the significance of the results obtained by the presented approaches, performs a thorough methodological comparison between them, in-depth result analysis, and proposes a novel ensemble method for enhanced recognition. Our analysis shows that surgical workflow analysis is not yet solved, and also highlights interesting directions for future research on fine-grained surgical activity recognition which is of utmost importance for the development of AI in surgery.

KEYWORDS: Surgical activity recognition, tool-tissue interaction, CholecT50, laparoscopic video, action recognition.

1. Introduction



Fig. 1. A cross-section of CholecT50 dataset for the CholecTriplet2021 challenge. Represented surgical action triplets are illustrated for different time-points during a laparoscopic cholecystectomy.

Activity recognition is the basis for the development of many potential applications in health, surveillance, manufacturing, sports, etc (Avci et al., 2010). In surgery, it can be leveraged to provide intra-operative context-awareness and decision support (Maier-Hein et al., 2017a). This could augment surgeons’ capabilities, fostering safety and efficiency in the operating room (OR) (Vercauteren et al., 2019).

Despite the vast literature in medical computer vision, a majority of the works tackle activity recognition at a very coarse-grained level such as phase (Twinanda et al., 2016), which does not provide an accurate picture of activities taking place. As an example, the *calot triangle dissection* phase in cholecystectomy contains a multitude of finer actions that would be relevant to recognize. The finer division of activity recognition such as step (Lecuyer et al., 2020), gesture (DiPietro et al., 2019) or action (Wagner et al., 2021) recognition leaves out details about the anatomy and thus, does not provide comprehensive details to describe the tool-tissue interaction. The detection of the used surgical instruments and their target anatomy as well as their fine-grained interaction details is necessary for the development of artificial intelligence (AI) that is safe for the patient (Nwoye, 2021). This opens a vista of opportunities to develop methods that recognize the elements involved in tool-tissue interaction while accounting for the relationship and multiple instances.

In general computer vision, human activity is modeled as triplets $\langle \text{human}, \text{verb}, \text{object} \rangle$ providing full-scale and expressive details on Human-object interaction (HOI) (Chao et al., 2015b). Translating this formalism to surgical vision, however, was not achieved until very recently, when Nwoye et al. (2020) presented the video recognition of surgical activities as triplets of the used instruments, actions performed, and the underlying target anatomy.

We present a critical study on surgical action triplet recognition in the form of a challenge termed **CholecTriplet2021**¹, to pave the way for research targeting fine-grained and detailed recognition of surgical activities from videos. This international challenge was organized as part of the Endoscopic Vision (EndoVis) Grand Challenge (Speidel et al., 2021) and hosted at MICCAI 2021. The challenge presented promising technologies for the detailed understanding of tool-tissue interactions in minimally invasive surgery. A total of 19 teams participated, the highest record since the inception of the EndoVis Grand Challenge series.

The challenge provided a platform for the scientific community to perform comparative benchmarking and validation of endoscopic vision algorithms in a promising direction in surgical activity recognition. We provided private access to a high-

quality large-scale surgical action triplet dataset, *CholecT50* (Nwoye et al., 2022), for both method development and validation. In addition to these contributions brought by the event itself, the challenge report presented here offers contributions of its own. After individual descriptions for each approach, we highlight several trends in an in-depth methodological comparison, providing a comprehensive overview of possible strategies for tackling the surgical action triplet problem. To further facilitate future research efforts, we also compile the implementation details of all featured submissions. We then set a new upper baseline (+4.3% AP to challenge methods) for the surgical action triplet recognition problem, by proposing a simple but effective algorithm ensembling models featured in the challenge. Finally, results are analyzed in depth: our quantitative analysis considers multiple metrics to cover all aspects of the triplet recognition problem. A rich selection of qualitative results is presented as well, to better understand the behavior of all the methods.

The paper is organized as follows: we position our work with respect to the related literature in the next section. Afterward, we present the challenge overview and setup including the used dataset and a summary of participation. This is followed by the methods presented at the challenge including the evaluation protocols, results, and extensive analysis of their performance. We conclude by highlighting the benefits of this study, insights, and prospects.

2. Related work

The CholecTriplet2021 challenge relates to several research topics, for which we present the relevant literature in the following paragraphs.

2.1. Activity recognition

One of the key concerns of computer vision is the recognition of human activities; a task for which a wide variety of approaches has been proposed, both in medical and non-medical domains. In those approaches, the definition of this task can be more or less granular, which is a key factor in determining its difficulty as well as its utility. Early activity recognition work in general computer vision involved broad, coarse-grained classification tasks: the 2012 version of the PASCAL VOC (Everingham et al., 2015) challenge proposed an action classification task on static images with 10 classes. HMDB-51 (Kuehne et al., 2011) and UCF-101 (Soomro et al., 2012) are collections of realistic action video clips extracted from YouTube. As datasets expanded, the number of classes increased; classes became more diverse but also finer-grained. For example, Kinetics (Carreira and Zisserman, 2017) features three classes related to cycling ("*riding a bike*", "*riding a mountain bike*", "*falling off a bike*") while HMDB-51 only contains one (Kuehne et al., 2011). Generally speaking, the advantage of having more granular classes is that they enable more detailed and informative descriptions.

In surgical computer vision, fine-grained activity descriptions are particularly valuable: as key components in concepts for context-aware surgical systems (Maier-Hein et al., 2017a,b), activity recognition algorithms tie into clinical outcomes. However, descriptions achieved by surgical vision algorithms have historically been coarse. The task of surgical phase recognition (Twinanda et al., 2016; Dergachyova et al., 2016; Funke et al.,

¹<https://cholectriple2021.grand-challenge.org/>

2018; Yu *et al.*, 2019; Zisimopoulos *et al.*, 2018; Hajj *et al.*, 2018; Czempiel *et al.*, 2020), one of the main research topics in this area, breaks down an entire surgery into a small number of chunks, each with a broadly defined role within the surgical procedure. Since phases can last several minutes, they may contain many individual actions; this lack of detail is ultimately what limits the utility of the phase information. For this reason, other studies addressed more granular classes: subdivisions of phases into steps can be found in Ramesh *et al.* (2021); Lecuyer *et al.* (2020) and recognition of action verbs has been explored in Rupprecht *et al.* (2016); Khatibi and Dezyani (2020). In the previous edition of EndoVis, an action recognition task was featured (Wagner *et al.*, 2021), modeling tool-tissue interaction to a certain degree: each class was defined by one verb.

Overall, a shift towards finer-grained activity recognition has become a major focus in recent research on activity recognition, especially for surgery.

2.2. Action triplet recognition

The finest level of granularity in visual activity understanding is currently achieved by decomposing actions into their individual components: who performs the action, what the action is, and what the action is performed on. This decomposition into action triplets of (*subject, verb, object*) is central in HOI (Chao *et al.*, 2015b) studies. Mallya and Lazebnik (2016) used CNN features extracted from humans and objects detected in frames. Chao *et al.* (2018) proposed a multi-stream architecture to model spatial relationships. Gkioxari *et al.* (2018)'s method was built around a FasterRCNN object detector to locate humans. Qi *et al.* (2018) introduced the Graph Parsing Neural Network (GPNN), representing human-object interactions with the adjacency matrix and node labels of a graph.

In surgical computer vision, action triplets (Katic *et al.*, 2014) in the form of (*surgical tool/instrument, action verb, target anatomy*) are used to describe tool-tissue interactions (TTI) (Nwoye, 2021). Early works used them as auxiliary annotations to improve surgical phase recognition (Katic *et al.*, 2014, 2015). The first method to actually perform action triplet recognition from videos was introduced by Nwoye *et al.* (2020) as the Triplet, featuring verb and target detection using instrument class activation guidance, as well as triplet association using a 3D interaction space. A new model (Nwoye *et al.*, 2022) with more advanced modeling of interactions between triplet components using an attention mechanism was later developed. Xu *et al.* (2021) proposed a cross-domain method for automatic surgical captions that capture semantic relationships, similarly to action triplets.

2.3. Action triplet datasets

The study of action triplet recognition requires datasets annotated with triplet components. In general computer vision, an early example is the HICO dataset (Chao *et al.*, 2015a). CAD-120 (Koppula *et al.*, 2013) is annotated with object affordances. V-COCO (Sadhu *et al.*, 2021), as an extension of the widely used MS-COCO (Lin *et al.*, 2014), added visual semantic role labels; the provided bounding box annotations also enabled HOI spatial detection. HICO later received an update in the form of HICO-DET (Chao *et al.*, 2018), similarly incorporating bounding boxes. HCVRD (Zhuang *et al.*, 2018) was proposed as a benchmark for detecting human-centered visual relationships, which includes action verbs among other types.

Ambiguous-HOI (Li *et al.*, 2020) collects difficult examples from several of the datasets previously mentioned, in order to form a challenging HOI benchmark.

In the surgical domain, datasets offering action triplet annotations describing tool-tissue interaction are much more recent. In a challenge organized by Wagner *et al.* (2021), a cholecystectomy video dataset with annotations for four surgical action verbs was provided. The SARAS-ESAD dataset (Bawa *et al.*, 2021) featured more refined classes with actions described by both the verb and the anatomy; bounding boxes for 21 action classes were provided as well. The first dataset with full annotations for each triplet component was introduced in Nwoye *et al.* (2020) as CholecT40. The expanded version renamed CholecT50 (Nwoye *et al.*, 2022) is employed in this challenge. Another dataset used in Xu *et al.* (2021) incorporated surgical captions; despite not explicitly following the action triplet formalism, the level of detail offered is similar.

2.4. Endoscopic vision challenges

As a relatively new problem, surgical action triplet recognition has only received little attention despite its potential. This is a major motivating factor for *CholecTriplet2021*: over the past few years, challenges have played an important role in the surgical computer vision community due to the exposure they can bring to interesting research topics, as well as their ability to introduce and compare a wide variety of original methods. For instance, the M2CAI 2016 challenge (Stauder *et al.*, 2016) featured two tasks - surgical phase recognition and tool presence detection for cholecystectomy. The first edition of the CATARACTS challenge (Hajj *et al.*, 2019), involved tool recognition on cataract surgery videos; later editions became part of the EndoVis (Endoscopic Vision) challenge series. Each iteration of EndoVis featured multiple sub-challenges, some of which focused on surgical activity understanding: segmentation and tracking of tools for colorectal surgery videos (2015), phase recognition for colorectal surgery using video and sensor data (2017), and finally, phase, tool and action recognition for cholecystectomy videos (Wagner *et al.*, 2021).

Outside of EndoVis subchallenges, the SARAS-ESAD (Bawa *et al.*, 2021) organized within the MIDL 2020 challenge presents a benchmark for surgical action detection using a large-scale video dataset (ESAD) offering another level of granularity (*verb, target*) as a single label) as well as bounding boxes for action localization.

3. CholecTriplet challenge

The goal of this challenge is to assess AI solutions for fine-grained surgical activity recognition. This recognition is modeled as a triplet (Nwoye *et al.*, 2020), with the following notation: (*instrument, verb, target*)

3.1. Task

The task is to develop a machine learning method to recognize these triplets directly from unseen surgical videos. In modeling surgical action triplet, a prevailing problem to tackle is the simultaneous detection of the correct instruments, verbs, and targets in every image frame and resolving their association. This is challenging since the involvement of a component in a triplet can be visually subtle. Hence, the challenge also assesses the detection of these individual components for a more

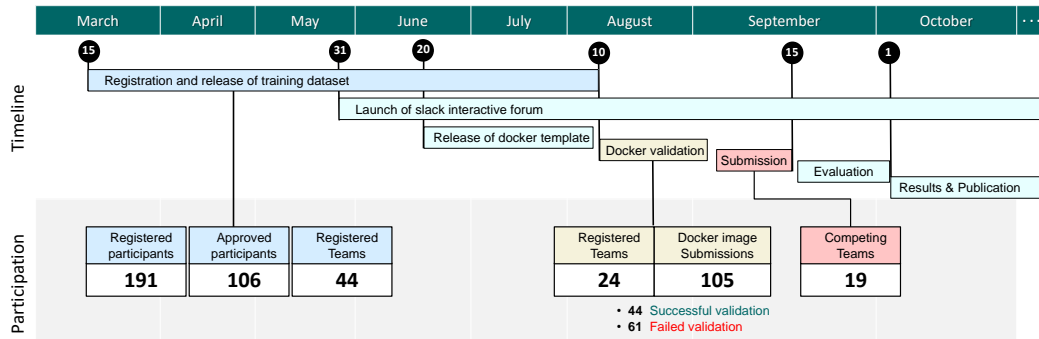


Fig. 2. CholecTriplet2021 challenge timeline and activities.

insightful analysis of a model’s understanding of the triplet’s composition.

3.2. Dataset

The dataset used for the challenge is *CholecT50* (Nwoye *et al.*, 2022), which consists of 50 video recordings of laparoscopic cholecystectomy that have been annotated with the binary presence of action triplets. At 1 frame per second (fps), the dataset reaches a total of 100.9K frames and 161K triplet instance labels. *CholecT50* consists of 100 action triplet classes composed from 6 instruments (*grasper*, *bipolar*, *hook*, *scissors*, *clipper*, *irrigator*), 10 verbs (*grasp*, *retract*, *dissect*, *coagulate*, *clip*, *cut*, *aspirate*, *irrigate*, *pack*, *null*), and 15 targets (*gallbladder*, *cystic-duct*, *cystic-artery*, *blood-vessel*, *fluid*, *abdominal-wall or cavity*, *liver*, *omentum*, *peritoneum*, *gut*, *specimen-bag*, *null*). The triplet labels we provided in form of a single binary presence for each triplet class (*instrument*, *verb*, *target*) taken as a whole. Binary labels for each component of the triplet were also provided.

On the data split, 45 out of 50 videos in the dataset were released to the participants for training and validation; those videos were part of the publicly released *Cholec80* dataset Twinanda *et al.* (2016). The remaining 5 videos, which were not public, were withheld from the participants as the testing set.

3.3. Challenge design

The challenge was conducted as part of EndoVis grand-challenge (Speidel *et al.*, 2021) in MICCAI 2021. The challenge proposal went through two rounds of open review between November 2020 and February 2021. By early March 2021, a call for participation was circulated and the challenge officially started on March 15, 2021, with a public release of the training data. Participation took place by web registration and signing a non-disclosure contract on the use of the challenge dataset.

By the end of May, a Slack² channel was created for the registered participants and their teams. A blog was provided with snippets of code and instructions for getting started with the challenge. Reference to the baseline methods and code for some specialized functions were also provided. A Docker³ submission template, development guide, and evaluation metrics were provided to the participants during their method development

phase. Participating teams were allowed to develop novel methods, fine-tune a state-of-the-art method, or improve on existing solutions. Entrants were allowed to pre-train their model on any third-party publicly available dataset. Developed methods were intended to take sequential image frames from videos, process, and return a vector probability for the 100 triplet classes for each image. Due to the possibility of multiple instances of the triplets, both the triplets and their three components were cast as a multi-label classification problem. The output probability vectors for the three components of the triplets were derived using a filtering algorithm proposed in Nwoye *et al.* (2022).

By the terminal point of the development phase, a validation process was initiated for users to ascertain that their Docker container had been built with the correct input/output format and was able to run without run-time errors on a set of randomly selected images from the training set. Teams were allowed to validate multiple times until an error-free Docker was obtained, but within a time-bound of three weeks which lasted till Sept 5, 2021. The final submission was performed once using only a validated Docker image container. These submissions were run and evaluated on the hidden test set by the challenge organizers with the outcome presented at the MICCAI 2021 satellite event. The challenge timeline is presented in Fig. 2.

We also provide a one-month post-challenge window (November 15 - December 15, 2021) during which all teams could re-evaluate their updated model if their prior submissions had inaccuracies.

3.4. Method submission

All Docker submissions were collected as saved Docker images uploaded to Dropbox⁴. Only Dockers strictly adhering to a predefined Docker format and input/output requirements were considered as valid submissions. In practice, this involved passing several automatic checks including ensuring compatibility with a defined directory structure and naming format. In terms of output format, automatic checks were conducted to ensure that an output probability value between 0 and 1 was written for each input frame and consequently each available video. Additionally, each submitted method was required to pass a causality check on a predefined, hidden subset of frames to ensure that future frames were not utilized. During the validation phase, participants were expected to make corrections to the submissions based on automatically generated feedback that was shared by email to ensure that only validated Docker images were submitted for final evaluation.

²<https://www.slack.com>

³<https://hub.docker.com>

⁴<https://www.dropbox.com>

Table 1. Demography of final participating teams.[†]

Asia				Europe					North America
China	Singapore	Japan	India	Germany	UK	Netherlands	Portugal	France	USA
8	2	1	1	3	2	1	1	1	3

[†] Some teams have affiliations in multiple countries.

3.5. Awards

The winning submission earned an NVIDIA GPU. Monetary prizes were also awarded to the top 3 competing teams.

3.6. Participation statistics

As presented in Fig. 2, an initial 44 teams registered by the end of the enrollment window, accounting for 106 approved individual participants out of 191 recorded registrations. During the Docker validation phase, 24 teams submitted a total of 105 Docker containers which were evaluated with a 41.9% success rate. On the deadline of September 15, 2021, a total of 19 teams' submissions were received for final evaluation. A baseline submission from the organizers brings this to a total of 20 teams. The organizer's submission consists of four different baseline methods. The participating teams were drawn from 10 countries across 3 continents. The demography of the final participating teams is presented in Table 1.

4. Methods

4.1. Baseline methods

Four baseline methods are provided by the challenge organizers (**Team CAMMA**): (1) MTL baseline (Nwoye *et al.*, 2020), (2) Tripnet (Nwoye *et al.*, 2020), (3) Attention Tripnet (Nwoye *et al.*, 2022), and (4) Rendezvous (Nwoye *et al.*, 2022). We present a brief overview of the baseline models below:

4.1.1. MTL baseline

The Multi-Task Learning baseline (Nwoye *et al.*, 2020) is built around a ResNet-18 backbone, which serves as the common visual feature extractor for three separate branches. Each branch is a 3-layer (2 convolutional, 1 fully connected) neural network responsible for recognizing one of the triplet components. The instrument branch differs from the other two with the addition of Global Max Pooling (GMP). A final fully-connected layer is used to combine the three components' prediction into a final triplet prediction.

4.1.2. Tripnet

The Tripnet (Nwoye *et al.*, 2020) also relies on multi-task learning from a ResNet-18 backbone but provides a stronger baseline thanks to its two characteristics. The first one is the Class Activation Guide (CAG): here the instrument branch, called the Weakly Supervised Localization (WSL) module, generates per-instrument Class Activation Maps, which are forwarded to a subnetwork in charge of verb and target detections as shown in Fig. 3. These maps are concatenated to the verb and target features. This additional information on instrument positions, in the form of concatenated instrument activation maps, provides clues for better detection of verbs and targets, since those are located at the tooltips.

The second innovation of the Tripnet is the 3D Interaction Space (3DIS), which associates the triplet components. Log

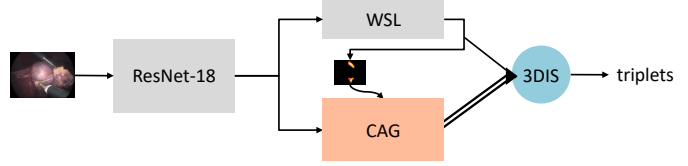


Fig. 3. Overview of Tripnet: a baseline method for action triplet recognition built on ResNet-18 backbone, weakly supervised localization (WSL) of surgical instruments, class activation guide (CAG) for verb-target recognition, and 3D interaction space (3DIS) for triplet association.

probabilities for each component are projected into a 3D grid of dimensions $n_I \times n_V \times n_T$ (number of instruments, verbs, and targets, respectively) by a trainable function. Each point in the 3D grid represents the probability of the triplet, as estimated by the 3DIS function.

4.1.3. Attention Tripnet

The Attention Tripnet (Nwoye *et al.*, 2022) enhances the previous Tripnet, by replacing the CAG with an attention model, the Class Activation Guided Attention Mechanism (CAGAM). While the CAG only proceeds by feature concatenation, the CAGAM computes additive enhancements for verb and target class maps based on the instrument maps. Those enhancements are obtained by two different attention mechanisms: a channel (instrument type) attention for verbs and a position attention for targets. Using attention in this manner explicitly steers the verb and target detections towards the correct locations, much more strongly than the concatenation in the Tripnet's CAG.

4.1.4. Rendezvous

Rendezvous (Nwoye *et al.*, 2022), or simply RDV, is the strongest baseline model. In contrast to the Attention Tripnet, triplet association is performed using a Transformer-like model built around a Multi-Head of Mixed Attention (MHMA), instead of the 3DIS. The MHMA operates on a set of four features: H_{IVT} , H_I , H_V , H_T for global triplet, instrument, verb, and target features respectively. H_{IVT} is drawn from a bottleneck layer, connected to an early layer of the ResNet-18 backbone. H_I is taken from the WSL's output, and the last two are taken from the CAGAM.

Inside the MHMA, H_{IVT} is processed by a self-attention head. A projection function maps it to three separate vectors - query, key, and value. A scaled dot product then transforms them into one refined feature. H_I is fed to a cross-attention head: the key and value are obtained by projecting H_I , but the query is the same one used in H_{IVT} 's self-attention head. H_V , H_T 's cross-attention heads operate in the same manner as H_I 's.

The resulting four refined features are concatenated; two convolutions and an AddNorm operation complete the MHMA, which outputs a refined version of H_{IVT} . Inside Rendezvous, a stack of 8 MHMAs is employed; the output of this stack is used to make the final prediction on the triplet.

Table 2. A complete cross-view of the MICCAI EndoVis CholecTriplet 2021 participating teams including their affiliations and presented methodologies.

Team	Affiliation(s)	Method
1 2Ai	Applied Artificial Intelligence Laboratory, Portugal	Surgical video analysis using an ensemble of multi-task recurrent convolutional networks (versions 1 & 2)
2 ANL-Triplet	Argonne National Laboratory, USA	ANL-Triplet: Exploiting temporal information for triplet recognition
3 Band of Broeders	Meander Medical Centre, The Netherlands Johnson & Johnson, USA	YoloV5 for surgical action triplet detection
4 CAMMA	University of Strasbourg, France (Organizers)	MTL Baseline, Tripnet, Attention Tripnet, and Rendezvous
5 CAMP	Technical University of Munich, Germany	EndoVisNet: Phase-guided temporal endoscopic action triplet classification
6 Casia Robotics	Institute of Automation, Chinese Academy of Sciences, China	Triplet translation embedding network with coordinate attention
7 Ceaiik	Indian Institute of Technology Kharagpur, India	Spatio-temporal learning of action triplets in surgical videos
8 CITI SJTU	Shanghai Jiao Tong University, China	Action triplet recognition via convolutional LSTMs and multi-task learning
9 Digital Surgery	Digital Surgery, a Medtronic Company, London, UK Wellcome/EPSRC Center for Interventional and Surgical Science, University College London, UK	TE-TAR: Temporal ensemble triplet action recognition
10 Lsgroup	Xiamen University, China	Feature fusion and weak locational information calculating in triplet classification multi-task
11 HFUT-MedIA	Hefei University of Technology, China	COEMNet: Correlation embedded multi-task network
12 HFUT-NUS	National University of Singapore, Singapore Hefei University of Technology, China	Multi-task learning based surgical interaction triplet recognition
13 J&M	Johnson & Johnson, US Meander Medical Centre, The Netherlands	Surgical action triplet recognition with Efficient-MSTCN
14 Med Recognizer	Chinese University of Hong Kong, China	Surgical action triplet recognition via temporal memory relation gradient network
15 MMLAB	National University of Singapore, Singapore	Temporal triplet net for triplet presence detection in surgical videos
16 NCT-TSO	National Center for Tumor Diseases Partner Site Dresden, Germany	Multi-task learning framework for action triplet recognition
17 SIAT CAMI	Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences, China	Multi-task mutual channel recurrent net for fine-grained surgical triplet recognition
18 SJTU-IMR	Shanghai Jiao Tong University, China	Tracking surgical actions with transformers and action label-guided fine-grained information aggregation
19 SK	Muroran Institute of Technology, Japan	Action triplet recognition with weakly-supervised attention of surgical instruments
20 Trequantista	University of Chicago, USA	Phase-aware multitasking action recognition model with adjustment for low-data triplet classes

4.2. Competing methods

4.2.1. Team 2Ai: *Surgical video analysis using an ensemble of multi-task recurrent convolutional networks*

Team 2Ai proposed a solution (version 1) for this challenge consisting of an ensemble of multi-task recurrent convolutional networks. Each model architecture consists of a multi-task recurrent convolutional network with four heads, where each branch targets one of four tasks, namely surgical instrument detection, verb recognition, phase identification, and target recognition. To extract generic visual features, they used a shared feature extractor for all four branches. Specifically, for each of the instrument detection and target recognition branches, a fully connected layer is connected to the backbone to compute the signals for both tasks. Here, a sigmoid activation layer is applied to produce the final predictions. In the branches for surgical action and phase recognition, long short-term memory (LSTM) units (Hochreiter and Schmidhuber, 1997) are connected to the backbone to leverage the temporal context of the current frame. A sigmoid and softmax layers are added to the end of the last LSTM units in the action and phase recognition branches, respectively. Finally, binary cross-entropy is used as a

loss function for the surgical instrument detection, verb recognition, and target recognition tasks, and cross-entropy is used for phase identification. The different networks, each with a different feature extraction backbone, are individually optimized and trained. Afterward, majority-vote ensemble is applied to combine the predictions resulting from the different networks. As a final step, a temporal smoothing technique is applied to each task to avoid temporally incoherent results. To meet the challenge output requirements, the individual signals are finally represented as surgical actions triplets. 2Ai version 2 is submitted post-challenge to correct the error in the output format of the initial model by changing the final output mapping from binary to probability scores. In this revised version, the ensemble is performed using an average of probabilities outputs from all the networks.

4.2.2. Team ANL-Triplet: *Exploiting temporal information for triplet recognition*

Team ANL-Triplet used three ResNet-18 backbones (He et al., 2016) for instrument, verb and target prediction. In a multi-task setup, these component predictions are concatenated before us-

ing a single convolutional layer and a fully connected layer to recognize the action triplets represented in the frame. To improve verb and target recognition performance, these models are initialized using the learned instrument model weight. Additionally, an independent ResNet-34 is trained to incorporate temporal information by predicting action triplets directly from features extracted from the current frame and 7 previous frames. In both triplet recognition approaches, a smoothing factor of $0.2, 0.2^2$ is used for predictions that predicted two and one components correctly, respectively, reducing the penalty for semantically closer predictions. The final prediction is made by taking the average of the two triplet predictions.

4.2.3. Team Band of Broeders: *YOLOv5 for surgical action triplet detection*

The Band of Broeders utilized manually generated bounding box annotations of the surgical triplets to train a YOLOv5 based object detector for action triplet prediction. They chose the YOLOv5 (?) for its state-of-the-art performance on various object detection datasets. Their network⁵ consists of three stages: (1) the backbone CSP-DenseNet (Wang et al., 2020a), which is used for its gradient efficiency, (2) the neck Path Aggregation Network (PA-Net) (Liu et al., 2018), which is used for multi-scale feature extraction, and (3) the head which produces the final bounding box predictions. These components allow for the building of neural networks with large receptive fields and a multi-scale view of the frame enabling the detection of objects of various sizes.

4.2.4. Team CAMP: *EndoVisNet: Phase-guided temporal endoscopic action triplet classification*

Team CAMP's methodology is based on the idea that the presence of an action in a given frame implies that only a smaller set of actions could occur in the immediate frames that follow. To leverage this temporal property, visual features from both the input frame and its t preceding frames are extracted using the SlowFast (Feichtenhofer et al., 2019) network to ensure that the extracted features are relevant and independent of the movement speed in an input video. After extraction, the temporal features are pooled into 1 dimension to make a final prediction. Feature classification is modeled using a similar approach to Tripnet (Nwoye et al., 2020), leveraging a multi-task learning approach, but with an extra branch for the phase detection, which is trained jointly on labels extracted from the Cholec80 dataset (Twinanda et al., 2016). The final action triplet classification scores are computed as a learned linear combination of the instrument, verb, and target prediction scores.

4.2.5. Team Casia Robotics: *Triplet translation embedding network with coordinate attention*

Team Casia Robotics' submission (see Fig. 4) builds on the work (Nwoye et al., 2020) and replaces the introduced 3D interaction space with a translation embedding module, following Zhang et al. (2017). Here, the translation embedding module couples the instrument, target, and verb features to produce the final triplet predictions. Specifically, they used a ResNet-18 (He et al., 2016) backbone for extracting spatial features followed by three sub-networks with convolutional layers

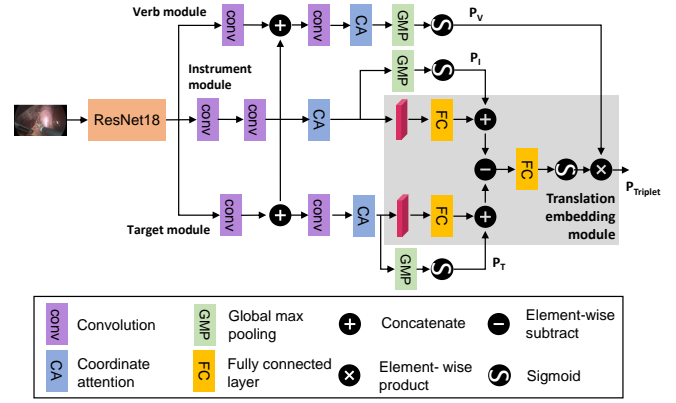


Fig. 4. Overview of the Casia Robotics submission.

and fully-connected layers that are used to predict instruments, verbs, and targets, respectively. To model the complex relationship between triplets, they modeled projection matrices, W_i and W_t , that translate learned representations from the feature space to a shared relation space. The interaction relation between the three learned component feature vectors f_t (target), f_i (instrument/tool), f_v (verb) is then represented as:

$$(W_t f_t - W_i f_i) \cdot f_v.$$

All invalid component combinations are then masked out to only generate predictions for the 100 relevant triplet classes. Further, a coordinate attention mechanism (Hou et al., 2021a) is used to incorporate positional information into each sub-network to more accurately locate each component of the triplet.

4.2.6. Team Ceaiik: *Spatio-temporal learning of action triplets in surgical videos*

Team Ceaiik's model architecture design is based on the triplet classification as a composition of three individual tasks namely instrument, verb, and target classifications (Nwoye et al., 2020). In the first stage of training, the model utilizes a ResNet-50 (He et al., 2016) to extract spatial features based on the results of Mishra et al. (2017); Mondal et al. (2019) and employs 3 classification heads to predict instrument, verb, and target class probabilities. A final classification head is employed to utilize the previously computed component probabilities, which are aggregated to make a triplet prediction. Further, in the second stage of training, an effort is made to incorporate certain temporal properties of action triplets in their model design. Keeping the learned ResNet-50 weights frozen from the first stage of training, the model extracts visual features which are passed through an LSTM module (Hochreiter and Schmidhuber, 1997) before instrument, verb, target, and triplet prediction are performed using a similar approach to the first stage.

4.2.7. Team CITI SJTU: *Action triplet recognition via convolutional LSTMs and multi-task learning*

Team CITI SJTU focused their method on modeling the sub-components and temporal coherence when predicting action triplets. Their multi-task deep learning network includes four branches with one main triplet branch and three auxiliary branches generating the recognition results for instruments, verbs, and targets, respectively. All the branches share the same

⁵non-competing method as it is trained on private annotations

ResNet-50 (He et al., 2016) network for feature extraction, followed by two convolutional Long Short Term Memory (ConvLSTM) layers (Shi et al., 2015) for modeling spatial-temporal relationships. After training, they used the triplet prediction branch to obtain the final triplet prediction for inference. The three auxiliary branches allow them to add fine-grained information for training to significantly boost triplet recognition performance.

4.2.8. Team Digital Surgery: TE-TAR: Temporal ensemble triplet action recognition

Team Digital Surgery’s proposed model, TE-TAR, consists of an ensemble of encoders, an LSTM model (Hochreiter and Schmidhuber, 1997), and a classification layer. The ensemble of encoders is composed of four HRNet32 backbones (Wang et al., 2020b) and a classification head that efficiently combines multi-scale information. Each ensemble is trained with a different subset of data for learning more diverse features. For each image, four features are generated in total (i.e., one per encoder), which are combined by summation. In addition, to allow the model to estimate the action of the instrument (verb) and the anatomy where the instrument is applied (target), they proposed the use of temporal information, using a sliding window approach with a window length of 10 frames. Features are encoded by the ensemble for all frames in the window and then fed to an LSTM. As they formulated the action triplet task as a classification problem, the final triplets are estimated by feeding the aggregated features to a linear classification layer.

4.2.9. Team Lsgroup: Feature fusion and weak locational information calculating in triplet classification multi-task

Team Lsgroup submitted a multi-task learning network with four sub-networks that are used to classify each component of the action triplet and the triplet itself. Their model is based on two fundamental assumptions: (1) The type and location of surgical instruments are critical factors for determining the verb component of the triplet (Nwoye et al., 2020), and (2) it is important to combine the feature of the instrument, verb, and target (Nwoye et al., 2020; Jin et al., 2021). Therefore, their approach is composed of a CNN backbone (feature extraction), weak instrument localization, component feature fusion sub-network, and triplet-possibility mapping using weighted averaging. Firstly, a ResNet-18 backbone is trained to jointly learn to perform instrument, verb, and target prediction using three classifier heads. To focus the predictions on the relevant portions of the image, weak instrument localization features (Nwoye et al., 2019) are concatenated to the learned features for verb and target heads before performing the final classification. The predicted probability vectors for each of the three tasks are then concatenated before a single fully connected layer is used to map this vector to the final triplet prediction vector.

4.2.10. Team HFUT-MedIA: COEMNet: Correlation embedded multi-task network

Team HFUT-MedIA participated with a correlation embedded multi-task network, named COEMNet (see Fig. 5). First, a multi-task learning network is trained for instrument, verb, and target recognition tasks. The learned features for instrument prediction are leveraged for better recognition of verbs and targets. Secondly, the correlation between all classes is modeled as a graph and used the statistical co-occurrence frequencies as the

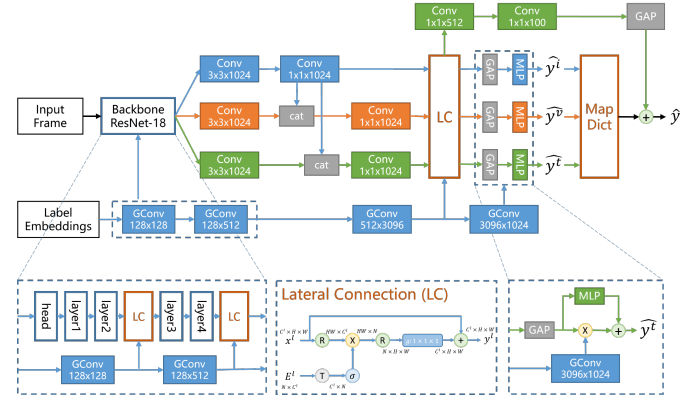


Fig. 5. Overview of the HFUT-MedIA submission.

adjacency matrix. A multi-layer graph convolutional network (GCN) is deeply integrated into their end-to-end network to efficiently learn the label classifiers and embed correlation information into feature representation (Wang et al., 2020c). Since associating the predicted instruments, verbs, and targets to form the right triplet predictions is complex when there are multiple triplet annotations in a single frame, a triplet adjustment task is added to the proposed network to minimize the association errors. The learned adjustment factors are applied to the original triplet prediction for more accurate recognition.

4.2.11. Team HFUT-NUS: Multi-task learning-based surgical interaction triplet recognition

Team HFUT-NUS’s method utilizes a ResNet-18 (He et al., 2016) backbone to extract features from endoscopic images followed by three parallel fully connected layers that are used to predict instrument, verb, and target, respectively. To further improve the precision of surgical interaction triplet detection, an attention mechanism is adopted to weigh the importance of instrument, verb, and target components. Finally, the weighted instrument vector, verb vector, and target vector are spliced together and fed into a fully connected layer with a sigmoid activation to predict the final surgical action triplet.

4.2.12. Team J&M: Surgical action triplet recognition with Efficient-MSTCN

Team J&M’s submission is a 2-stage network that is trained to first extract relevant spatial features and then utilize the temporal context of each frame to make action triplet predictions for a frame. In their proposed model, EfficientNetV2-M (Tan and Le, 2021) is used as the feature extractor. It is trained to make predictions using a single frame. Afterward, the weights are frozen and a Multi-Stage Temporal Convolutional Network (MS-TCN) (Farha and Gall, 2019) is employed to refine the model predictions.

4.2.13. Team Med Recognizer: Surgical action triplet recognition via temporal memory relation gradient network

Team Med Recognizer’s competing model is a temporal memory relation gradient network, in which a stem module first extracts spatial features from each frame, and then splits into three branches, with different temporal supportive information integrated to represent the action triplet. The temporal lengths are set as 10, 5, and 5 for ‘verb’, ‘instrument’, and ‘target’, respectively, given that different tasks require different amounts

of temporal context. The spatial-temporal feature extraction model is developed based on Jin *et al.* (2021), which is originally designed for phase recognition, to leverage the long-range and multi-scale temporal patterns in the video. The features are then fed into the classifiers to generate the prediction probabilities of the three tasks. The obtained probabilities from three branches are then integrated to produce the final one. Post-processing methods are then applied to account for the label imbalance, where higher weights are assigned to the classes with fewer data samples. The temporal information is also further utilized to weight the predicted probabilities of the previous frames when producing the results of the current frame. Dropout is used at both the training and the testing time.

4.2.14. Team MMLAB: *Temporal triplet net for triplet presence detection in surgical videos*

Team MMLAB proposed the Temporal Triplet Net (TTN) which consists of DenseNet (Huang *et al.*, 2017) and a Graph Convolutional Network (GCN) that utilizes both spatial and temporal features for the recognition of action triplets in surgical videos. The DenseNet is used as an image classification model to extract spatial features for each labeled frame for an input video. To utilize the temporal information, the GCN is incorporated to capture the temporal relationships among the continuous frames of a video sequence using the features extracted for each frame. Specifically, the representation for each frame is regarded as the node of the graph, and the similarity between each pair of nodes is regarded as the edge of the graph.

4.2.15. Team NCT-TSO: *Multi-task learning framework for action triplet recognition*

Team NCT-TSO's method is similar to Nwoye *et al.* (2020) designed for the tasks of verb, target, and triplet recognition. In their proposed method, a separate model is trained for instrument recognition using the Cholec80 dataset (Twinanda *et al.*, 2016). The instrument recognition model is based on a convolutional neural network (CNN) which uses the ResNet-50 (He *et al.*, 2016) as its backbone and spatial pooling to learn class-specific feature maps of the instruments in a weakly supervised manner (Durand *et al.*, 2017). These instrument maps are subsequently fed to the verb and target paths of the triplet recognition network. The verb and target paths share the same ResNet-50 model as their backbone, followed by two convolutional layers for each path. The verb and target paths are also trained to learn class-specific feature maps and use wildcat spatial pooling (Durand *et al.*, 2017) on these maps for the prediction of labels. The instrument logits from the pre-trained instrument model, along with the verb and target logits, are subsequently used to learn the triplets using a 3D interaction space as proposed in the work of Nwoye *et al.* (2020).

4.2.16. Team SIAT CAMI: *Multi-task mutual Channel Recurrent Net for Fine-grained Surgical Triplet Recognition*

Team SIAT CAMI used MT-MCLNet (see Fig. 6), a multi-task surgical triplet recognition network with multi-label mutual channel loss (Chang *et al.*, 2020) to extract local fine-grained features and aggregate temporal information on verb and triplet branch. When a sequence of video images is fed into the network, the backbone ResNet-50 module (He *et al.*, 2016) first extracts a global 2048-dimension spatial feature for each image in the sequence. To extend the global features to each

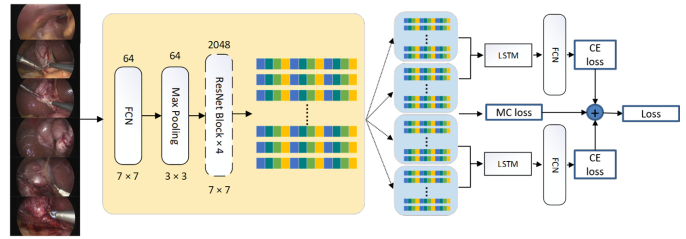


Fig. 6. Overview of the SIAT CAMI submission.

task, different 1×1 convolutions are applied to generate 2040-dimension features for instrument and target branches, 2000 for verb and triplet. Then, the 2040/2000-dimension spatial features are fed into two LSTM modules to capture 512-dimension temporal motion information. Then, 4 fully connected layers are utilized to produce the final classification output. The overall loss function is a weighted sum of standard cross-entropy loss and a mutual channel loss, which is intended to encourage learning a diverse range of discriminative features.

4.2.17. Team SJTU-IMR: *Tracking surgical actions with transformers and action label-guided fine-grained information aggregation*

Although each triplet is unique, different triplets may share certain elements. Team SJTU-IMR hypothesized that this observation is critical to modeling instrument-tissue interactions as it can establish a link between different actions with shared elements. To this end, they presented a two-stage transformer-based learning framework to solve the surgical action tracking problem. Specifically, the features learned from the first stage, which conducts frame-level action recognition via a Swin-S network architecture (Liu *et al.*, 2021), are taken as the input to the second stage, which performs sequence-level action recognition via masked transformers (Vaswani *et al.*, 2017). At both stages, the same multi-task learning strategy of combining coarse-grained action recognition with fine-grained information aggregation is employed. Their fine-grained information aggregation is guided by coarse-grained action labels so that the networks can put more emphasis on features modeling instrument-tissue interactions.

4.2.18. Team SK: *Action triplet recognition with weakly-supervised attention of surgical instruments*

Team SK followed an instrument-centric approach to action triplet recognition by first performing instrument recognition that is then used to condition the verb, target, and consequently, triplet recognition tasks. Using the first 4 convolutional layers of a ResNet-18 (He *et al.*, 2016) as their backbone, the first sub-network generates attention maps using 2 convolutional layers (layer 5 of ResNet-18 + 1×1 convolution) to localize the position of the instruments in the image. These attention maps are implicitly learned through weak supervision with instrument labels. A second sub-network then uses a convolutional layer to predict a fixed number of channels corresponding to each instrument which is multiplied by the previously learned attention maps to appropriately weigh different parts of the image based on instrument location. Finally, a 1×1 convolution, global average pooling, and softmax operation are performed to obtain the triplet probabilities.

4.2.19. Team Trequantista: *Phase-aware multitasking surgical action recognition model with adjustment for low-data triplet classes*

Team Trequantista's entry used multi-task learning on 5 tasks, triplet, instrument, verb, target, and phase prediction. The phase prediction ground truth made use of Cholec80 annotations (Twinanda et al., 2016). Using a ResNet-18 and ResNet-34 (He et al., 2016) as their primary models, their work focuses on finding the optimal hyperparameters and appropriately post-processing their predictions specifically for the challenge metric (mAP). Following an observation that predicting under-represented classes has a strongly negative effect on the overall metric, the triplet predictions probability for the 63 least represented triplets are adjusted to compensate for this effect based on the following heuristic:

$$P(\text{Triplet}) = P(\text{Instrument}) * 0.03P(\text{Verb}) + 0.97P(\text{Target}). \quad (1)$$

4.3. Theoretical comparative analysis

Presented methods are broadly classified into 5 categories:

4.3.1. Multi-task learning: *Rank 1-5, 7-8, 10-12, 14-19*

Multi-task learning methods aim to improve triplet recognition performance by learning a shared representational space by training a model to perform multiple related tasks. Unsurprisingly, given the availability of additional annotated information, all but 3 teams employed multi-task learning in their model design. The most common formulation is to learn to predict the triplet components, instrument-verb-target, in addition to the triplet itself. Interestingly, this formulation has come in two flavors: 1) Predicting the triplet as an explicitly modeled association of the 3 predicted components (2Ai, ANL Triplet, CAMP, HFUT-MedIA, HFUT-NUS, lsgroup, Med Recognizer, NCT-TSO, SJTU-IMR, and 2) Predicting the triplets and 3 components using a shared backbone that implicitly learns useful features from the other tasks for triplet recognition (CITI SJTU, Trequantista, SIAT-CAMI). Given that surgical phases (Twinanda et al., 2016) are defined by and consequently signal the occurrence of certain actions, surgical phase recognition is highly related to the task of action triplet recognition. Three teams (Trequantista, 2Ai, CAMP) leverage this fact by incorporating spatial phase annotations for the challenge training videos that are publicly available in the Cholec80 dataset (Twinanda et al., 2016). All three methods do so in multi-task learning setups, with phase and triplet component recognition included as relevant tasks to boost triplet recognition performance. Finally, the Band of Broeders entry treated the triplet recognition tasks as an object detection task, positing that localization information is critical to effectively recognizing the action triplets represented in an image. They manually generated bounding boxes and instrument labels, which are then assigned to their corresponding triplets, to facilitate model training by learning to both detect and localize triplets.

4.3.2. Temporal modeling: *Rank 2, 4-6, 8-11, 18-19*

While surgical action triplet recognition can be done on single frames, temporal models processing information from several frames at a time are interesting solutions to investigate: object permanence, motion, and surgical workflow are indeed temporal concepts that can inform action recognition. The many teams (11 out of 19) choosing this direction proposed a

diverse range of temporal modeling methods, which can be described according to three main traits. The first is the choice of temporal architecture. Some methods used non-trainable operations to aggregate information from static image models across multiple frames (ANL-Triplet, LSGroup, CAMP, MMLAB); most entrants, however, resort to sequence models running on features extracted by CNN. Among them, LSTM-based recurrent neural networks are a clear trend, with various forms appearing in 6 submissions (2AI, CEAIK, Digital Surgery, CITI-SJTU, Med Recognizer, SIAT-CAMI). Team CITI-SJTU in particular used a pair of ConvLSTMs. Besides LSTMs, other more recent models appeared as well: TCNs (J&M, Med Recognizer), and Transformers (SJTU-IMR).

The second notable trait is the model's temporal range. In some methods, this range is very short: 2 frames for team LSGroup, 4 frames for team CEAIK, 5 frames for team ANL-Triplet, and team MMLAB. Longer periods are featured in submissions by team 2AI (10 frames), team Digital Surgery (10 frames), team CITI-SJTU (16 frames), and team CAMP (32 frames); team Med Recognizer used "short video clips" for the SV-RCNet, 10-frame clips for the verb memory bank, and 5-frame clips for the instrument memory bank. Team SJTU-IMR feed their model in chunks of 200 frames, which can cover large workflow sections. One team, J&M, used concatenated features from all video frames, making the entire history available to the model at any given timestep.

The third trait distinguishing temporal methods is end-to-end training. Some methods feature temporal layers that are trained separately (CEAIK, Digital Surgery, J&M, Med Recognizer, SJTU-IMR), while others (2AI, ANL-Triplet, CITI-SJTU, SIAT-CAMI, LSGroup, CAMP) trained all parts simultaneously.

4.3.3. Attention mechanism/transformer: *Rank 3,5,9,10,12,14*

Attention mechanisms are methods designed to modulate a model's input or internal representation, to highlight parts that are important for making predictions. The use of attention can greatly improve surgical action triplet recognition, as shown by Nwoye et al. (2022) in the two approaches of Attention Tripnet and Rendezvous - these two models reappear in this challenge as baselines. All methods from the challenge featured in this category used some form of spatial attention since areas surrounding tooltips are particularly informative. A few of these methods explicitly used instrument maps as the source of attention (Attention Tripnet, RDV, LSGroup, SK). The other forms of spatial attention featured are team SIAT-CAMI's "channel-wise attention", team Casia Robotics's "coordinate attention" (Hou et al., 2021b) and team SJTU-IMR's Swin-S vision Transformer. Two entries used additional attention components of different types: semantic, in the RDV baseline's multi-head of mixed attention; and temporal, in the masked Transformer appearing in the second stage of team SJTU-IMR's method.

4.3.4. Graph convolutional networks: *Rank 3, 13*

Two teams, HFUT-MedIA and MMLAB, made use of GCN and CNN to better recognize action triplets represented in a given image, using two different strategies. Team MMLab used a GCN primarily to leverage temporal relationships using spatial features extracted using a CNN in the first stage of training. Here, the extracted features for each frame are treated as a graph node and the similarity between each pair of nodes

features an edge between two nodes. Team HFUT-MedIA, in contrast, does not base its graph on CNN-extracted spatial features but rather used it to effectively incorporate triplet co-occurrence distribution statistics into various stages of a discriminative CNN to predict action triplet probabilities.

4.3.5. Ensemble models: Rank 1, 6, 16

Ensembling is a commonly used strategy for limiting noisy predictions, by combining outputs from independently trained models. Three teams employed this type of approach, with distinct variations on the ensembling concept. The choice of operation for merging outputs varies between teams: summation (Digital Surgery), averaging (2AI version 2), majority vote (2AI version 1), ad-hoc heuristics (Trequartista). Early or late fusion is another differentiating trait: teams Trequartista and 2AI merged final probabilities, while team Digital Surgery combined features from various feature extractors before applying an LSTM model. Finally, various ensemble sizes and architectures are used: 3 ResNets (Trequartista), 4 HRNets (Digital Surgery), and an ensemble of 6 CNNs (team 2AI).

4.4. Ensemble predictions

For optimal performance on the dataset and as a summary of the challenge benchmark study, we ensemble the predictions of 7 top models (with triplet recognition AP above 30.0%): Trequartista, HFUT-MedIA, SIAT-CAMI, RDV, ANL-Triplet, CITI-SJTU, and Digital Surgery. This helps minimize errors due to noise, bias, and variance while improving the stability, reliability, and accuracy of predictions. In this work, we experiment with 6 ensemble methods.

As shown in Fig. 7, we start with simple **averaging** of the probability scores of the different models (Fig. 7a). This is extended to **weighted averaging** (Fig. 7b) where the weights are computed as the performance ratio of each model against the others. We also perform **soft voting** (Fig. 7c) between the maximum (presence) and minimum (absence) probability scores per class at a threshold of 0.5. Our first trainable method uses a two-layer network to learn a **deep ensemble** (Fig. 7d) of the challenge networks predictions. Lastly, we focus on learning the model averaging weights and in turn use this for a **deep weighted ensemble** (Fig. 7e). In this case, we learn two types of weights: (1) a vector of N weights for the N models, and (2) a matrix of $N \times C$ weights for per class (C) weights of N models. The latter option, which is a **deep per-class weighted ensemble**, is designed to utilize the strength of each model in recognizing the different categories of the triplets. All the deep ensemble models are trained on triplets probability (Y_{IVT}) predictions of the selected models on the training dataset.

4.5. Implementation details

The implementation details of all submissions as well as the baselines are summarized in Table 3.

5. Evaluation

5.1. Metrics

To evaluate the performance of the presented models on surgical action triplet recognition, we use the average precision (AP) metric, computed as the area under the precision-recall curve; this is recommended as the standard practice for the

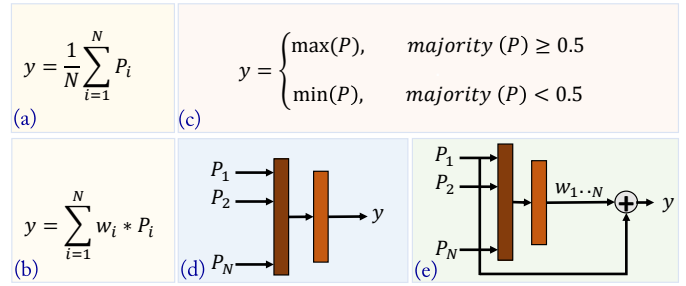


Fig. 7. Ensemble method for the model predictions: (a) averaging, (b) weighted averaging, (c) majority soft voting, (d) deep ensemble, and (e) deep weighted ensemble.

CholecT50 dataset. Using AP metrics, we assess the models' capacity for predicting the correct triplet components: AP_I , AP_V , AP_T , and the correct triplet associations: $AP_I V$, $AP_I T$, $AP_I VT$. The $AP_I VT$ evaluates the complete instrument-verb-target combination and hence serves as the primary metric in the challenge. We also analyze the quality of the model predictions by computing the topK performance scores where $K \in [5, 10, 15, 20]$ and the average over topK@[5:20] at intervals of 5 steps. All the evaluation scores are computed using the `ivtmetrics` library⁶ (Nwoye and Padoy, 2022).

5.2. Evaluation protocol

We maintain an online inference strategy to represent the real-time usage of developed methods in the OR. For uniform evaluation of all models - including those not multitasking the triplet components - we collect only the probability scores for triplet classes (Y_{IVT}) per frame as the model output. The individual component predictions are filtered from Y_{IVT} following the disentanglement function proposed in Nwoye *et al.* (2022). The video-specific AP scores are computed per category in a test video, then averaged categorically across all test videos. We obtain the mean AP by averaging the category APs.

For method ranking, the challenge evaluation is performed on 94 valid triplet classes excluding the 6 null classes in the dataset. The null classes, which include $\langle \text{grasper}, \text{null-verb}, \text{null-target} \rangle$, $\langle \text{bipolar}, \text{null-verb}, \text{null-target} \rangle$, $\langle \text{hook}, \text{null-verb}, \text{null-target} \rangle$, $\langle \text{scissors}, \text{null-verb}, \text{null-target} \rangle$, $\langle \text{clipper}, \text{null-verb}, \text{null-target} \rangle$, and $\langle \text{irrigator}, \text{null-verb}, \text{null-target} \rangle$, are possible only when an instrument is idle or performing an action that is not part of the 100 considered triplet classes. A total absence of a triplet, $\langle \text{null-instrument}, \text{null-verb}, \text{null-target} \rangle$, is not an explicit class in a multi-label classification problem.

6. Results and discussion

In this section, we provide a quantitative and qualitative overview of all methods featured in the challenge, including baselines and models evaluated post-challenge. In total there are 24 models: 4 are baseline models from the challenge organizers, 19 are competing models and 1 is a post-challenge submission. The baseline models provide lower bound performances for monitoring and assessing the improvement and effectiveness of newer implementations. The analysis of their results builds a foundation for the CholecT50 dataset, establishing it as a validated reference benchmark.

⁶<https://pypi.org/project/ivtmetrics>

Table 3. Methodological and implementation details.

Team	Architecture	Backbone	Multi-task	Temporal component	Attention	Output components	Post-processing	Data augmentation	Transfer learning	Optimizer	Loss
Trequartista	CNN	ResNet18 ResNet34	Instrument verb target triplet phase	N/A	N/A	Probability vectors for phases, Instruments, verbs, targets, triplets	63 rare triplets inferred using ad-hoc formula.	Fast AI default image augmentation	Cholec80 phase	Adam	Focal Loss + pair Loss focusing on ranking
2AI	CNN ensemble-LSTM	EfficientNetB0 EfficientNetB4 ResNet50 ResNest50 ResNest101 SENet	Instrument verb target phase	LSTM	None	Probability vectors for phases, instruments, verbs, targets	Temporal smoothing	crop, horizontal flip; color jitter, 10 deg. rotation	Cholec80 phase	SGD	CE (Cross-entropy)
SIAT-CAMI	CNN-LSTM	ResNet50	Instrument verb target triplet	LSTM	Channel-wise attention	Probability vectors for Instruments, verbs, targets, triplets	N/A	crop, horizontal flip, rotation, color jitter	N/A	Adam	Mutual channel loss + CE
HFUT-MedA	CNN-GCN	ResNet18	Instrument verb target triplet	N/A	N/A	Probability vectors for triplets	N/A	rotation, flip, color jitter, gaussian blur, erasing, cutmix	Imagenet	AdamW	Weighted CE
RDV	CNN-CAGAM-MHMA	ResNet18	instrument verb target triplet	N/A	Spatial (CAGAM) + semantic (MHMA)	Probability vectors for Instruments, verbs, targets, triplets	N/A	Random scaling, brightness	Imagenet	Momentum SGD	Weighted CE
CITI-SJTU	CNN-ConvLSTM	ResNet50	Instrument verb target triplet	ConvLSTM	N/A	Probability vectors for Instruments, verbs, targets, triplets	N/A	horizontal flip, rotation, color jitter	Imagenet	SGD	Weighted CE
ANL Triplet	CNN	ResNet18 ResNet34	Instrument verb target triplet	Previous 4 frames injection	N/A	Probability vectors for Instruments, verbs, targets, triplets	Temporal averaging	FastAI	Unspecified pretraining for Resnet34	AdamW	CE
Digital Surgery	CNN ensemble-LSTM	HRNet32	N/A	LSTM	N/A	Probability vector for triplets	N/A	illumination, color, blur, noise	Imagenet	Adam	Smoothed NLL + flattened CE
Casia Robotics	CNN	ResNet18	Instrument verb target	N/A	Coordinate attention	Probability vectors for Instruments, verbs, targets, triplets	N/A	rotation, flip, patch masking	N/A	Adam	CE
LSGroup	CNN (4 subnetworks)	ResNet18	Instrument verb target triplet	Previous frame injection	Class Active Mapping	Probability vectors for triplets	N/A	horizontal flip	Imagenet.	SGD	CE
J&M	CNN-TCN	EfficientNetV2-M	N/A	MS-TCN / MS-TCN++	N/A	Probability vectors for triplets	N/A	N/A	Imagenet	SGD (EfficientNetV2-M), Adam (MS-TCN)	CE
Attention Triplet	CNN-CAGAM-3DIS	ResNet18	Instrument verb target triplet	N/A	Spatial (CAGAM)	Probability vectors for Instruments, verbs, targets, triplets	N/A	Flip, brightness	Imagenet	Momentum SGD	Weighted CE
CEAIK	CNN-LSTM	ResNet50	instrument verb target triplet	LSTM	N/A	Probability vector for triplets	N/A	Blur, brightness, contrast, rotation, sun flare	N/A	Adam	Multi-label soft margin loss
SJTU-IMR	2-stage Transformer	Swin-S	instrument verb target triplet	Masked transformers (stage 2)	Transformer multi-head self-attention	Probability vectors for triplets	N/A	crop, flip, color shift, temporal warp	N/A	AdamW	CE
Tripnet	CNN-CAG-3DIS	ResNet18	instrument verb target triplet	N/A	N/A	Probability vectors for Instruments, verbs, targets, triplets	Invalid triplet masking	Flip, brightness	Imagenet	Momentum SGD	Weighted CE
SK	CNN	Resnet18	Instrument triplet	N/A	Instrument maps	Probability vector for instruments, triplets	N/A	Shift, scaling, rotation, color jitter, blur, noise, crop	Imagenet	Adam	Weighted CE
MMLAB	CNN-Graph conv	DenseNet121	No	Previous 4 frames injection	N/A	Probability vector for triplets	N/A	Horizontal flip	Imagenet.	SGD (DenseNet), Adam (GCN)	CE
Band of Broeders	YOLOv5 object detector	CSP-DenseNet	Triplet detection & localization	N/A	N/A	Bounding boxes, triplet class probabilities, objectness confidence scores	Overlapping bounding boxes suppression	YOLOv5 augmentation	COCO	SGD	CE
MTL baseline	CNN	ResNet18	Instrument verb target triplet	N/A	N/A	Probability vectors for Instruments, verbs, targets, triplets	N/A	Flip, brightness	Imagenet	Momentum SGD	Weighted CE
NCT-TSO	CNN	ResNet50	Target triplet	N/A	N/A	Probability vectors for triplets	Invalid triplet masking	N/A	Cholec80	Adam	Weighted CE
HFUT-NUS	CNN	ResNet18	Instrument verb target	N/A	N/A	Probability vectors for Instruments, verbs, targets, triplets	N/A	Horizontal flip, rotation	N/A	Adam	Weighted CE
CAMP	CNN	SlowFast50	Instrument, verb target phase	SlowFast	N/A	Probability vectors for Instruments, verbs, targets, triplets	N/A	Scale, crop, horizontal flip	Training on Phase labels from Cholec80	Adam	Weighted CE
Med Recognizer	SVRCNet-TMRNet	ResNet50	Instrument verb target	LSTM, Memory bank, temporal variation layer	N/A	Probability vectors for triplets	N/A	Crop, flip	Cholec80	SGD	CE

Table 4. Performance summary of the presented methods across all task divisions

Team	Component detection			Triplet association			Challenge ranking
	AP_I	AP_V	AP_T	AP_{IV}	AP_{IT}	AP_{IVT}	
Trequantista	79.9	52.9	46.4	<u>39.0</u>	41.9	38.1	1
2AI: Version 2 †	79.8	50.1	42.8	35.2	<u>42.4</u>	<u>36.9</u>	-
SIAT CAMI	82.1	<u>51.5</u>	<u>45.5</u>	37.1	43.1	35.8	2
HFUT-MedIA	77.1	46.7	37.8	33.1	35.9	32.9	3
RDV (CAMMA) ‡	77.5	47.5	37.7	39.4	39.6	32.7	-
CITI SJTU	67.8	37.1	34.8	29.9	33.0	32.0	4
ANL Triplet	73.6	47.3	40.5	32.6	37.1	31.9	5
Digital Surgery	<u>80.8</u>	50.0	41.1	35.1	35.7	31.7	6
Casia Robotics	<u>72.6</u>	43.9	31.2	30.7	30.6	26.7	7
Lsgroup	73.8	44.3	34.9	31.4	31.9	26.3	8
J&M	69.4	46.7	39.2	28.9	28.8	25.6	9
Attention-Tripnet (CAMMA) ‡	77.1	43.4	30.0	32.3	29.7	25.5	-
Ceaiik	68.9	40.5	30.9	27.5	28.4	25.2	10
SJTU-IMR	72.6	42.5	34.1	29.2	26.4	24.8	11
Tripnet (CAMMA) ‡	74.6	42.9	32.2	27.0	28.0	23.4	-
SK	52.6	30.4	20.2	25.8	21.0	18.4	12
MMLAB	50.1	31.8	31.6	20.6	22.1	18.1	13
Band of Broeders ¶	63.4	35.2	26.2	19.7	18.6	16.0	-
MTL baseline (CAMMA) ‡	48.6	27.8	19.8	18.5	15.5	13.7	-
NCT-TSO	27.3	15.7	12.3	13.6	11.7	10.4	14
2AI: Version 1	46.2	24.4	20.5	13.3	12.3	10.0	15
HFUT-NUS	34.1	20.2	13.5	16.0	11.2	09.8	16
CAMP	30.4	19.5	11.8	13.2	09.7	09.3	17
Med Recognizer	20.6	12.8	10.1	07.0	04.5	04.2	18
Mean ± standard deviation (stdev)	62.5±18.9	37.7±12.41	30.2±11.0	26.3±8.9	26.5±11.3	23.3±9.9	

bold = best score and underlined = second best. Not eligible for award: † post-challenge submission, ‡ organizers' baselines, ¶ used non-public third-party dataset.

6.1. Summary of the quantitative results

For a concise overview, we first summarize the AP scores on both the component detection and triplet association in Table 4. On the instrument component, half of the presented models achieved an AP score higher than 70% with the highest score of 82.1% by team SIAT-CAMI and the average performance (with standard deviation) of $62.5 \pm 18.9\%$. These results suggest the tremendous progress made in the use of deep learning models for surgical instrument recognition in laparoscopic videos, as 18 out of the 24 presented methods recognized the instruments at a performance higher than 50%. The verb recognition peaked at 52.9% AP with only four teams achieving a higher than 50% score. It is observed that these four teams either leveraged temporal information or exploited phase labels, which is also a temporal task. This suggests a strong correlation between surgical phases and actions and the quality of temporal feature modeling in tracking activity workflow. Improving the verb performance is a promising direction for future work likely by exploring a better-conditioned range of temporal dependencies attuned to triplet timings. The bulk of the AP scores from the competing teams falls within 30-40% with a mean of $37.7 \pm 12.1\%$.

Target, being the most difficult component, was recognized at a maximum AP of 46.4% by team Trequantista. The low performance on this sub-task can be attributed to the *instrument-centric* nature of the underlying target which makes its recognition very challenging. The mean performance at the challenge was $30.2 \pm 11.0\%$. As evidenced by the submissions, tackling surgical target recognition still is not straightforward.

On the association part, an interesting observation is that top-performing models recognized the instrument-target (AP_{IT}) pair better than they recognized the instrument-verb (AP_{IV})

pair, despite the larger number of classes for the former. Their drop in performance from AP_T to AP_{IT} is much lower than from AP_V to AP_{IV} . These deep learning models were more likely to recognize the correct operating instruments given an underlying target (a spatial relationship) than given their actions (a temporal relationship). The reverse is the case for lower-performing models. Here, it is easier to recognize instrument-verb pairs as most instruments perform specific actions and the verbs have a lesser number of classes than targets, and therefore their combinations (IV or IT). Many instruments can act on a wide range of targets - including unintended contact - making the problem more challenging for less advanced models. Overall these observations give a different perspective when interpreting the strength of the proposed models in understanding tool-tissue interactions. As shown in Table 4, team SIAT-CAMI obtained the highest AP_{IT} score while the baseline Rendezvous Nwoye *et al.* (2022) still retained the state-of-the-art (SOTA) performance on AP_{IV} .

The complete triplet recognition was best achieved at an AP of 38.1% by team Trequantista topping the challenge leaderboard. The second was a model submitted post-challenge by 2AI as an improvement of their challenge competing method. Meanwhile, team SIAT-CAMI claimed the runner-up prize with an AP of 35.8% while team HFUT-MedIA took the third-place prize with an AP of 32.9%. Three other teams (CITI-SJTU, ANL Triplet, and Digital Surgery) and the baseline Rendezvous achieved an AP higher than 30% which are promising performances for 100 class triplet recognition tasks. The mean performance for the triplet recognition recorded at the CholecTriplet 2021 challenge was $23.3 \pm 9.9\%$ suggesting room for improvement on this challenging task.

Table 5. Top K accuracy of the triplet predictions

Team	Top 5	Top 10	Top 15	Top 20	Top {5:20}
RDV ‡	69.35	84.38	89.93	93.24	84.23
Tripnet ‡	67.89	83.99	90.76	93.65	<u>84.07</u>
HFUT-MedIA	65.05	85.35	<u>91.75</u>	<u>93.59</u>	83.94
Attention-Tripnet ‡	66.86	82.49	91.85	93.25	83.61
Trequantista	<u>68.50</u>	82.40	88.24	92.29	82.86
Ceaiik	66.02	81.34	89.74	93.40	82.63
Digital-Surgery	65.97	81.56	88.78	92.92	82.31
HFUT-NUS	65.71	84.18	88.68	90.43	82.25
SIAT-CAMI	66.58	81.93	88.59	91.84	82.24
SJTU-IMR	66.50	81.88	84.19	84.89	79.37
ANL-Triplet	52.12	83.65	89.19	91.37	79.08
CITI-SJTU	54.95	78.61	88.96	92.41	78.73
SK	48.08	79.46	90.41	92.12	77.52
2AI: Version 2 †	64.87	76.17	82.54	86.26	77.46
Casia-Robotics	59.11	75.03	84.44	90.41	77.25
MMLAB	60.53	76.57	82.72	86.67	76.62
Lsgroup	62.88	73.03	78.64	81.98	74.13
J&M	56.09	66.36	72.36	76.90	67.93
MTL-baseline ‡	45.59	52.45	56.65	59.77	53.62
Med-Recognizer	30.26	43.69	58.24	68.05	50.06
Band-of-Broeders ¶	39.86	40.34	41.14	48.27	42.40
2AI	29.44	30.18	32.11	41.14	33.22
CAMP	08.73	17.91	21.25	25.71	18.40
NCT-TSO	04.88	11.78	16.59	21.93	13.80
Mean ± stdev	53.1±18.3	67.5±22.3	73.9±23.3	77.9±22.0	68.1±21.2

bold = best score and underlined = second best. Not eligible for award: ‡ organizers' baselines,

† post-challenge submission, ¶ used non-public third-party dataset.

6.2. TopK accuracy on surgical action triplet recognition

Due to the large number of classes and the high semantic overlap in the triplet classes, we also evaluate the topK of the presented models. This metric measures the ability of a model to predict the exact triplets within its top K confidence scores. We analyze the top 5, 10, 15, 20, and average across these four thresholds, topK@[5:20] as shown in Table 5. The obtained results describe the model's confidence in its predictions with the Rendezvous model obtaining a 69.35% accuracy score as the best model when the top 5 confident predictions are considered. Similarly, HFUT-MedIA, Attention Tripnet, and Tripnet models produce the best results at top 10, 15, and 20 confidence scores respectively.

On average, the Rendezvous model can correctly recognize the triplets at a performance of 84.23% when top confident predictions are taken into account. The average performance of all the presented models at the challenge is 53.1±18.3%. On this metric, it is not surprising to see the challenge winner, Trequantista, scoring lower in top K accuracy because the model uses a mathematical operation to suppress the less confident predictions. These low confidence scores, when taken into account by the AP metric, lower the performance of other models. With topK focusing only on top confident predictions, the models are not penalized by their less confident predictions. This accuracy metric is highly informative and more usable when thresholding predictions to binary values to obtain the unique IDs of the predicted triplets. It also suggests that most of the presented models would ordinarily obtain higher scores with fewer triplet classes, less class similarity, and less semantic overlap.

Meanwhile, the top K accuracy increases with the K tolerance as seen in Table 5.

Table 6. Per-class performance on instrument presence detection

Team	Grasper	Bipolar	Hook	Scissors	Clipper	Irrigator	Mean
SIAT CAMI	95.8	92.8	97.5	94.9	81.1	28.7	82.1
Digital Surgery	94.9	94.2	<u>98.4</u>	<u>92.8</u>	85.7	16.6	<u>80.8</u>
Trequantista	95.1	91.3	98.1	86.3	81.5	25.4	79.9
2AI: Version 2 †	96.8	88.2	98.3	88.4	81.5	23.5	79.8
RDV ‡	95.1	90.1	98.2	89.0	79.5	10.6	77.5
HFUT-MedIA	93.0	83.1	95.9	84.7	81.4	22.2	77.1
Attention Tripnet ‡	95.4	87.9	98.6	88.5	78.8	10.8	77.1
Tripnet ‡	86.7	82.3	97.6	79.4	80.3	19.4	74.6
Lsgroup	91.6	85.3	96.8	76.3	76.5	14.0	73.8
ANL Triplet	88.3	68.6	96.6	84.0	<u>82.4</u>	19.8	73.6
Casia Robotics	92.0	88.2	97.7	67.7	72.1	15.7	72.6
SJTU-IMR	85.7	89.7	97.4	65.7	76.6	18.7	72.6
J&M	91.7	72.9	96.6	48.4	76.7	28.7	69.4
Ceaiik	88.1	84.9	98.0	52.5	66.8	21.1	68.9
CITI SJTU	92.4	92.1	66.6	62.7	78.4	12.6	67.8
Band of Broeders ¶	91.6	55.5	94.2	66.8	66.8	03.6	63.4
SK	57.3	72.5	30.6	61.5	70.6	22.3	52.6
MMLAB	86.9	44.2	88.8	28.6	31.0	19.5	50.1
MTL baseline ‡	81.5	58.9	93.2	13.8	37.8	04.4	48.6
2AI: Version 1	83.1	57.7	91.4	11.7	28.1	03.5	46.2
HFUT-NUS	49.8	80.0	58.3	03.3	06.5	05.9	34.1
CAMP	61.2	15.6	72.0	03.5	20.4	08.9	30.4
NCT-TSO	39.0	81.7	25.7	05.1	05.6	05.7	27.3
Med Recognizer	56.3	09.2	47.0	03.0	04.1	03.3	20.6
Mean ± stdev	82.9±16.7	73.6±23.1	84.7±22.4	56.6±33.2	60.4±28.4	15.2±8.1	62.5±18.9

bold = best score and underlined = second best. Not eligible for award: ‡ organizers' baselines,

† post-challenge submission, ¶ used non-public third-party dataset.

6.3. Per-class component detection AP

Beyond the broad overview of the ranked performances, we present a detailed analysis of per-class performance for each component task.

On surgical instrument presence detection, the most frequently used instruments, *hook* and *grasper*, are the most correctly detected, as shown in Table 6. Their recognition APs are above 90% for half of the methods and their overall mean performances are above 82%. The suction irrigation device (*irrigator*) is only used when the field is unclear, resulting in a low usage frequency and mean performance of 15.2%. The scissors with the highest standard deviation of ±33.2 is the most complicated instrument to recognize as it often confounded with other instruments such as *grasper*, *bipolar*, and *clipper*.

Table 7 presents the per-class performance for the verbs. The most frequently used verbs such as *grasp*, *retract*, *dissect* are detected above 50.0% by the top models and above 40% on average. Verbs such as *dissect*, *coagulate*, *clip*, *cut*, which have the strongest affinity with a particular instrument class, are detected above 70% by the top models and with a higher average challenge performance. This confirms that triplets are instrument-centric. The average performance for *cut* is approximately 50% of the top team score. This is likely affected by the low detection of the performing instrument, *scissors*. In cases where an instrument has multiple frequent verbs, the performance tends to spread out over those verbs according to their prevalence: for example $retract \approx grasp \gg pack$ for *grasper*, $aspirate \gg irrigate$ for *irrigator*, etc. *Irrigate* is the least detected verb, likely due to its temporal nature as it can only be distinguished from *aspirate* based on the temporal dynamics of the fluid. Remarkably, team HFUT-MedIA leveraged graph convolution networks, notable for temporal action detection, to better detect this verb.

Table 7. Per-class performance on verb recognition.

Team	Grasp	Retract	Dissect	Coagulate	Clip	Cut	Aspirate	Irrigate	Pack	Null	Mean
Trequantista	54.0	<u>55.5</u>	79.6	70.6	80.7	81.6	20.9	01.8	<u>48.8</u>	30.2	52.9
SIAT CAMI	<u>56.0</u>	<u>46.7</u>	69.2	72.8	81.2	<u>74.2</u>	<u>28.1</u>	01.9	<u>48.7</u>	31.7	<u>51.5</u>
2AI: Version 2 †	<u>56.0</u>	46.8	<u>78.3</u>	68.5	80.6	<u>72.8</u>	<u>17.6</u>	01.8	44.8	28.6	50.1
Digital Surgery	53.2	45.7	68.3	70.6	86.2	74.1	18.4	03.0	49.3	26.4	50.0
RDV ‡	52.4	48.5	73.2	69.5	80.2	70.5	09.9	01.2	37.1	28.2	47.5
ANL Triplet	44.6	60.7	73.7	62.2	<u>82.4</u>	69.1	11.4	00.9	37.0	26.9	47.3
HFUT-MedIA	56.8	41.8	68.2	64.0	81.6	67.2	17.2	16.2	20.4	<u>30.3</u>	46.7
J&M	47.4	44.9	75.0	<u>71.6</u>	77.9	36.4	37.8	<u>11.2</u>	37.7	<u>24.3</u>	46.7
Lsgroup	50.3	46.0	73.5	65.1	76.6	63.1	17.3	00.5	21.1	25.3	44.3
Casia Robotics	46.8	43.7	72.1	65.8	71.2	63.8	12.0	08.4	23.5	27.8	43.9
Attention Tripnet ‡	53.2	39.4	71.4	65.1	79.2	68.4	09.1	02.8	18.1	22.9	43.4
Tripnet ‡	48.9	48.0	70.2	67.5	79.4	60.2	19.6	00.9	08.7	21.5	42.9
SJTU-IMR	57.6	47.1	74.8	69.9	76.6	43.0	14.5	00.4	10.4	26.1	42.5
Ceaiik	50.0	43.5	75.4	61.7	66.1	40.8	14.1	03.3	23.1	23.1	40.5
CITI SJTU	54.5	45.2	72.9	66.9	67.0	04.2	12.1	00.5	20.6	23.7	37.1
Band of Broeders ¶	44.1	35.9	68.9	55.0	67.9	52.0	02.7	00.4	01.2	20.0	35.2
MMLAB	49.0	38.0	61.0	43.3	32.4	19.1	11.5	03.9	33.8	23.4	31.8
SK	47.4	23.0	53.3	50.8	70.7	02.3	17.9	00.5	12.7	22.7	30.4
MTL baseline ‡	37.7	38.7	66.6	47.2	39.7	12.5	05.0	00.3	09.3	17.9	27.8
2AI: Version 1	43.5	34.9	58.1	39.9	28.9	13.6	02.5	00.4	01.2	18.7	24.4
HFUT-NUS	40.9	21.7	39.4	59.6	06.4	04.5	03.6	01.1	04.2	18.1	20.2
CAMP	42.6	26.0	45.5	21.4	23.5	06.2	06.3	00.3	04.7	16.9	19.5
Naive CNN ‡	34.4	34.2	55.5	16.4	07.8	03.4	06.3	00.8	03.4	19.0	18.3
NCT-TSO	24.9	16.5	21.5	55.6	05.2	05.6	07.2	00.6	00.8	17.8	15.7
Med Recognizer	35.6	21.1	32.3	07.4	01.7	01.9	02.5	00.4	01.9	21.8	12.8
Mean ± stdev	47.8±7.8	40.0±11.2	64.3±15.2	58.0±16.3	60.1±28.1	42.0±29.2	13.3±8.5	2.6±3.9	21.6±16.7	23.9±4.3	37.7±12.1

bold = best score and underlined = second best. Not eligible for award: † post-challenge submission, ‡ organizers' baselines, ¶ used non-public third-party dataset.

Table 8. Per-class performance on target recognition.[§]

Team	Gallbladder	Cystic-duct	Cystic-artery	Blood-vessel	Fluid	Abdominal-wall or cavity	Liver	Omentum	Peritoneum	Gut	Specimen-bag	Null	Mean
Trequantista	<u>91.4</u>	<u>66.8</u>	30.6	33.0	20.9	<u>60.8</u>	82.5	00.7	<u>37.5</u>	08.4	<u>89.2</u>	30.2	46.4
SIAT CAMI	89.4	64.3	<u>31.2</u>	15.4	<u>28.1</u>	67.3	79.5	00.6	<u>32.6</u>	<u>16.3</u>	<u>85.4</u>	31.7	<u>45.5</u>
2AI: Version 2 †	89.4	67.3	<u>26.1</u>	33.0	<u>17.6</u>	43.1	75.0	00.8	24.4	<u>14.8</u>	<u>89.2</u>	28.6	42.8
Digital Surgery	93.3	61.2	29.4	02.1	18.4	30.9	82.5	00.6	39.9	10.5	92.6	26.4	41.1
ANL TRIPLET	86.5	65.5	26.9	27.4	11.4	35.9	69.9	00.4	33.9	18.3	78.6	26.9	40.5
J&M	87.3	43.1	27.3	03.8	37.8	56.7	83.0	00.3	14.5	15.2	73.6	24.3	39.2
HFUT-MedIA	85.0	61.0	24.5	25.0	17.2	26.5	69.1	00.9	15.5	08.1	86.1	<u>30.3</u>	37.8
RDV ‡	88.2	54.7	32.0	18.3	09.9	27.7	70.9	00.5	25.2	08.4	83.5	28.2	37.7
Lsgroup	88.3	51.4	21.6	22.6	17.3	24.9	63.4	00.5	20.4	08.1	70.7	25.3	34.9
CITI SJTU	84.8	55.2	23.6	04.8	12.1	44.8	64.7	01.0	06.7	13.5	78.1	23.7	34.8
SJTU-IMR	83.6	54.7	24.7	08.9	14.5	34.2	67.0	00.8	09.6	05.2	76.2	26.1	34.1
Tripnet ‡	82.2	49.4	24.7	07.0	19.6	16.0	68.6	01.2	09.1	02.4	80.6	21.5	32.2
MMLAB	79.6	45.3	18.1	14.2	11.5	10.4	51.3	00.6	26.9	16.0	77.6	23.4	31.6
Casia Robotics	84.3	48.9	16.7	05.9	12.0	21.2	59.5	03.9	06.0	12.1	71.8	27.8	31.2
Ceaiik	84.3	44.7	21.0	05.4	14.1	21.4	55.1	01.2	12.0	08.7	76.0	23.1	30.9
Attention Tripnet ‡	77.4	42.5	23.1	13.8	09.1	22.2	41.0	<u>02.0</u>	18.5	03.6	79.7	22.9	30.0
Band of Broeders ¶	83.9	42.6	08.2	01.7	02.7	03.7	74.6	00.3	03.7	01.1	67.1	20.0	26.2
2AI: Version 1	76.9	22.1	05.8	01.7	02.5	00.9	45.6	00.3	03.7	01.1	63.7	18.7	20.5
SK	71.5	33.4	12.4	03.9	17.9	03.5	24.9	01.1	02.7	03.5	42.4	22.7	20.2
MTL baseline ‡	76.2	21.3	11.2	01.9	05.0	03.0	42.5	00.3	04.9	02.2	48.3	17.9	19.8
HFUT-NUS	52.6	12.3	05.2	03.1	03.6	03.8	40.9	01.1	04.7	02.5	11.7	18.1	13.5
NCT-TSO	40.9	10.6	12.1	16.2	07.2	00.9	33.7	00.3	02.5	00.7	02.8	17.8	12.3
CAMP	58.6	10.8	05.0	02.7	06.3	01.5	14.6	00.6	02.6	04.3	16.6	16.9	11.8
Med Recognizer	56.4	12.0	04.9	01.6	02.5	00.9	10.1	00.5	03.0	00.7	05.3	21.8	10.1
Mean ± stdev	78.8±13.5	43.4±19.0	19.4±9.1	11.4±10.3	13.3±8.5	23.4±20.3	57.1±21.3	0.9±0.7	15.0±12.3	7.7±5.7	64.5±27.9	23.9±4.3	30.2±11.0

bold = best score and underlined = second best. Not eligible for award: † post-challenge submission, ‡ organizers' baselines, ¶ used non-public third-party dataset. § shows only the targets in test videos.

Finally, we analyze per-class performance on target recognition which appears to be the most challenging component to correctly detect. As shown in Table 8, the *gallbladder* and *specimen-bag* are the most recognized targets, with the top models exceeding 90.0% AP. Their average performance across all methods is above 64.0%.

Other targets such as *liver*, *cystic-duct*, *abdominal wall* and *cavity*, are moderately detected. This is likely due to their obvious nature and clearer boundaries compared to the less de-

tected ones. Interactions with them are easier to ascertain than interactions with much smaller structures such as *cystic-artery* and other *blood-vessels*. Within the *cystic-pedicle*, the *cystic-duct* is the most detected tubular structure. The *cystic artery* in itself is hard to differentiate from other *blood-vessels*. This shows how deceptively complicated the task of anatomical target detection could be. The *peritoneum*, which covers the entire cavity, appears as a transparent layer making it difficult to identify. The heavily super-classed *omentum*, *peritoneum*, and *gut*

Table 9. Wilcoxon signed-rank test of the competing teams for rank stability.

p-value	Proposed method																			
	SK	MedR	Digital Surgery	Band of Broeders	Trequartista	MMLAB	Lsgroup	CAMP	NCT-TSO	HFUT-Media	Casia Robotics	SIAT-CAMI	HFUT-NUS	Ceaiik	SITU-IMR	J & M	CITI-SITU	ANL-Triplet	2AI	
SK	-	0.001	0.001	0.050	0.001	0.076	0.003	0.024	0.001	0.001	0.006	0.001	0.339	0.013	0.026	0.005	0.001	0.001	0.001	0.001
MedR	0.001	-	0.001	0.001	0.001	0.001	0.001	0.001	0.624	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.624
Digital Surgery	0.001	0.001	-	0.001	0.927	0.001	0.005	0.001	0.001	0.106	0.002	0.600	0.001	0.001	0.002	0.005	0.053	0.452	0.001	0.001
Band of Broeders	0.050	0.001	0.001	-	0.001	0.001	0.001	0.539	0.001	0.001	0.001	0.001	0.327	0.001	0.001	0.001	0.001	0.001	0.001	0.001
Trequartista	0.001	0.001	0.927	0.001	-	0.001	0.001	0.001	0.001	0.092	0.001	0.716	0.001	0.002	0.005	0.006	0.151	0.374	0.001	0.001
MMLAB	0.076	0.001	0.001	0.001	0.001	-	0.665	0.001	0.001	0.019	0.873	0.001	0.016	0.412	0.600	0.524	0.080	0.018	0.001	0.001
Lsgroup	0.003	0.001	0.005	0.001	0.001	0.665	-	0.001	0.001	0.002	0.624	0.001	0.002	0.802	0.699	0.682	0.053	0.001	0.001	0.001
CAMP	0.024	0.001	0.001	0.539	0.001	0.001	0.001	-	0.009	0.001	0.001	0.001	0.172	0.001	0.001	0.001	0.001	0.001	0.001	0.009
NCT-TSO	0.001	0.624	0.001	0.001	0.001	0.001	0.001	0.009	-	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001
HFUT-Media	0.001	0.001	0.106	0.001	0.092	0.019	0.002	0.001	0.001	-	0.002	0.039	0.001	0.021	0.039	0.096	0.399	0.509	0.001	0.001
Casia Robotics	0.006	0.001	0.002	0.001	0.001	0.873	0.624	0.001	0.001	0.002	-	0.003	0.001	0.946	0.699	0.495	0.101	0.009	0.001	0.001
SIAT-CAMI	0.001	0.001	0.600	0.001	0.716	0.001	0.001	0.001	0.001	0.039	0.003	-	0.001	0.001	0.009	0.003	0.119	0.399	0.001	0.001
HFUT-NUS	0.339	0.001	0.001	0.327	0.001	0.016	0.002	0.172	0.001	0.001	0.001	0.001	-	0.002	0.001	0.002	0.001	0.001	0.001	0.001
Ceaiik	0.013	0.001	0.001	0.001	0.002	0.412	0.802	0.001	0.001	0.021	0.946	0.001	0.002	-	0.802	0.785	0.219	0.011	0.001	0.001
SITU-IMR	0.026	0.001	0.002	0.001	0.005	0.600	0.699	0.001	0.001	0.039	0.699	0.009	0.001	0.802	-	0.374	0.096	0.023	0.001	0.001
J & M	0.005	0.001	0.005	0.001	0.006	0.524	0.682	0.001	0.001	0.096	0.495	0.003	0.002	0.785	0.374	-	0.2100	0.053	0.001	0.001
CITI-SITU	0.001	0.001	0.053	0.001	0.151	0.080	0.053	0.001	0.001	0.399	0.101	0.119	0.001	0.219	0.096	0.210	-	0.322	0.001	0.001
ANL-Triplet	0.001	0.001	0.452	0.001	0.374	0.018	0.001	0.001	0.001	0.509	0.009	0.399	0.001	0.011	0.023	0.053	0.322	-	0.001	0.001
2AI	0.001	0.624	0.001	0.001	0.001	0.001	0.001	0.009	-	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001

Significant p -value ≤ 0.05

Not significant p -value > 0.05

are the least detected in this case.

For single-task objectives, models proposed by SIAT-CAMI and Digital Surgery have the best performances and seem the most suitable for surgical instrument recognition. Both teams utilized a very deep feature extraction backbone and multi-scale aggregation of local and global features for final classification. On the other hand, models presented by Trequartista, SIAT-CAMI, and 2AI would be the best bet for modeling surgical action recognition which could be attributed to their advanced exploitation of the video temporal information. They also show promising performances on the target recognition task.

Taking all the per-component analysis into account, it is observed that the target recognition rate has a larger impact on the overall triplet recognition. This is most likely due to the target being the most challenging component to be correctly recognized and classified by the presented models. A model would have a higher probability of recognizing the other components if it can correctly recognize the complex targets.

For the joint task, there is no clear trade-off in the sub-task modeling, instead, balancing their scores guarantees a better overall triplet recognition performance as is the case in the 2AI version 2 model (Tables 6,7,8 and 4).

6.4. Result ranking stability

To measure the rank stability, we employ Wilcoxon signed-rank test as a non-parametric alternative to the dependent samples t-test since our data is not multivariate normal and teams' predictions could present many outliers. Using this, we test each team's method against a null hypothesis (H_0) to ascertain the statistical significance of its performance over others. The H_0 states that the difference between the proposed method and the alternative methods has a mean signed-rank of zero. Each H_0 test of one team against another produces a p -value between 0 and 1 as a measure of its level of statistical significance with a smaller value showing stronger evidence to reject the null hypothesis. Typically, a $p \leq 0.05$ is considered statistically significant.

To perform this analysis, we sample $N = 30$ random batches of 100 consecutive frames to simulate a video clip evaluation and perform the Wilcoxon test by iteratively using each competing team as a proposed method against the rest of the teams as its alternatives. The obtained p -values, tabulated in Table 9, show that closely-ranked teams do not significantly improve each other methods (accept H_0), whereas distantly ranked teams show a significant difference in their model performances (reject H_0). At a 5% confidence level, we conclude that the CholecTriplet2021 challenge has assembled teams with both similar and diverse methods in both modeling and performance.

6.5. Model ensemble results

Table 10 demonstrates that combining decisions from multiple models indeed helps to improve their overall performance. A simple averaging leads to an additional 0.8% gain in triplet recognition AP and +1.1% AP when the models' contributions are weighted by their individual strengths. Although an ensemble builds a strong learner from a group of weak learners, voting is surprisingly ineffective in this regard. In comparison to non-trainable alternatives, training ensemble techniques help to better minimize noise, bias, and variance errors, resulting in superior performances. Learning a new averaging weight appears more efficient than training totally new prediction models (Deep ensemble), as shown in Table 10. Learning the model weights per task category rather than general class weights per model is even more fascinating and reaches the highest mAP of 42.4% for surgical action triplet recognition. In general, performance improvement with the model ensemble is obtained across all the six sub-tasks.

6.6. Qualitative results

To better analyze the quality of the detections, we visualize the top K predictions for each model on sample frames with K triplet instances. An easy example in Fig. 8 showcases an image frame with a triplet of the best-recognized instrument (*hook*), the best recognized verb (*dissect*), and best-recognized target (*gallbladder*). This case is correctly detected by $\approx 88\%$

Table 10. Performance summary of the ensemble methods in comparison with the top 7 methods at the challenge.

Method	Component detection			Triplet association			
	AP_I	AP_V	AP_T	AP_{IV}	AP_{IT}	AP_{IVT}	
Top Challenge Teams	Trequantista	79.9	<u>52.9</u>	46.4	39.0	41.9	38.1
	SIAT CAMI	82.1	51.5	<u>45.5</u>	37.1	43.1	35.8
	HFUT-MedIA	77.1	46.7	37.8	33.1	35.9	32.9
	RDV ‡	77.5	47.5	37.7	39.4	39.6	32.7
	CITI SJTU	67.8	37.1	34.8	29.9	33.0	32.0
	ANL Triplet	73.6	47.3	40.5	32.6	37.1	31.9
	Digital Surgery	80.8	50.0	41.1	35.1	35.7	31.7
Model Ensemble	Averaging	82.4	52.9	44.7	40.4	43.5	38.9
	Weighted Averaging	82.5	53.1	44.9	40.5	43.8	39.2
	Soft Voting	79.6	46.7	42.6	35.3	39.8	35.1
	Deep ensemble	71.4	37.3	28.6	30.5	30.5	30.3
	Deep weighted ensemble	81.9	51.5	44.0	39.3	43.1	<u>40.5</u>
Deep per-class weighted ensemble	81.4	52.2	46.4	40.0	42.9	42.4	

bold = best score and underlined = second best. ‡ = organizers' baseline

	Image + ground truth	Team	Top prediction	Team	Top prediction	Team	Top prediction		
(a.)		2AI: Version 1	Hook, dissect, gallbladder	CITI SJTU	Hook, dissect, gallbladder	CAMMA: MTL ‡	Hook, dissect, gallbladder		
		2AI: Version 2 †	Hook, dissect, gallbladder	Digital Surgery	Hook, dissect, gallbladder	NCT-TSO	Hook, dissect, gallbladder		
		ANL Triplet	Hook, dissect, gallbladder	HFUT-MedIA	Hook, dissect, gallbladder	CAMMA: RDV ‡	Hook, dissect, gallbladder		
		CAMMA: Attention Tripnet ‡	Hook, dissect, gallbladder	HFUT-NUS	Grasper, retract, gallbladder	SIAT CAMI	Hook, dissect, gallbladder		
		Band of Broeders	Hook, dissect, gallbladder	J&M	Hook, dissect, gallbladder	SJTU-IMR	Hook, dissect, gallbladder		
		CAMP	Grasper, retract, cystic-plate	Lsgroup	Hook, dissect, gallbladder	SK	Hook, cut, peritoneum		
		Casia Robotics	Hook, dissect, gallbladder	Med Recognizer	Hook, dissect, gallbladder	Trequantista	Hook, dissect, gallbladder		
		Ceaiik	Hook, dissect, gallbladder	MMLAB	Hook, dissect, gallbladder	CAMMA: Tripnet ‡	Hook, dissect, gallbladder		
		(b.)		2AI: Version 1	Grasper, retract, liver Bipolar, coagulate, liver	CITI SJTU	Grasper, retract, liver hook, cut, peritoneum	CAMMA: MTL ‡	Grasper, retract, liver Bipolar, coagulate, liver
				2AI: Version 2 †	Grasper, retract, liver Bipolar, coagulate, liver	Digital Surgery	Grasper, retract, liver Bipolar, coagulate, liver	NCT-TSO	Clipper, clip, blood-vessels Scissors, cut, liver
ANL Triplet	Grasper, retract, liver Bipolar, coagulate, liver			HFUT-MedIA	Grasper, retract, liver Bipolar, coagulate, liver	CAMMA: RDV ‡	Grasper, retract, liver Bipolar, coagulate, liver		
CAMMA: Attention Tripnet ‡	Grasper, retract, liver Bipolar, coagulate, liver			HFUT-NUS	Grasper, retract, liver Bipolar, coagulate, liver	SIAT CAMI	Grasper, retract, liver Bipolar, coagulate, liver		
Band of Broeders	Grasper, retract, liver Bipolar, coagulate, liver			J&M	Grasper, retract, liver Bipolar, coagulate, liver	SJTU-IMR	Grasper, retract, liver Bipolar, coagulate, liver		
CAMP	Grasper, retract, cystic-plate Bipolar, retract, cystic-pedicle			Lsgroup	Grasper, retract, liver Bipolar, coagulate, liver	SK	Grasper, retract, gallbladder Hook, cut, peritoneum		
Casia Robotics	Grasper, retract, liver Bipolar, coagulate, liver			Med Recognizer	Grasper, dissect, gallbladder Hook, dissect, gallbladder	Trequantista	Grasper, retract, liver Bipolar, coagulate, liver		
Ceaiik	Grasper, retract, liver Bipolar, coagulate, liver			MMLAB	Grasper, retract, liver Bipolar, coagulate, liver	CAMMA: Tripnet ‡	Grasper, retract, liver Bipolar, coagulate, liver		
(c.)				2AI: Version 1	Clipper, clip, cystic-duct	CITI SJTU	Hook, cut, peritoneum	CAMMA: MTL ‡	Clipper, null-verb, null-target
				2AI: Version 2 †	Clipper, clip, cystic-artery	Digital Surgery	Clipper, clip, cystic-artery	NCT-TSO	Scissors, cut, adhesion
		ANL Triplet	Grasper, retract, omentum	HFUT-MedIA	Clipper, clip, cystic-duct	CAMMA: RDV ‡	Clipper, clip, cystic-artery		
		CAMMA: Attention Tripnet ‡	Clipper, clip, cystic-duct	HFUT-NUS	Hook, dissect, gallbladder	SIAT CAMI	Clipper, clip, cystic-artery		
		Band of Broeders	Clipper, clip, cystic-artery	J&M	Clipper, clip, cystic-artery	SJTU-IMR	Clipper, clip, cystic-duct		
		CAMP	Clipper, clip, cystic-artery	Lsgroup	Clipper, clip, cystic-artery	SK	Grasper, retract, gallbladder		
		Casia Robotics	Clipper, clip, cystic-artery	Med Recognizer	Hook, dissect, gallbladder	Trequantista	Clipper, clip, cystic-artery		
		Ceaiik	Clipper, clip, cystic-duct	MMLAB	Scissors, cut, cystic-duct	CAMMA: Tripnet ‡	Clipper, clip, cystic-artery		

Fig. 8. Qualitative results visualizing triplet predictions: a cross-section of teams' top k predictions on an input image depicting k action triplets: (a.) easy case, (b.) moderate case, and (c.) difficult case. † = post-challenge submission, ‡ = organizers' baselines, ¶ = used non-public third-party dataset. green = incorrect prediction red = incorrect prediction.

of the teams. A moderately difficult case of a frame with multiple triplets recorded seven incorrect predictions. An image frame showing a single triplet involving a *cystic-artery* proves to be difficult for more than half of the teams. The incorrect prediction is mostly on the target component of the triplet.

Finally, we analyze the triplet recognition showing a sequential flow of tool-activity recognition. Here, we group all the triplets by their instruments and use a range of colors to show their variations and transitions in temporal order. The ground-truth flow shows the level of simplicity and complexity of many triplets with regard to their affinities with instruments. As shown in Table 11, some triplets can be easily inferred by the instrument information such as scissors (likely cutting either *cystic-artery*, *cystic-duct*, *blood-vessels* or some

adhesion), clipper (likely clipping the respective anatomical targets). However, instruments such as *grasper*, *hook*, *bipolar*, etc, perform multiple triplets which are not easily deduced from instrument presence alone. The overcrowded coloring of *grasper* shows that it is used for various actions/targets and often for short sequences/intervals. To better understand the model's capacity to consistently approximate the ground-truth labels, we present their utilized methodologies in the first 8 columns. Even with the complexity of actions performed by the grasper, it is better recognized by top methods. Triplets with bipolar appear very complex to all the methods. A closer attempt is by team ANL-Triplet which leveraged temporal modeling. Triplets with the hook are less difficult, however, their transition is mostly unclear to all the models. Triplets with scissors

Table 11. Qualitative results regarding the quality of triplet recognition on the CholecT50 dataset. The results are obtained on a concatenation of 3 different surgical videos from the testing set. The methodologies employed by each model are indicated in columns 1-7. The triplet flows are categorized by their instrument, and the color shades illustrate their varying interactions (verbs) on different targets.

TEAM	Multi-task Learning	Temporal Modeling	Attention Mechanism	Graph Convolution	Ensemble Methods	+ Phase Labels	+ Spatial Labels	Action Triplets Performed Using Instrument:					
								Grasper	Bipolar	Hook	Scissors	Clipper	Irrigator
Groundtruth													
Trequantista	✓				✓	✓							
2AI (version 2) [†]	✓	✓			✓	✓							
SIAT-CAMI	✓	✓	✓										
HFUT-MedIA	✓			✓									
RDV [‡]	✓	✓											
CITI SJTU	✓	✓											
ANL Triplet	✓	✓											
Digital Surgery	✓			✓									
Casia Robotics	✓	✓											
Lsgroup	✓	✓											
J&M		✓											
Attention Tripnet [‡]	✓	✓											
Ceaiik	✓	✓											
SJTU-IMR	✓	✓	✓										
Tripnet [‡]	✓												
SK	✓	✓											
MMLAB				✓									
Band of Broeders [¶]	✓					✓							
MTL Baseline [‡]	✓												
NCT-TSO	✓												
2Ai	✓	✓			✓	✓							
HFUT-NUS	✓												
CAMP	✓	✓			✓								
Med Recognizer	✓	✓		✓									

Not eligible for award: [†] post-challenge submission, [‡] organizers' baselines, [¶] used non-public third-party dataset.

are the best-approximated predictions due to their limited variability. Specifically, models utilizing attention and/or temporal modeling are the best at this triplet. The use of phase labels mimics methods utilizing temporal modeling. Triplets with the clipper are another class with limited variability and are closely approximated by most of the models. Triplets with the irrigator, although limited in variability, appear very confusing to all the presented methods.

6.7. Limitations

This challenge and analysis present several limitations, which should be addressed in future iterations. Firstly, as with most challenges, given that the participants are not constrained in terms of modeling, comparisons made between different modeling approaches must be treated as indications rather than facts. Participating submissions employed varying degrees of focus on components such as hyperparameter tuning, thus limiting their comparability. We also note that different submissions use a wide range of parameter counts, potentially due to resource constraints or research priorities, which could greatly affect model performance. Another limitation of this challenge was the enforcement of the causality constraint. Participants were asked to ensure that their submission only made use of past frames to make predictions; however, enforcing this at inference time is computationally impractical as each input image must be processed individually along with its entire temporal context to ensure that there is no leakage of information. Alternative strategies could have been to decide on a fixed and limited temporal context (n frames) that would be used to make a prediction or to allow acausal predictions. Bearing in mind that the value of these systems may lie in real-time systems, we opted to allow an unlimited past context and designed a hidden subset of the test set that was used to test causality.

7. Conclusion

As the finest-grained and most comprehensive description of surgical activities for computer vision, surgical action triplets carry significant clinical value. Before this challenge, however, this problem received little attention from the community; in that regard, *CholecTriplet2021* was a success. With a record-high 19 submissions, three of which surpassed the state of the art, this event featured a diverse range of approaches, providing a solid methodological foundation for future research efforts. Most importantly, this challenge showed that surgical action triplet recognition remains an open challenge, with several promising directions to explore. The use of attention, already studied before *CholecTriplet2021*, can be expanded, as shown by several submissions. This challenge also saw the very first uses of graph convolutions and temporal models for surgical action triplets.

Finally, future work should focus on refining the spatial aspect of this problem: locating action triplets in addition to simply detecting their presence would provide much richer information on the surgery. In that sense, full bounding box annotations would be a considerable step forward for research on tool-tissue interaction. Owing to the similarity of tissue manipulation across surgical procedures, action triplet recognition can be adapted to other surgical procedures by following the triplet labeling formalism to annotate the specific procedure data. Transfer learning from existing triplet models or pretraining on the

CholecT50 dataset could potentially benefit model convergence on the new data.

Acknowledgment

The organizers would like to thank the IHU and IRCAD research teams for their help with the initial data annotation during the CONDOR project. We also thank the EndoVis 2021 organizing committee for providing the platform for this challenge. Specifically, we thank Stefanie Speidel, Lena Maier-Hein, and Danail Stoyanov.

Funding

This work was supported by French state funds managed within the Investissements d'Avenir program by BPI France (project CONDOR) and by the ANR under references: Labex CAMI [ANR-11-LABX-0004], DeepSurg [ANR-16-CE33-0009], IHU Strasbourg [ANR-10-IAHU-02], and National AI Chair AI4ORSafety [ANR-20-CHIA-0029-01]. Software validation and evaluation were performed with servers managed by CAMMA, as well as HPC resources from Unistra Mésocentre and GENCI-IDRIS [Grant 2021-AD011011638R1, 2021-AD011011640R1]. Awards for the challenge were sponsored by NVIDIA and Medtronic plc.

Participating teams would like to acknowledge the following funding: **NCT-TSO**: Federal Ministry of Health, Germany (BMG) - part of the SurgOmics project. **CAMP**: Ph.D. Fellowship at CAMP Chair at TU Munich. **CASIA Robotics**: National Key Research and Development Program of China [Grant 2017YFB1302704], the National Natural Science Foundation of China [Grant U1713220], the Beijing Science and Technology Plan [Grant Z191100002019013], and the Youth Innovation Promotion Association of the Chinese Academy of Sciences [Grant 2018165]. **HFUT-MedIA**: National Natural Science Foundation of China [No. 91846107]. **SIAT-CAMI**: Guangdong Key Area Research and Development Program [2020B010165004] and Shenzhen Key Basic Science Program [JCYJ20180507182437217]. **MMLAB**: Shun Hing Institute of Advanced Engineering [SHIAE project #BME-p1-21, 8115064] at the Chinese University of Hong Kong (CUHK) and Singapore Academic Research Fund [Grant R397000353114]. **CITI-SJTU and SJTU-IMR**: Shanghai Municipal Science and Technology Commission [20511105205] and National Natural Science of China [U20A20199]. **2AI**: NORTE-01-0145-FEDER-000045" and "NORTE-01-0145-FEDER-000059", supported by Northern Portugal Regional Operational Programme (NORTE 2020), under the Portugal 2020 Partnership Agreement, through the European Regional Development Fund (FEDER), FCT and FCT/MCTES [project: UIDB/05549/2020, UIDP/05549/2020]. Also, Fundação para a Ciência e a Tecnologia (FCT), Portugal and the European Social Found, European Union, for funding support through the "Programa Operacional Capital Humano" POCH [PhD Grants SFRH/BD/136721/2018, SFRH/BD/136670/2018]. **HFUT-NUS**: the Chinese Scholarship Council [No. 202006690025]. **Ls-group**: Ministry of Science and Technology of the People's Republic of China [2021ZD0201900, 2021ZD0201903]. **Digital Surgery**: Wellcome/EPSRC Centre for Interventional and Surgical Sciences (WEISS) [203145Z/16/Z]; Engineering and Physical Sciences Research Council (EPSRC) [EP/P027938/1, EP/R004080/1, EP/P012841/1]; The Royal Academy of Engineering Chair in Emerging Technologies scheme.

References

- Avci, A., Bosch, S., Marin-Perianu, M., Marin-Perianu, R., Havinga, P., 2010. Activity recognition using inertial sensing for healthcare, wellbeing and

- sports applications: A survey, in: 23th International conference on architecture of computing systems 2010, VDE. pp. 1–10.
- Bawa, V.S., Singh, G., Kaping'a, F., Skarga-Bandurova, I., Oleari, E., Leporini, A., Landolfo, C., Zhao, P., Xiang, X., Luo, G., Wang, K., Li, L., Wang, B., Zhao, S., Li, L., Stabile, A., Setti, F., Muradore, R., Cuzzolin, F., 2021. The SARAS endoscopic surgeon action detection (ESAD) dataset: Challenges and methods. *CoRR abs/2104.03178*. URL: <https://arxiv.org/abs/2104.03178>, [arXiv:2104.03178](https://arxiv.org/abs/2104.03178).
- Carreira, J., Zisserman, A., 2017. Quo vadis, action recognition? A new model and the kinetics dataset, in: 2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21–26, 2017, IEEE Computer Society. pp. 4724–4733. URL: <https://doi.org/10.1109/CVPR.2017.502>, doi:10.1109/CVPR.2017.502.
- Chang, D., Ding, Y., Xie, J., Bhunia, A.K., Li, X., Ma, Z., Wu, M., Guo, J., Song, Y.Z., 2020. The devil is in the channels: Mutual-channel loss for fine-grained image classification. *IEEE Transactions on Image Processing* 29, 4683–4695.
- Chao, Y., Liu, Y., Liu, X., Zeng, H., Deng, J., 2018. Learning to detect human-object interactions, in: 2018 IEEE Winter Conference on Applications of Computer Vision, WACV 2018, Lake Tahoe, NV, USA, March 12–15, 2018, IEEE Computer Society. pp. 381–389. URL: <https://doi.org/10.1109/WACV.2018.00048>, doi:10.1109/WACV.2018.00048.
- Chao, Y., Wang, Z., He, Y., Wang, J., Deng, J., 2015a. HICO: A benchmark for recognizing human-object interactions in images, in: 2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7–13, 2015, IEEE Computer Society. pp. 1017–1025. URL: <https://doi.org/10.1109/ICCV.2015.122>, doi:10.1109/ICCV.2015.122.
- Chao, Y.W., Wang, Z., He, Y., Wang, J., Deng, J., 2015b. Hico: A benchmark for recognizing human-object interactions in images, in: Proceedings of the IEEE International Conference on Computer Vision, pp. 1017–1025.
- Czempiel, T., Paschali, M., Keicher, M., Simson, W., Feussner, H., Kim, S.T., Navab, N., 2020. Tecno: Surgical phase recognition with multi-stage temporal convolutional networks, in: Medical Image Computing and Computer Assisted Intervention - MICCAI 2020 - 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part III, Springer. pp. 343–352. URL: https://doi.org/10.1007/978-3-030-59716-0_33, doi:10.1007/978-3-030-59716-0_33.
- Dergachyova, O., Bouget, D., Huauilmé, A., Morandi, X., Jannin, P., 2016. Automatic data-driven real-time segmentation and recognition of surgical workflow. *Int. J. Comput. Assist. Radiol. Surg.* 11, 1081–1089. URL: <https://doi.org/10.1007/s11548-016-1371-x>, doi:10.1007/s11548-016-1371-x.
- DiPietro, R., Ahmidi, N., Malpani, A., Waldram, M., Lee, G.I., Lee, M.R., Vedula, S.S., Hager, G.D., 2019. Segmenting and classifying activities in robot-assisted surgery with recurrent neural networks. *International journal of computer assisted radiology and surgery* 14, 2005–2020.
- Durand, T., Mordan, T., Thome, N., Cord, M., 2017. Wildcat: Weakly supervised learning of deep convnets for image classification, pointwise localization and segmentation, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 642–651.
- Everingham, M., Eslami, S.M.A., Gool, L.V., Williams, C.K.I., Winn, J.M., Zisserman, A., 2015. The pascal visual object classes challenge: A retrospective. *Int. J. Comput. Vis.* 111, 98–136. URL: <https://doi.org/10.1007/s11263-014-0733-5>, doi:10.1007/s11263-014-0733-5.
- Farha, Y.A., Gall, J., 2019. Ms-tn: Multi-stage temporal convolutional network for action segmentation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 3575–3584.
- Feichtenhofer, C., Fan, H., Malik, J., He, K., 2019. Slowfast networks for video recognition, in: Proceedings of the IEEE/CVF international conference on computer vision, pp. 6202–6211.
- Funke, I., Jenke, A., Mees, S.T., Weitz, J., Speidel, S., Bodenstedt, S., 2018. Temporal coherence-based self-supervised learning for laparoscopic workflow analysis, in: OR 2.0 Context-Aware Operating Theaters, Computer Assisted Robotic Endoscopy, Clinical Image-Based Procedures, - and - Skin Image Analysis - First International Workshop, OR 2.0 2018, 5th International Workshop, CARE 2018, 7th International Workshop, CLIP 2018, Third International Workshop, ISIC 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 16 and 20, 2018, Proceedings, Springer. pp. 85–93. URL: https://doi.org/10.1007/978-3-030-01201-4_11, doi:10.1007/978-3-030-01201-4_11.
- Gkioxari, G., Girshick, R.B., Dollár, P., He, K., 2018. Detecting and recognizing human-object interactions, in: 2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18–22, 2018, Computer Vision Foundation / IEEE Computer Society. pp. 8359–8367. URL: http://openaccess.thecvf.com/content_cvpr_2018/html/Gkioxari_Detecting_and_Recognizing_CVPR_2018_paper.html, doi:10.1109/CVPR.2018.00872.
- Hajj, H.A., Lamard, M., Conze, P., Cochener, B., Quellec, G., 2018. Monitoring tool usage in surgery videos using boosted convolutional and recurrent neural networks. *Medical Image Anal.* 47, 203–218. URL: <https://doi.org/10.1016/j.media.2018.05.001>, doi:10.1016/j.media.2018.05.001.
- Hajj, H.A., Lamard, M., Conze, P., Roychowdhury, S., Hu, X., Marsalkaite, G., Zisimopoulos, O., Dedmari, M.A., Zhao, F., Prellberg, J., Sahu, M., Galdran, A., Araújo, T., Vo, D.M., Panda, C., Dahiya, N., Kondo, S., Bian, Z., Quellec, G., 2019. CATARACTS: challenge on automatic tool annotation for cataract surgery. *Medical Image Anal.* 52, 24–41. URL: <https://doi.org/10.1016/j.media.2018.11.008>, doi:10.1016/j.media.2018.11.008.
- He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 770–778.
- Hochreiter, S., Schmidhuber, J., 1997. Long short-term memory. *Neural computation* 9, 1735–1780.
- Hou, Q., Zhou, D., Feng, J., 2021a. Coordinate attention for efficient mobile network design, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 13713–13722.
- Hou, Q., Zhou, D., Feng, J., 2021b. Coordinate attention for efficient mobile network design, in: IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19–25, 2021, Computer Vision Foundation / IEEE. pp. 13713–13722. URL: https://openaccess.thecvf.com/content/CVPR2021/html/Hou_Coordinate_Attention_for_Efficient_Mobile_Network_Design_CVPR_2021_paper.html.
- Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K.Q., 2017. Densely connected convolutional networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 4700–4708.
- Jin, Y., Long, Y., Chen, C., Zhao, Z., Dou, Q., Heng, P.A., 2021. Temporal memory relation network for workflow recognition from surgical video. *IEEE Transactions on Medical Imaging* 40, 1911–1923.
- Katic, D., Julliard, C., Wekerle, A., Kenngott, H., Müller-Stich, B.P., Dillmann, R., Speidel, S., Jannin, P., Gibaud, B., 2015. Lapontospm: an ontology for laparoscopic surgeries and its application to surgical phase recognition. *Int. J. Comput. Assist. Radiol. Surg.* 10, 1427–1434. URL: <https://doi.org/10.1007/s11548-015-1222-1>, doi:10.1007/s11548-015-1222-1.
- Katic, D., Wekerle, A., Gärtner, F., Kenngott, H., Müller-Stich, B.P., Dillmann, R., Speidel, S., 2014. Knowledge-driven formalization of laparoscopic surgeries for rule-based intraoperative context-aware assistance, in: Information Processing in Computer-Assisted Interventions - 5th International Conference, IPCAI 2014, Fukuoka, Japan, June 28, 2014, Proceedings, Springer. pp. 158–167. URL: https://doi.org/10.1007/978-3-319-07521-1_17, doi:10.1007/978-3-319-07521-1_17.
- Khatibi, T., Dezyani, P., 2020. Proposing novel methods for gynecologic surgical action recognition on laparoscopic videos. *Multim. Tools Appl.* 79, 30111–30133. URL: <https://doi.org/10.1007/s11042-020-09540-y>, doi:10.1007/s11042-020-09540-y.
- Koppala, H.S., Gupta, R., Saxena, A., 2013. Learning human activities and object affordances from RGB-D videos. *Int. J. Robotics Res.* 32, 951–970. URL: <https://doi.org/10.1177/0278364913478446>, doi:10.1177/0278364913478446.
- Kuehne, H., Huang, H., Garrote, E., Poggio, T.A., Serre, T., 2011. HMDB: A large video database for human motion recognition, in: Metaxas, D.N., Quan, L., Sanfeliu, A., Gool, L.V. (Eds.), IEEE International Conference on Computer Vision, ICCV 2011, Barcelona, Spain, November 6–13, 2011, IEEE Computer Society. pp. 2556–2563. URL: <https://doi.org/10.1109/ICCV.2011.6126543>, doi:10.1109/ICCV.2011.6126543.
- Lecuyer, G., Ragot, M., Martin, N., Launay, L., Jannin, P., 2020. Assisted phase and step annotation for surgical videos. *Int. J. Comput. Assist. Radiol. Surg.* 15, 673–680. URL: <https://doi.org/10.1007/s11548-019-02108-8>, doi:10.1007/s11548-019-02108-8.
- Li, Y., Liu, X., Lu, H., Wang, S., Liu, J., Li, J., Lu, C., 2020. Detailed 2d-3d joint representation for human-object interaction, in: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13–19, 2020, Computer Vision Foundation / IEEE. pp. 10163–10172. URL: https://openaccess.thecvf.com/content_CVPR_2020/html/Li_Detailed_2D-3D_Joint_Representation_for_Human-Object_Interaction_CVPR_2020_paper.html, doi:10.1109/CVPR42600.2020.01018.
- Lin, T., Maire, M., Belongie, S.J., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L., 2014. Microsoft COCO: common objects in context, in: Fleet, D.J., Pajdla, T., Schiele, B., Tuytelaars, T. (Eds.), Computer Vision - ECCV 2014 - 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V, Springer. pp. 740–755. URL: https://doi.org/10.1007/978-3-319-10602-1_48, doi:10.1007/

- 978-3-319-10602-1_48.
- Liu, S., Qi, L., Qin, H., Shi, J., Jia, J., 2018. Path aggregation network for instance segmentation, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 8759–8768.
- Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B., 2021. Swin transformer: Hierarchical vision transformer using shifted windows, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 10012–10022.
- Maier-Hein, L., Vedula, S., Speidel, S., Navab, N., Kikinis, R., Park, A., Eisenmann, M., Feussner, H., Forestier, G., Giannarou, S., Hashizume, M., Katić, D., Kennigott, H., Kranzfelder, M., Malpani, A., März, K., Neumuth, T., Padoy, N., Pugh, C., Jannin, P., 2017a. Surgical data science for next-generation interventions. *Nature Biomedical Engineering* 1. doi:10.1038/s41551-017-0132-7.
- Maier-Hein, L., Vedula, S., Speidel, S., Navab, N., Kikinis, R., Park, A., Eisenmann, M., Feussner, H., Forestier, G., Giannarou, S., et al., 2017b. Surgical data science: Enabling next-generation surgery. *Nature Biomedical Engineering* 1, 691–696.
- Mallya, A., Lazebnik, S., 2016. Learning models for actions and person-object interactions with transfer to question answering, in: Leibe, B., Matas, J., Sebe, N., Welling, M. (Eds.), *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part I*, Springer. pp. 414–428. URL: https://doi.org/10.1007/978-3-319-46448-0_25, doi:10.1007/978-3-319-46448-0_25.
- Mishra, K., Sathish, R., Sheet, D., 2017. Learning latent temporal connectionism of deep residual visual abstractions for identifying surgical tools in laparoscopy procedures, in: 2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), IEEE. pp. 2233–2240.
- Mondal, S.S., Sathish, R., Sheet, D., 2019. Multitask learning of temporal connectionism in convolutional networks using a joint distribution loss function to simultaneously identify tools and phase in surgical videos. arXiv preprint arXiv:1905.08315.
- Nwoye, C., 2021. Deep Learning Methods for the Detection and Recognition of Surgical Tools and Activities in Laparoscopic Videos. Ph.D. thesis. URL: <http://icube-publis.unistra.fr/8-Nwoy21>.
- Nwoye, C.I., Gonzalez, C., Yu, T., Mascagni, P., Mutter, D., Marescaux, J., Padoy, N., 2020. Recognition of instrument-tissue interactions in endoscopic videos via action triplets, in: *International Conference on Medical Image Computing and Computer-Assisted Intervention*, Springer. pp. 364–374.
- Nwoye, C.I., Mutter, D., Marescaux, J., Padoy, N., 2019. Weakly supervised convolutional lstm approach for tool tracking in laparoscopic videos. *International journal of computer assisted radiology and surgery* 14, 1059–1067.
- Nwoye, C.I., Padoy, N., 2022. Data splits and metrics for benchmarking methods on surgical action triplet datasets. arXiv preprint.
- Nwoye, C.I., Yu, T., Gonzalez, C., Seeliger, B., Mascagni, P., Mutter, D., Marescaux, J., Padoy, N., 2022. Rendezvous: Attention mechanisms for the recognition of surgical action triplets in endoscopic videos. *Medical Image Analysis* 78, 102433. URL: <https://www.sciencedirect.com/science/article/pii/S1361841522000846>, doi:<https://doi.org/10.1016/j.media.2022.102433>.
- Qi, S., Wang, W., Jia, B., Shen, J., Zhu, S., 2018. Learning human-object interactions by graph parsing neural networks, in: Ferrari, V., Hebert, M., Sminchisescu, C., Weiss, Y. (Eds.), *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part IX*, Springer. pp. 407–423. URL: https://doi.org/10.1007/978-3-030-01240-3_25, doi:10.1007/978-3-030-01240-3_25.
- Ramesh, S., Dall'Alba, D., Gonzalez, C., Yu, T., Mascagni, P., Mutter, D., Marescaux, J., Fiorini, P., Padoy, N., 2021. Multi-task temporal convolutional networks for joint recognition of surgical phases and steps in gastric bypass procedures. *International Journal of Computer Assisted Radiology and Surgery*, 1–9.
- Rupprecht, C., Lea, C., Tombari, F., Navab, N., Hager, G.D., 2016. Sensor substitution for video-based action recognition, in: 2016 IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2016, Daejeon, South Korea, October 9-14, 2016, IEEE. pp. 5230–5237. URL: <https://doi.org/10.1109/IROS.2016.7759769>, doi:10.1109/IROS.2016.7759769.
- Sadhu, A., Gupta, T., Yatskar, M., Nevatia, R., Kembhavi, A., 2021. Visual semantic role labeling for video understanding, in: *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021, Computer Vision Foundation / IEEE*. pp. 5589–5600. URL: https://openaccess.thecvf.com/content/CVPR2021/html/Sadhu_Visual_Semantic_Role_Labeling_for_Video_Understanding_CVPR_2021_paper.html.
- Shi, X., Chen, Z., Wang, H., Yeung, D.Y., Wong, W.K., Woo, W.c., 2015. Convolutional lstm network: A machine learning approach for precipitation nowcasting. *Advances in neural information processing systems* 28.
- Soomro, K., Zamir, A.R., Shah, M., 2012. Ucf101: A dataset of 101 human actions classes from videos in the wild. ArXiv abs/1212.0402.
- Speidel, S., Maier-Hein, L., Stoyanov, D., Bodenstedt, S., Wagner, M., Müller, B., Chen, J., Müller, B., Mathis-Ullrich, F., Scheikl, P., Bernal, J., Histache, A., Fernandes-Esparrach, G., Dray, X., Bano, S., Casella, A., Vasconcelos, F., Moccia, S., Nwoye, C., Alapatt, D., Vardazaryan, A., Padoy, N., Huaultme, A., Harada, K., Jannin, P., Zia, A., Bhattacharyya, K., Liu, X., Wang, Z., Jarc, A., 2021. Endoscopic vision challenge 2021. URL: <https://doi.org/10.5281/zenodo.4572973>, doi:10.5281/zenodo.4572973.
- Stauder, R., Ostler, D., Kranzfelder, M., Koller, S., Feußner, H., Navab, N., 2016. The TUM lapcholo dataset for the M2CAI 2016 workflow challenge. CoRR abs/1610.09278. URL: <http://arxiv.org/abs/1610.09278>, arXiv:1610.09278.
- Tan, M., Le, Q., 2021. Efficientnetv2: Smaller models and faster training, in: *International Conference on Machine Learning*, PMLR. pp. 10096–10106.
- Twinanda, A.P., Shehata, S., Mutter, D., Marescaux, J., De Mathelin, M., Padoy, N., 2016. Endonet: a deep architecture for recognition tasks on laparoscopic videos. *IEEE transactions on medical imaging* 36, 86–97.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I., 2017. Attention is all you need. *Advances in neural information processing systems* 30.
- Vercateren, T., Unberath, M., Padoy, N., Navab, N., 2019. Cai4cai: the rise of contextual artificial intelligence in computer-assisted interventions. *Proceedings of the IEEE* 108, 198–214.
- Wagner, M., Müller-Stich, B.P., Kisilenko, A., Tran, D., Heger, P., Mündermann, L., Lubotsky, D.M., Müller, B., Davitashvili, T., Capek, M., Reinke, A., Yu, T., Vardazaryan, A., Nwoye, C.I., Padoy, N., Liu, X., Lee, E.J., Disch, C., Meine, H., Xia, T., Jia, F., Kondo, S., Reiter, W., Jin, Y., Long, Y., Jiang, M., Dou, Q., Heng, P.A., Twick, I., Kirtac, K., Hosgor, E., Bolmgren, J.L., Stenzel, M., von Siemens, B., Kennigott, H.G., Nickel, F., von Frankenberg, M., Mathis-Ullrich, F., Maier-Hein, L., Speidel, S., Bodenstedt, S., 2021. Comparative validation of machine learning algorithms for surgical workflow and skill analysis with the heichole benchmark. arXiv:2109.14956.
- Wang, C.Y., Liao, H.Y.M., Wu, Y.H., Chen, P.Y., Hsieh, J.W., Yeh, I.H., 2020a. Cspnet: A new backbone that can enhance learning capability of cnn, in: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pp. 390–391.
- Wang, J., Sun, K., Cheng, T., Jiang, B., Deng, C., Zhao, Y., Liu, D., Mu, Y., Tan, M., Wang, X., et al., 2020b. Deep high-resolution representation learning for visual recognition. *IEEE transactions on pattern analysis and machine intelligence* 43, 3349–3364.
- Wang, Y., He, D., Li, F., Long, X., Zhou, Z., Ma, J., Wen, S., 2020c. Multi-label classification with label graph superimposing, in: *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 12265–12272.
- Xu, M., Islam, M., Lim, C.M., Ren, H., 2021. Learning domain adaptation with model calibration for surgical report generation in robotic surgery, in: *IEEE International Conference on Robotics and Automation, ICRA 2021, Xi'an, China, May 30 - June 5, 2021, IEEE*. pp. 12350–12356. URL: <https://doi.org/10.1109/ICRA48506.2021.9561569>, doi:10.1109/ICRA48506.2021.9561569.
- Yu, T., Mutter, D., Marescaux, J., Padoy, N., 2019. Learning from a tiny dataset of manual annotations: a teacher/student approach for surgical phase recognition. *IPCAI*.
- Zhang, H., Kyaw, Z., Chang, S.F., Chua, T.S., 2017. Visual translation embedding network for visual relation detection, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5532–5540.
- Zhuang, B., Wu, Q., Shen, C., Reid, I., van den Hengel, A., 2018. Hcprd: a benchmark for large-scale human-centered visual relationship detection, in: *Thirty-Second AAAI Conference on Artificial Intelligence*.
- Zisimopoulos, O., Flouty, E., Luengo, I., Giataganas, P., Nehme, J., Chow, A., Stoyanov, D., 2018. Deepphase: Surgical phase recognition in CATARACTS videos, in: Frangi, A.F., Schnabel, J.A., Davatzikos, C., Alberola-López, C., Fichtinger, G. (Eds.), *Medical Image Computing and Computer Assisted Intervention - MICCAI 2018 - 21st International Conference, Granada, Spain, September 16-20, 2018, Proceedings, Part IV*, Springer. pp. 265–272. URL: https://doi.org/10.1007/978-3-030-00937-3_31, doi:10.1007/978-3-030-00937-3_31.

CRedit authorship contribution statement

C.I. Nwoye : Conceptualization, Data Curation, Data Analysis and Interpretation, Methodology, Software, Investigation, Validation, Evaluation, Formal Analysis, Visualization, Writing - Original Draft, Writing - Review & Editing, Challenge Organization, Resources.

D. Alapatt : Conceptualization, Investigation, Validation, Evaluation, Formal Analysis, Visualization, Writing - Original Draft, Writing - Review & Editing, Challenge Organization, Resources.

T. Yu : Data Curation, Investigation, Formal Analysis, Writing - Original Draft, Writing - Review & Editing, Visualization.

A. Vardazaryan : Conceptualization, Investigation, Validation, Writing - Review & Editing, Challenge Organization, Resources.

F. Xia, Z. Zhao, T. Xia, F. Jia, Y. Yang, H. Wang, D. Yu, G. Zheng, X. Duan, N. Getty, R. Sanchez-Matilla, M. Robu, L. Zhang, H. Chen, J. Wang, B. Zhang, B. Gerats, S. Raviteja, R. Sathish, R. Tao, S. Kondo, W. Pang, H. Ren, J.R. Abbing, M.H. Sarhan, S. Bodenstedt, N. Bhasker, B. Oliveira, H. Torres, L. Ling, F. Gaida, T. Czempliel, Y. Jiang, Y. Long, J. Vilaça, P. Morais, J. Fonseca, R.M. Egging, I.N. Wijma, C. Qian, G. Bian, Z. Li, V. Balasubramanian, D. Sheet, I. Luengo, Y. Zhu, S. Ding, J. Aschenbrenner, N.E van der Kar, M. Xu, M. Islam, L. Seenivasan, A.C. Jenke, D. Stoyanov : Methodology, Software, Writing - Review & Editing.

C. Gonzalez, B. Seeliger, P. Mascagni : Data Curation, Writing - Review & Editing.

D. Mutter : Data Curation, Writing - Review & Editing, Supervision.

N. Padoy : Conceptualization, Writing - Review & Editing, Supervision, Challenge Organization, Resources, Funding Acquisition, Project Administration.