# Multitask Learning in Minimally Invasive Surgical Vision: A Review

Oluwatosin Alabi[a], Tom Vercauteren[a], Miaojing Shi[b,*]

[a]*School of Biomedical Engineering & Imaging Sciences, King's College London*
[b]*College of Electronic and Information Engineering, Tongji University*

## ABSTRACT

Minimally invasive surgery (MIS) has revolutionized many procedures and led to reduced recovery time and risk of patient injury. However, MIS poses additional complexity and burden on surgical teams. Data-driven surgical vision algorithms are thought to be key building blocks in the development of future MIS systems with improved autonomy. Recent advancements in machine learning and computer vision have led to successful applications in analyzing videos obtained from MIS with the promise of alleviating challenges in MIS videos.

Surgical scene and action understanding encompasses multiple related tasks that, when solved individually, can be memory-intensive, inefficient, and fail to capture task relationships. Multitask learning (MTL), a learning paradigm that leverages information from multiple related tasks to improve performance and aid generalization, is well-suited for fine-grained and high-level understanding of MIS data.

This review provides an overview of the current state-of-the-art MTL systems that leverage videos obtained from MIS. Beyond listing published approaches, we discuss the benefits and limitations of these MTL systems. Moreover, this manuscript presents an analysis of the literature for various application fields of MTL in MIS, including those with large models, highlighting notable trends, new directions of research, and developments.

## 1. Introduction

Minimally invasive surgeries (MIS) have become increasingly popular due to their benefits, such as reduced blood loss, less pain, faster recovery times, and fewer post-surgical complications (Jaffray, 2005). However, MIS utilizes indirect vision with a limited field of view from the endoscope, making it challenging for surgeons to interpret events on the surgical scene accurately (Islam et al., 2019a). To overcome this, computer-aided intervention techniques that augment the information obtained through minimally invasive cameras have been proposed (Vercauteren et al., 2020; Ward et al., 2021).

Deep learning has demonstrated remarkable success in providing solutions to computer vision tasks for MIS, such as instrument classification (Wang et al., 2017; Mishra et al., 2017; Al Hajj et al., 2018), instrument segmentation (Garcia-Peraza-Herrera et al., 2017; Milletari et al., 2018; Pakhomov and Navab, 2020; Islam et al., 2019b), and surgical scene depth estimation (Luo et al., 2019; Xiao et al., 2020; Wei et al., 2022; Luo et al., 2022). Conventionally, separate models would be trained for these related tasks, which can be computationally impractical and inefficient. Therefore, there is a need for deep learning approaches that leverage information from different tasks to improve both performance and efficiency.

Multitask learning (MTL) is a machine learning approach

---
*Corresponding author.
*e-mail:* mshi@tongji.edu.cn (Miaojing Shi)

that seeks to improve generalisation performance by leveraging domain-specific information from multiple related tasks, with shared parameters reducing overfitting, memory footprint and improving regularisation (Caruana, 1997). MTL has been successfully applied in various fields such as natural language processing (Collobert and Weston, 2008), computer vision (Girshick, 2015), and robotics (Kalashnikov et al., 2021). In the context of MIS, MTL has primarily been utilized to enhance scene understanding. By simultaneously solving multiple tasks and leveraging knowledge from all of them, MTL offers an efficient approach for solving multiple tasks in MIS. To ensure the widespread adoption of computer-aided intervention solutions for MIS, it is crucial to develop systems that can comprehensively understand the surgical scene, rather than addressing individual tasks separately.

In this survey, we aim to review literature that utilizes MTL to solve multiple tasks in MIS on a case-by-case basis. We identify prevailing trends and offer valuable insights into these trends. We restrict our scope to MTL in MIS. For a broader survey of MTL in deep learning and deep learning in MIS, we refer readers to (Vandenhende et al., 2022; Zhang and Yang, 2022) and (Rivas-Blanco et al., 2021) respectively. Additionally, we restrict the scope of this manuscript to methods that utilize videos and/or images obtained from the MIS camera to solve multiple tasks where each task provides meaningful and relevant outputs. We also include papers that predict a primary task with a meaningful and relevant output, along with other auxiliary outputs that are meant to guide the learning of the primary task. Lastly, with the advent of large models, we summarize previous works that utilize large models with capabilities to solve multiple tasks in MIS.

We exclude literature on open surgeries and external operating room cameras, as our focus is primarily on techniques that utilize endoscopic videos in surgical scenarios, along with the challenges that come with this scenario and the solutions offered to these challenges.

The remainder of this survey paper is divided into four sections. Section 2 provides a quick introduction to popular deep MTL techniques commonly employed in MTL research for vision, emphasizing their key features and concepts. Section 3 reviews how MTL has been applied to solve multiple tasks for surgical scene understanding, examining each work extensively and identifying the current trends in various research areas. In Section 4, an overview of publicly available datasets that can be used to advance MTL research in MIS is presented. Section 5 presents our insights on the current state of MTL research in MIS, discusses potential directions for future work, and concludes this survey by summarising the main contributions of this paper.

## 2. Deep MTL for vision

This section presents a short introduction to widely used deep MTL techniques in the field of computer vision research. We present a short examination of these techniques and representative papers that showcase the utilization of each technique. By presenting the foundations and applications of MTL techniques,
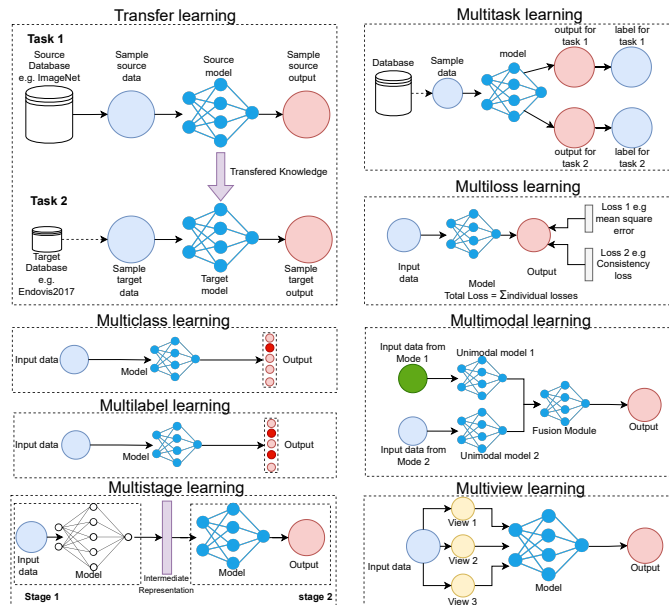


**Fig. 1:** Diagrams illustrating the key concepts and core components for transfer learning, multitask learning, multiloss learning, multiclass learning, multilabel learning, multistage learning, multimodal learning, and multiview learning.

this section aims to equip readers with a basic understanding of the subject, setting the stage for the upcoming sections that analyze methods employed for solving multiple vision tasks in the context of MIS.

### 2.1. MTL and other learning paradigms

MTL is a learning paradigm that trains a single model to perform multiple tasks simultaneously by sharing information across them, leading to better representations useful for all tasks.

The idea of using information from a different task or different outputs to improve representation is not unique to MTL and shares similarities with other learning paradigms, which can make it challenging to clearly differentiate between them. Related learning paradigms include transfer learning, multiview learning, multiloss learning, multilabel learning, multiclass learning, and multistage learning, multimodal learning. A diagrammatic representation of the paradigms discussed in this section can be seen in Fig. 1. Despite their similarities, each paradigm has its unique features and objectives, and understanding their differences and similarities is essential for their proper application.

*Transfer learning* enhances target task performance using knowledge from a related source task, while MTL trains multiple tasks concurrently. *Multi-loss learning* employs multiple loss functions for neural network training. This approach finds use in both MTL and single-task learning. *Multilabel learning* refers to problems where single data instances can have multiple class labels. *Multiclass learning* classifies input data into one of several classes, whereas MTL concurrently optimizes multiple related tasks, sharing the same input data. *Multiview learning* optimizes multiple data representations (views) with

the same output. *Multistage learning* processes inputs through multiple stages. Multitask networks can have single or multiple stages based on the design. *Multimodal learning* involves learning over diverse data modalities. Multimodal learning is usually done with a separate encoder for each different modality and a fusion module. Multitask learning can have inputs of different modalities.

### 2.2. MTL concepts

There are different ways of categorizing MTL core concepts. For our purpose, we focus on four concepts that recurrently emerge in the MTL literature for MIS: parameter sharing, optimization and task-balancing, auxiliary objectives, and data-efficient approaches.

#### 2.2.1. Parameter sharing and feature representation

The mainstream multitask paradigm assumes that tasks are related as the knowledge required for solving these tasks is connected. Following this assumption, features produced for the same input sample on multiple tasks which are related should also be related. Hence, common feature representation learnt for multiple tasks would have more informative and robust feature representations as they take into account multiple tasks (Zhang and Yang, 2022). Additionally, each task acts as a regularization for other tasks, smoothening out noise, as features which are utilized for multiple different tasks are less likely to overfit to training samples (Zhang and Yang, 2022).

Historically, the methods for sharing feature representations in deep neural networks are classified into hard and soft parameter sharing (Vandenhende et al., 2022). A diagrammatic representation of soft and hard parameter sharing can be seen in Fig. 2. Hard parameter sharing refers to architectures where tasks are to be jointly learnt utilizing the same weights and biases for some layers. These layers are aptly called *shared layers*. The other layers which are not shared are called *task-specific* layers. Hard parameter sharing is commonly interpreted using an encoder-decoder architecture where the encoder is shared, and the decoder is task-specific. Some network architectures for hard parameter sharing also include methods to facilitate information transfer between different decoders (Xu et al., 2018; Liu et al., 2019; Zhang et al., 2018). If $x_i$ represents input samples, $h$ denotes hidden layers, and $y_{i,t}$ represent outputs for the $i^{th}$ sample on the $t$ task, $f_{sh}(.)$ represents shared layers and $f_{task}(.)$ represents the task-specific layers, then hard parameter sharing can be summarized as

$$h_i = f_{sh}(x_i) \quad y_{i,t} = f_{task}(h_i) \tag{1}$$

On the other hand, soft parameter sharing does not directly share layers for each task. Instead, a separate model is used for each task. Soft parameter sharing instead shares parameters by adjusting the weights and biases of different task models based on information from other task models, leading to model weights and biases, which are functions of representations from different tasks. If $M_t(.)$ is a model for task $t$ and $M_t$ produces features $f_t$, then soft parameter sharing can be written as
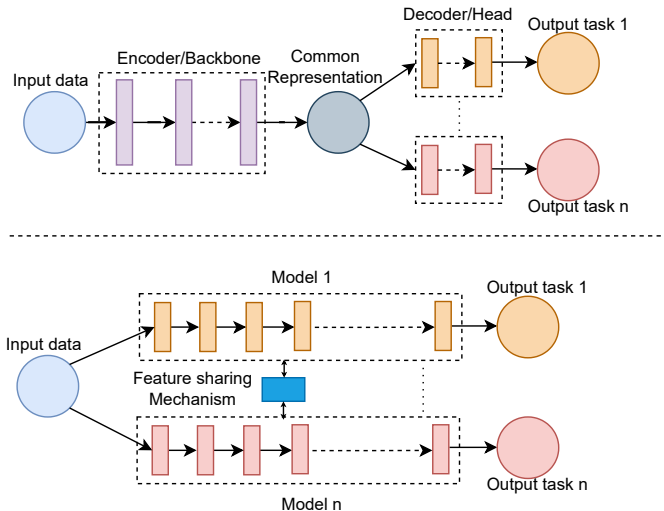
$$M_t = R(f_1, f_2...f_t) \tag{2}$$



**Fig. 2:** Top: hard parameter sharing in deep neural networks for multitask learning, featuring a shared encoder/backbone with a common representation and separate decoders or heads. Right: soft parameter sharing with separate models per task and specialized feature-sharing mechanism.

where $R$ is a feature-sharing mechanism which determines how features are shared.

A prototypical example that uses hard parameter sharing is UberNet (Kokkinos, 2017). UberNet is proposed as a general-purpose computer vision model for multiple tasks. The network architecture is a shared multi-resolution VGG (Simonyan and Zisserman, 2015) with a branch network from each layer of the VGG for each task. The output of each task branch is fused with branches for the same branch from other layers and finally fused for all resolutions to produce predictions for each task. Other examples of hard parameter-sharing architectures include (Xu et al., 2018; Teichmann et al., 2018; Heuer et al., 2021; Liu et al., 2019; Zhang et al., 2018),

Cross-stitch network (Misra et al., 2016) is a good example of soft parameter sharing. This paper is concerned with the design decision about when to split a network into task-specific networks and general/shared parameters. To solve this, the authors propose a unified splitting architecture, which is a type of soft parameter sharing architecture. Different separate networks are connected with the cross-stitch module across various layers of the separate networks. This is done by learning a linear combination of the input activation maps from both tasks. Other examples of soft parameter sharing architectures include (Gao et al., 2019; Ruder et al., 2019)

#### 2.2.2. Optimization and task balancing

Designing networks with shared parameters and learning these parameters together raises an important question: What is the right way to optimise both tasks jointly?

Fig. 3 provides a brief overview of the various methods of optimization methods discussed in this section. LibMTL(Lin and Zhang, 2023) is a python repository with adaptable code implementation for the common multitask architectures and optimization strategies discussed here.

For hard parameter sharing with shared encoders and task-specific decoders, as shown in Fig. 2, we can write the baseline equation for the multitask loss while omitting batching considerations for simplicity, as Equation (3). The gradient descent equation for the shared parameters can be written as Equation (4), and the gradient descent equation for the task-specific parameters as Equation (5):

$$L_{total} = \sum_{i=0}^{N} L_i (\theta_t, \theta_s, X, Y) \tag{3}$$

$$\theta_s = \theta_s - \alpha \sum_{i=0}^{N \cdot} \frac{\partial L_i}{\partial \theta_s} (\theta_t, \theta_s, X, Y) \tag{4}$$

$$\theta_{t(a)} = \theta_{t(a)} - \alpha \frac{\partial L_{i(a)}}{\partial \theta_{t(a)}} (\theta_t, X, Y) \tag{5}$$

where $N$ is the number of tasks, $X$ is the input, $Y$ is the ground-truth, $\alpha$ is the learning rate, $L_i$ is the loss for the $i$th task, $a$ refers to the current task been considered, $\theta_s$ denotes the shared parameters while $\theta_t$ the task specific parameters.

From Equation (4), we can deduce that the network weights in shared layers are affected by multiple supervision signals. Hence, the best way of balancing different task signals in the shared layers has been explored in several works (Sener and Koltun, 2018; Cipolla et al., 2018; Guo et al., 2018).

Equation (3) shows the multitask loss for a network called linear scalarization - when the loss functions for $N$ tasks are simply added together. Linear scalarization is the most utilised task-balancing method (Hu et al., 2023). It involves assigning a scalar weight to each task and optimising the scaled addition as can be seen in (3), (7), and (8)

$$L_{total} = \sum_{i=0}^{N} w_i L_i (\theta_t, \theta_s, X, Y) \tag{6}$$

$$\theta_s = \theta_s - \alpha \sum_{i=0}^{N \cdot} w_i \frac{\partial L_i}{\partial \theta_s} (\theta_t, \theta_s, X, Y) \tag{7}$$

$$\theta_{t(a)} = \theta_{t(a)} - \alpha w_{i(a)} \frac{\partial L_{i(a)}}{\partial \theta_{t(a)}} (\theta_t, \theta_s, X, Y) \tag{8}$$

where $w_i$ is the weight assigned to a particular task, and all other notations remain the same. Despite its simplicity, linear scalarization works surprisingly well, especially when combined with techniques like grid search (Xin et al., 2022).

Another example of a task balancing method is the automatic dynamic tuning of the weight assigned to each task at different points during training. Cipolla et al. (2018) propose a method to automatically weight losses during multitask optimization using the uncertainty in the predictions of each task. The authors reformulate the linear scalarization loss function to include a homoscedastic loss uncertainty term for each task, an uncertainty term independent of the inputs but depends only on the task problem. This method ascribes an uncertainty scalar for each task, and depending on the value of this uncertainty scalar, the gradient magnitudes of each task are weighted. Instead of using uncertainty to automatically determine loss weights,

Chen et al. (2018) utilize each task's gradient values to estimate the rate of convergence of each task and corrects this rate of convergence to be equal by adjusting the weights associated with each task in linear scalarization formulation. Guo et al. (2018) introduce the notion of dynamic task prioritization to prioritize more difficult tasks. Guo et al. (2018) observe that imbalance in task difficulty can result in prioritization of the easy tasks during optimization and propose to measure task difficulty as inversely proportional to the task performance. Task performance is measured by key performance metrics (KPIs) like accuracy for classification and IoU for segmentation. Dynamic task prioritization automatically prioritizes more difficult tasks by adaptively adjusting the weight of each task's loss objective based on task KPIs.

A different task-balancing idea is the direct adjustment of the gradients for each task calculated during back-propagation instead of tuning weights to give priority to different tasks. These methods claim to be able to tackle problems with multitask optimization, such as negative interference. Yu et al. (2020a) propose the PCGrad method, which uses cosine similarity to check if the directions of gradients for tasks trained for multitask learning are conflicting, which can lead to negative interference. If the gradients are conflicting, the PCGrad method projects the task gradients of a task to the normal plane of the other task to remove the conflict and optimizes with this projection. Further investigation into conflicting gradients by Wang et al. (2021) reveals that the occurrence of conflicting gradients during training as measured by cosine similarity is very sparse. Hence, Wang et al. (2021) design a more proactive method for ensuring that there are no conflicting gradients called "gradient vaccine". Wang et al. (2021) note that the adjusted gradients using PCGrad project gradients to normals, making a cosine similarity of 0 the target in PCGrad. Gradient vaccine sets a targeted cosine similarity, which is greater than 0. By selecting a targeted cosine similarity, gradient vaccine ensures that both conflicting and non-conflicting gradients that are still very dissimilar can be projected to the target cosine similarity plane, ensuring frequent gradient update.

Another key task-balancing idea is the interpretation of multitask optimization as multi-objective optimization, which requires a Pareto optimal solution (Sener and Koltun, 2018; Lin et al., 2019). A Pareto optimal solution is reached when one objective function cannot be improved without sacrificing another objective. Pareto optimal solutions lie on a Pareto front, which is a set of optimal solutions in the space of objective functions in multi-objective optimization problems. MTMO (Sener and Koltun, 2018) and PTML (Lin et al., 2019) are examples of papers that try to figure out the best way to optimize multiple tasks together with the Pareto principle. Unlike linear scalarization, which combines multiple loss functions via a system of fixed weights and optimizes the multiple losses as a single loss function, Pareto optimization is a multi-objective optimization problem. Sener and Koltun (2018) establishes the foundation for a Pareto optimal solution in MTL by mathematically deriving an upper bound for the multi-objective optimization problem. This upper bound is then optimized with a multi-gradient descent algorithm. PTML pushes this idea further by gener-
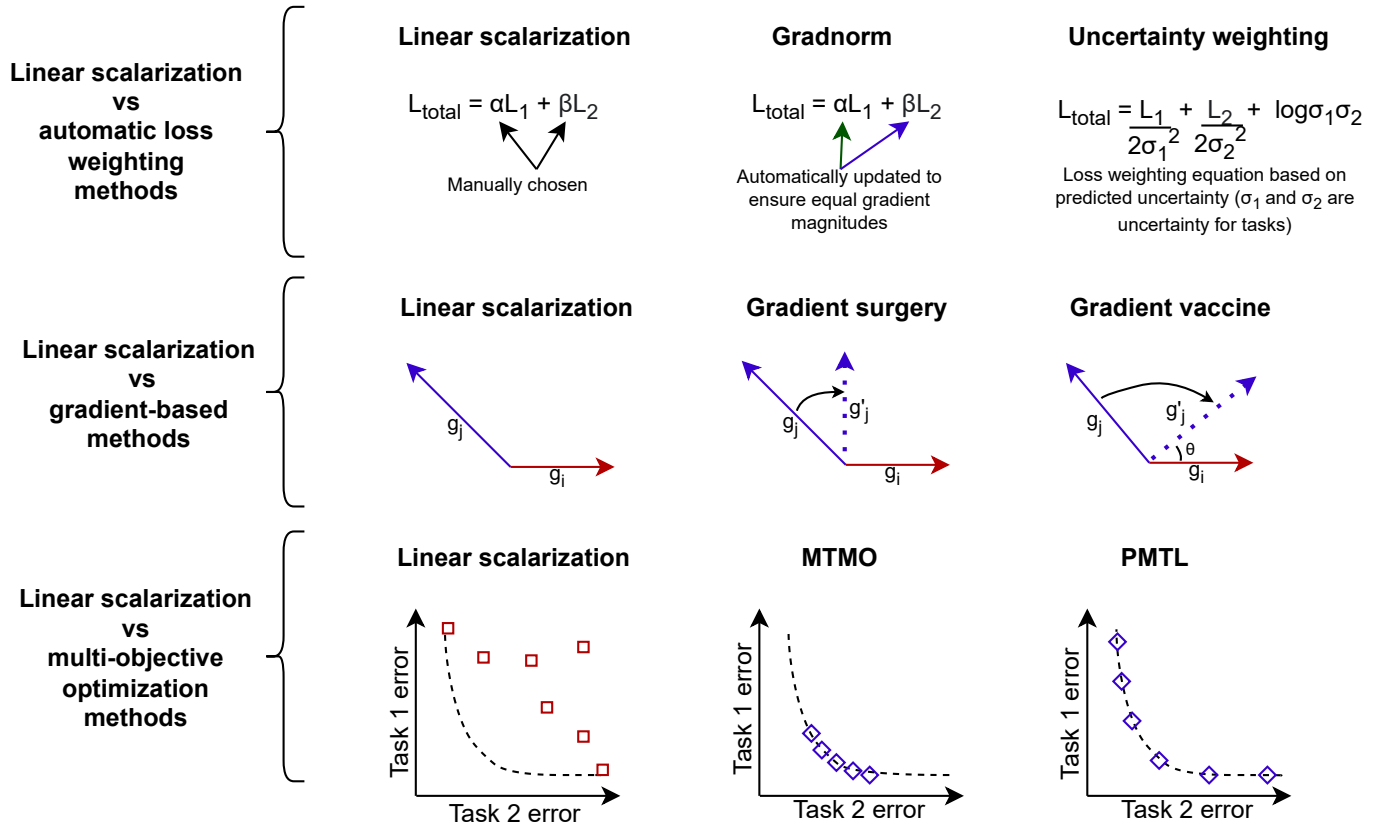
**Fig. 3:** An overview of optimization techniques for multitask learning as discussed in Section 2.2.2. The classical method, linear scalarization, involves manually weighting the loss functions of all tasks. The first row illustrates linear scalarization against automatic loss weighting methods that dynamically adjust weights during training, such as updating weights to ensure gradient magnitude consistency (Chen et al., 2018) and defining a loss weighting equation based on predicted uncertainty of each task (Cipolla et al., 2018). The second row illustrates gradient-based approaches in comparison to linear scalarization. Gradient-based methods directly modify gradients to mitigate negative transfer, achieved by projecting conflicting gradients to the normal plane of the gradient of another task (Yu et al., 2020a) or ensuring gradients are at a target angle to each other (Wang et al., 2021). The third row illustrates linear scalarization and multi-objective optimization techniques. MTMO (Sener and Koltun, 2018) ensures that solutions are on the Pareto front, while PTML (Lin et al., 2019) enables the selection of Pareto front solutions, favouring specific tasks.

alizing the mathematical formulations from Sener and Koltun (2018). It does so by relaxing the Pareto optimal condition and allows for solutions that are as Pareto optimal as possible while still favouring certain tasks.

### 2.2.3. Auxiliary objectives

In MTL, auxiliary tasks are additional tasks that are learned simultaneously with the primary task to improve the overall performance of a model on the primary task. The idea behind auxiliary tasks is that they can provide helpful information or constraints to help the model learn a rich and robust representation of the input data, which in turn benefits the primary task (Liebel and Körner, 2018). Zhang et al. (2014) proffer a method that optimizes a main task in their application, facial landmark detection, together with auxiliary tasks like head pose estimation, facial expression recognition and age estimation. The authors propose a deep network to extract shared features with different classifier heads for each task.

Moreover, auxiliary tasks can also be used to extract features or information that are needed for the main task. This is popular for primary tasks that are complex and can be logically broken into various sub-tasks, with information about each of the aux-

iliary subtasks task being useful for solving the more complex primary task. The MaskFormer (Cheng et al., 2021) is an example of network design that follows this principle in which instance or semantic segmentation is the primary task, but the primary task is broken down into two stages. The first stage has two auxiliary tasks for binary mask prediction and corresponding classification for each mask. The second stage uses masks and classification features generated from the earlier stage to produce the final segmentation output.

### 2.2.4. Data efficient approaches

Building large-scale computer vision datasets is resource intensive (Liao et al., 2021). This is even more noticeable for networks that use the MTL paradigm, which requires labelling for multiple different tasks per sample. There are few publically available datasets designed specifically for working with multiple tasks (Zamir et al., 2018; Yu et al., 2020b). Hence, problems such as limited annotations for multiple tasks and the availability of different task labels in different datasets, are often faced by researchers face when they want to utilize MTL. Consequently, there have been research works to develop MTL techniques to solve these challenges.

A method by which multitask learning can be used to tackle issues with limited annotations for multiple tasks is with better representation learning via multitask learning. Representation learning involves learning the representation of the data that makes it easier to extract useful information for downstream tasks that require a few labelled data (Bengio et al., 2013). Learning representations that are task-agnostic over a range of tasks would be preferred to learning single-task representations as these representations would be more robust to noise (Nguyen et al., 2022). Self-supervised MTL methods for representation learning utilize self-supervised tasks which do not require new annotations and learn representations for multiple tasks together. These representations can then be utilized for downstream tasks with fewer amounts of data. An example of such work is from Doersch and Zisserman (2017) who utilize four self-supervised vision tasks (relative position, colourization, the exemplar task, and motion segmentation) for pretraining instead of a single pretraining task. Doersch and Zisserman (2017) demonstrate that the representations learnt using their method are competitive to ImageNet pretrained representations without the need for labelling. Ghiasi et al. (2021) presents Multi-Task Self-Training (MuST), which is a method to harness the knowledge from specialized teacher models for different tasks to create a large multitask dataset with pseudo-labels generated from these specialized teachers. The large pseudo-label dataset is then used to train a multi-task network. The authors show that the representation generated in this multitask network generalizes well to downstream tasks.

There are also examples in the literature for combining information in different datasets to perform MTL. Li et al. (2022d) proffer a method for combining multiple datasets collected on similar distributions but annotated for different tasks. The authors propose a method that leverages task relations between task pairs to supervise MTL task pairs jointly. They map a task pair to a supervised joint space to enable information sharing between two tasks. Then, they propose a supervised learning loss for the task with known labels and a consistency loss in the joint task space to train tasks with unknown labels. Dorent et al. (2021) tackles the problem of learning from domain-shifted datasets, each with single task-specific annotations. Specifically, for the brain tissue segmentation task, there are datasets with segmentation of brain structures in healthy brains, and there are datasets which segment pathologies like brain lesions and tumours, but there are no datasets with both types of labels. The authors derive an upper bound of the loss for the joint probability problem.

## 3. MTL in MIS

This section delves into the applications of the multitask paradigm in the context of surgical scene understanding, where it tackles the challenge of simultaneously learning multiple subtasks to enhance both efficiency and accuracy. Additionally, it explores the utility of auxiliary tasks to augment primary tasks in surgical scene understanding. In the subsequent subsections, we describe the deployment of MTL for surgical scene understanding across six distinct categories: MTL for perceptual tasks, MTL for tracking and control, MTL for surgical workflow analysis, MTL for surgical skill assessment, MTL for report generation, and large models for solving multiple tasks.

### 3.1. MTL for perceptual tasks

In the context of this report, *perceptual tasks* refer to tasks such as image segmentation, object detection, motion estimation, and depth estimation. These tasks are focused on extracting essential visual information from images or videos, with an emphasis on revealing the spatial layout, motion, semantic, and depth relationships of objects within the scene. Perceptual tasks provide valuable information about spatial layout, motion, and depth relationships in the scene. The outcomes of perceptual tasks play a critical role in providing insights into the surgical scene, serving as the foundational elements for more advanced computer vision processes, including action recognition and object tracking. By enhancing the performance of perceptual tasks, researchers can better prepare input data in a way that significantly eases and enhances the accuracy of subsequent higher-level tasks.

MTL for perceptual tasks in minimally invasive scenes seeks to identify task combinations that can improve problem-solving accuracy (Huang et al., 2022a; Islam et al., 2020a; Sanchez-Matilla et al., 2021; Qin et al., 2020a; Bhattarai et al., 2023; Wang et al., 2022a) or efficiency, taking into account factors such as inference speed (Sanchez-Matilla et al., 2021; Islam et al., 2020a), memory usage (Sanchez-Matilla et al., 2021; Islam et al., 2020a), and annotation requirements (Sanchez-Matilla et al., 2021). Furthermore, the existing body of literature demonstrates the presence of various distinct auxiliary tasks that can provide valuable guidance for optimizing perceptual tasks (Qin et al., 2020a; Bhattarai et al., 2023; Wang et al., 2022a).

In their work, Huang et al. (2022a) present an approach that concurrently tackles the depth estimation and binary tool segmentation tasks in laparoscopic images. The methodology employs a U-Net-like architecture with a shared encoder and two separate decoders, one dedicated to each task. Laparoscopic stereo images serve as the model's input. In this framework, the segmentation decoder follows a conventional U-Net decoder style (Ronneberger et al., 2015). On the other hand, the depth estimation decoder generates disparity maps using the left-right consistency unsupervised depth estimation method, as outlined in Godard et al. (2017). Notably, the unsupervised depth estimation approach is advantageous for endoscopic surgical datasets, which often lack ground-truth depth labels. The results of this study indicate significant improvements in both instrument binary tool segmentation and depth estimation tasks.

Sanchez-Matilla et al. (2021) introduce an efficient and memory-friendly approach to MTL, which simultaneously addresses classification, instrument-type segmentation, and bounding box detection. The authors recognize the challenges associated with obtaining segmentation labels and, as a result, devise a solution that supports weak supervision. This method harnesses the power of an EfficientNet backbone (Koonce and Koonce, 2021). Unlike Huang et al. (2022a), which advocates for distinct decoders for each task, the proposed model employs

a single backbone with straightforward task heads for each objective. During training, the emphasis is placed on leveraging bounding box and classification labels, with only limited instrument segmentation labels. In cases where segmentation labels are unavailable, class activation maps are employed for guidance and supervision. The outcomes reported by the authors are notable, demonstrating that even with just one per cent of the segmentation labels, they achieve commendable results in the surgical tool segmentation task.

In their work illustrated in Fig. 4, Islam et al. (2020a) introduce a real-time network designed for solving instrument detection and instrument-type segmentation in robotic surgery, referred to as the Attention pruned MTL network (AP-MTL). AP-MTL adopts a U-Net-like architecture with shared encoders and task-specific decoders. The decoder for object detection follows the design of the multi-box single shot detector (SSD) (Liu et al., 2016), while the segmentation network is a variant of the standard U-Net segmentation decoder. It incorporates custom squeeze excitation modules (Hu et al., 2018) known as the spatial channel squeeze excitation module (scSE). Additionally, Islam et al. (2020a) presents a custom optimization method known as Asynchronous Task-aware Optimization (ATO). ATO optimizes different parts of the proposed network separately to ensure that correlated tasks converge at the same point, even if they do so at varying speeds. Following the optimization step, ATO introduces regularization to promote a more generalized network and ensures the smooth flow of gradients. To implement the ATO algorithm, gradients for each task are calculated and optimized separately, simulating an *attach, optimize, and detach* mechanism for the task decoders.

In their work, Baby et al. (2023) introduce a modification to the standard Mask-RCNN framework (He et al., 2017) for addressing instance segmentation in surgical instruments. The authors identify an issue where the classification of predicted masks often produces inaccurate results, despite the bounding box detection and mask segmentation tasks generally being performed correctly. To address this challenge, the authors propose the incorporation of a dedicated classification module to decouple the classification task from the region proposal and mask prediction processes. This new module takes as input multiscale features extracted from the feature extractor and the predicted instance masks. These predicted instance masks are employed for multiscale masking, and the results at different scales are subsequently merged and passed to a classification head.

Zhao et al. (2022) propose TraSeTR, a method for leveraging tracking cues to enhance surgical instrument segmentation. This approach employs a transformer-based architecture that bears similarities to the DETR architecture (Carion et al., 2020) for predicting instrument class, bounding box, and binary segmentation classes. Zhao et al. (2022) improve on the standard DETR architecture by utilizing queries from prior frames for the instrument detection in the current frame. These queries are encoded with previous instrument information and serve as a form of tracking signal. These queries are applied to the transformer decoder, and identity matching between the previous queries and current queries are used as tracking cues. In addition, the authors also apply a contrastive query learning strat-
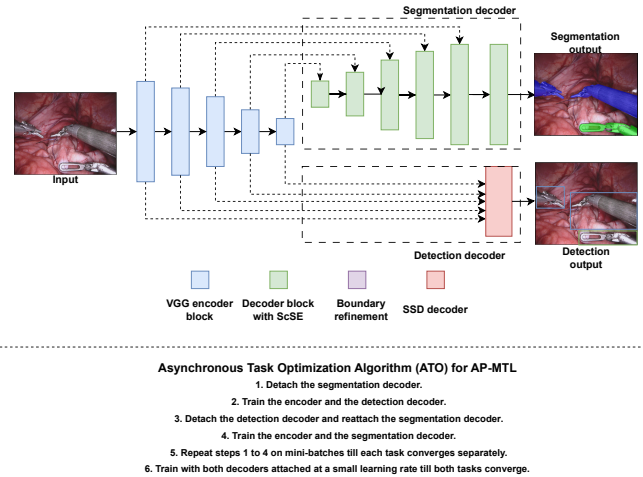


**Fig. 4:** Illustration of the Attention Pruned Multitask Learning (AP-MTL) Network and the optimization method used for training this network (Islam et al., 2020a). The top image shows an encoder-decoder network with skip connections for its segmentation and detection decoders. A summary of the Asynchronous Task Optimization (ATO) for obtaining convergence for both tasks in the AP-MTL network is provided at the bottom.

egy to reshape the query feature space and alleviate difficulties in identity matching. The authors report improved performance on instrument segmentation using their approach.

These papers (Islam et al., 2020a; Sanchez-Matilla et al., 2021; Huang et al., 2022a; Baby et al., 2023; Zhao et al., 2022) demonstrate cases where the utilization of multiple tasks leads to improvements in the performance of all tasks. Psychogyios et al. (2022) highlights that MTL can face challenges as well. The authors present a method for jointly learning both disparity estimation and surgical instrument segmentation, but the results show a decrease in performance. The architecture employed in Psychogyios et al. (2022) is U-Net-like, featuring a shared encoder and separate decoder heads for each task, including segmentation. The model undergoes pretraining using the Flythings3D dataset (Mayer et al., 2016) for both the encoder and the disparity decoder. Following pretraining, Psychogyios et al. (2022) present results from various training schemes. These include MTL training alone, training on one task and fine-tuning on the other, and training with MTL followed by fine-tuning on a single task. Interestingly, the authors observed decreased performance when utilizing plain MTL for their two tasks. Instead, their experiments revealed that for marginal improvements in the segmentation task, it is necessary to first train the model using MTL and then fine-tune solely on segmentation. Conversely, for training the disparity task, it proves more effective to train the disparity task alone without segmentation.

In the previous papers, the focus is on enhancing the performance of two or more tasks by training them together. However, an alternative approach is to prioritize one primary task and introduce another task as an auxiliary task (Qin et al., 2020a; Bhattarai et al., 2023; Huang et al., 2022b; Wang et al., 2022a). The role of the auxiliary task is to guide the training of the primary perceptual task. These auxiliary tasks can take different forms. They can be real tasks with specific objectives (Huang
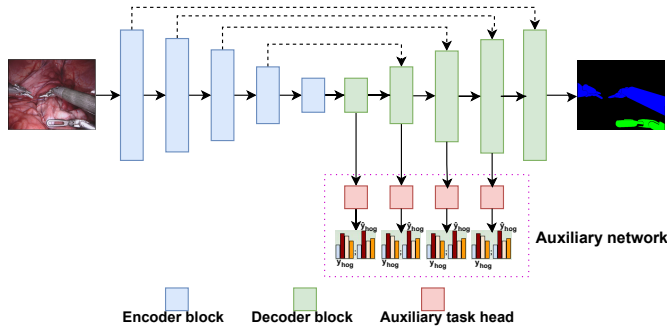
**Fig. 5:** Histogram of gradient multitask learning (HoG-MTL) (Bhattarai et al., 2023) demonstrates the use of an unconventional auxiliary task of predicting the histogram of gradients of input images image as an auxiliary task for semantic segmentation.

et al., 2022b), or derived tasks constructed from existing information in the image or labels (Wang et al., 2022a; Bhattarai et al., 2023; Qin et al., 2020a). In some cases, auxiliary tasks are used to inject domain-specific or prior knowledge into neural networks, enriching their capacity to learn and generalize (Wang et al., 2022a; Huang et al., 2022b).

In their work, Qin et al. (2020a) introduce the concept of contour supervision, a form of boundary prediction to serve as an auxiliary task for semantic segmentation. Contour supervision involves the creation of object outlines or contours for class instances within a semantic segmentation, with the network tasked with predicting these contour maps. The authors argue that contour prediction is valuable for localizing precise edges of the segmentation mask and providing information about the outer shape of objects.

In a different approach illustrated in Fig. 5, Bhattarai et al. (2023) suggest employing the prediction of the Histogram of Gradients (HoG) of the image as an auxiliary task to guide the learning of the segmentation task. The authors contend that leveraging image features as pseudo-labels in image classification is a well-established practice, and the histogram of gradient is a widely used handcrafted feature for object detection. Therefore, learning HoG can serve as valuable supervision for the segmentation task. To implement this approach, the authors propose a network that utilizes either the U2-Net (Qin et al., 2020b) or a U-Net with a single encoder and single decoder with deep supervision (Lee et al., 2015). Auxiliary heads are integrated into the architecture to predict the histogram of gradients at various decoder levels.

In the paper Huang et al. (2022b), the authors focus on depth estimation and adopt the standard left-right consistency methodology, as introduced in Godard et al. (2017). However, they introduce a novel auxiliary task by considering the 3D left-right consistency of 3D point clouds. The authors argue that while left-right consistency in disparity images is valuable, an additional source of information can be derived from the left-right consistency of point clouds, generated using predicted disparity images and information about laparoscopic stereo cameras. Their method follows a two-stage MTL approach. In the first stage, a standard depth estimation architecture predicts 2D left and right disparity images. Subsequently, these disparity

images, along with the focal length, the stereo distance, and predetermined blind masks for outlier removal, are used to generate 3D point clouds. The 2D depth loss functions include an appearance matching loss, a smoothness loss, and a left-right disparity consistency loss. In addition, the authors employ the iterative closest point (ICP) algorithm (Besl and McKay, 1992) to create a loss function for 3D left-right consistency by using the final residual registration error after ICP minimization. These losses are combined using linear scalarization. Incorporating the 3D auxiliary task is a key feature, as it enables the integration of information about the stereo camera into the optimization process of the neural network.

In endoscopic submucosal dissection (ESD), dissection landmarks are crucial for marking the boundary between a lesion and normal tissues. ESD involves creating landmarks around the lesion to label the boundary between the lesion and normal tissues (Ono et al., 2021). Wang et al. (2022a) present a neural network designed for detecting dissection landmarks. This approach stands out by introducing an auxiliary task for capturing the spatial relationship between each dissection landmark. Essentially, the authors seek to enhance the detection of dissection landmarks by incorporating domain knowledge that dictates alignment along a curve. To represent the spatial relationships between the nearest landmark neighbours on the curve, an edge map is generated. This edge map is then transformed into a heatmap that preserves these spatial relationships. The authors propose a shape-aware relation network based on the U-Net architecture featuring multiple decoders. This network functions as an MTL system that predicts both landmark positions and the heatmap of landmark spatial relationships.

Multitask learning for perceptual tasks has also been used for pituitary surgery. Das et al. (2023) present Pituitary Anatomy Identification Network (PAINet) and a newly introduced dataset for endoscopic pituitary surgery. They tackle the challenging task of identifying ten anatomical structures, with two prominent ones addressed through semantic segmentation and the other eight through centroid prediction. Their approach employs a U-Net architecture with two different task heads for each task for joint learning of the semantic segmentation and centroid prediction task.

### 3.2. MTL for tracking and control

#### 3.2.1. Pose estimation

Pose estimation is the process of determining the spatial configuration of objects or entities within a given scene, with a focus on estimating their positions and orientations. Pose estimation can be categorized into two branches based on the coordinate system in which the pose is measured: 2D pose estimation and 3D pose estimation. 2D pose estimation is particularly well-suited for scenarios where depth information is not essential or available. The goal is to accurately localize key points or landmarks that define the orientation of the object of interest within the image. On the other hand, 3D pose estimation involves estimating the full three-dimensional spatial configuration of objects in the scene. This estimation includes both the 2D position and depth information, along with orientation, making it a more challenging task. Accurate pose estima-

tion is of paramount importance in various applications related to endoscopic surgeries, including instrument tracking, automatic surgical camera control, augmented reality applications, and robotics. It is worth noting that the literature in the field of endoscopic surgeries often uses the 2D coordinate space (Li et al., 2022a; Laina et al., 2017; Hasan et al., 2021; Islam et al., 2019c, 2020b). Still, it may use algebraic and projective geometry techniques to derive 3D pose estimation if needed (Hasan et al., 2021). An exception to this practice among the papers reviewed is Li et al. (2022b), which directly operates in the 3D space.

In their paper, Laina et al. (2017) introduce a 2D instrument pose estimation method, employing the MTL paradigm. The two tasks trained for this purpose are segmentation and a variant of the instrument keypoint detection task. To achieve this, the authors reframe the 2D pose landmark detection task as a regression problem focusing on instrument heatmaps. The model utilized in this approach is a variant of the U-Net-like architecture with two different decoders: a segmentation decoder and a heatmap regression decoder. The output from the segmentation task is combined with the feature maps in the heatmap regression branch, a design choice that contributes to the improvement of predicted regression maps.

Hasan et al. (2021) present an MTL method for instrument presence detection, segmentation, and 2D pose estimation. The 2D pose estimations are represented as 'geometric primitives', which in their work refer to heat maps of critical components of the instruments, such as edge lines, mid-lines, and the tool-tip, as depicted in Fig. 6. The geometric primitives approach is chosen as they easily facilitate the calculation of 3D instrument poses. Unlike conventional 2D pose estimation methods that predict landmark positions (Li et al., 2022a) or heatmaps (Laina et al., 2017), this approach predicts a geometric primitive map of the relevant tool parts, including the tool edge, shaft mid-line, and tool-tip. The network architecture employed is U-Net-like, featuring a single encoder, direct instrument presence detection connected to the encoder, and four decoders for segmentation, edge-line primitive mapping, shaft mid-line primitive mapping, and tool-tip primitive mapping. The 2D geometric primitive map, once predicted, is then used in combination with prior information about the instrument, such as the radius of the tool shaft and the length of the tool head, to calculate the required 3D pose using algebraic and projective geometry techniques, followed by refinement.

Li et al. (2022a) introduce an alternative method for predicting 2D pose estimation while leveraging the MTL paradigm. The authors highlight a common challenge in the field - the abundance of datasets for tasks like surgical instrument segmentation, contrasted with the scarcity of datasets for pose estimation. To address this issue, they propose a multitask semisupervised approach to pose estimation. Their method predicts three different outputs: instrument segmentation, instrument number and instrument landmark keypoint detection. The instrument segmentation task is used as a spatial constraint, and the instrument number prediction task as a global constraint for the instrument landmark keypoint detection task. The semi-supervised technique employed here is the mean-
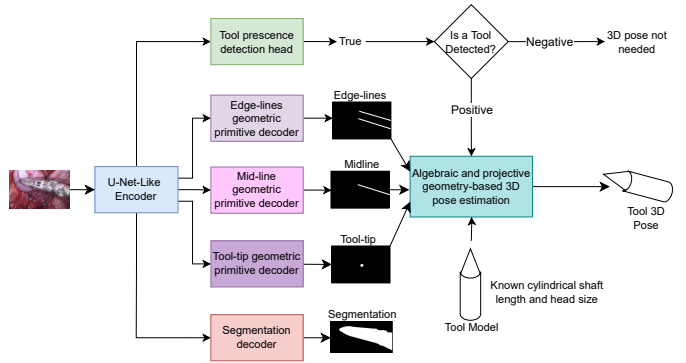


**Fig. 6:** The Augmented Reality Tool Network (ART-NET) (Hasan et al., 2021) presents a system to produce 3D pose estimations of instruments for augmented reality and 3D length measurement applications. As seen in the diagram, ART-Net simultaneously learns to predict tool presence, edge-lines, mid-lines, tool-tip and segmentation. The geometric primitives predicted are combined with prior information to produce the required Tool 3D pose.

teacher framework for semi-supervised learning (Tarvainen and Valpola, 2017). Both student and teacher models take the form of feature pyramidal networks (Lin et al., 2017), each with task-specific heads for the three tasks trained in conjunction.

As evident from the literature on instrument/camera tracking methods in MIS, such as (Hasan et al., 2021; Li et al., 2022a; Laina et al., 2017), the predominant approach involves predicting some form of 2D pose estimation rather than 3D pose estimation. This predicted 2D pose is usually used to facilitate other tasks such as tracking by prediction of 2D poses in a video sequence, as exemplified in Laina et al. (2017), or enabling augmented reality overlays, as demonstrated in Hasan et al. (2021). Another closely related and prevalent task is motion prediction for cameras, which also plays a significant role in the context of these tracking methods.

### 3.2.2. Camera motion prediction

The motion prediction task for cameras has been explored along two primary avenues: scanpath prediction and camera imitation learning. Scanpath prediction involves forecasting the potential camera paths based on the most captivating elements within the camera's current field of view. It draws inspiration from a theory in human gaze fixation, which posits a direct connection between human gaze scanpath and the attentional priority of objects in an image, often referred to as object saliency. In simpler terms, humans tend to focus on the most vital information first before shifting their gaze to less significant details. In the context of MIS, scanpath prediction tasks are centred around predicting the most salient objects in a scene and plotting a path from the most salient to the least salient areas. Foulsham (2019) provide a good reference for further insights on scanpaths, saliency, and their correlations.

Islam et al. (2019c) introduce a method for estimating saliency in surgical scenes, specifically by training in conjunction with segmentation as an auxiliary task. A notable challenge in this context is the absence of ground truth saliency datasets for MIS. To address this issue, the authors propose a novel method for generating saliency maps. The approach assumes
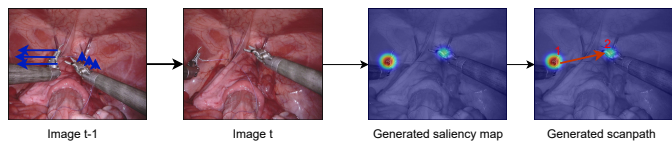
Image t-1          Image t          Generated saliency map          Generated scanpath

**Fig. 7:** Illustration of the heuristic used for the generation of saliency maps and scanpaths in Islam et al. (2019c). A saliency map is generated (using the code provided by Islam et al. (2019c) and images from EndoVis2017 (Allan et al., 2019)), which correlates to the motion of fixation points and the size of the fixation points, which are assumed to be located on the instrument wrist and claspers. The scanpath is then assumed to be the movement between the most salient to the least salient instrument.



**Fig. 8:** A visual representation of the different levels of granularities in surgical video workflow analysis. This representation for surgical video workflow analysis is adopted from the framework for holistic analysis of surgical videos in Valderrama et al. (2022).

that the most intriguing parts of surgical instruments are the wrist and claspers, and fixation points are located on these instrument segments. Additionally, the movement of instruments is utilized as a key indicator of saliency. Instruments with more significant movement are assigned higher saliency values than those with minimal movement. A scanpath is then defined as the movement from the most salient to the least salient instrument. Fig. 7 illustrates the process of generating a scanpath from images as presented in Islam et al. (2019c). The architectural framework employed is based on a U-Net structure featuring an encoder and two decoders, one dedicated to saliency map prediction and the other to instrument class segmentation. Attention modules are incorporated in the saliency map prediction branch to suppress irrelevant regions and highlight salient features.

Building upon the foundation laid by Islam et al. (2019c), the subsequent paper on scanpath prediction with MTL, known as ST-MTL (Islam et al., 2020b), maintains a similar architectural framework. However, the authors assert that the correct prediction of saliency maps requires information from the current frame and insights from previous frames. To address this temporal relationship, they introduce a convolutional Long Short-Term Memory (ConvLSTM) (Shi et al., 2015) into the saliency prediction decoder. The utilization of ConvLSTM modules enhances the model's ability to capture temporal dependencies and leverage information from past frames, contributing to more accurate saliency predictions. Moreover, Islam et al. (2020b) incorporate the ATO optimization method developed in Islam et al. (2020a), as discussed in Section 3.1.

Another promising avenue in the realm of motion prediction is camera imitation learning. This approach emphasises emulating the camera movements observed during surgeries performed by surgeons or their assistants. This technique is particularly relevant in the context of MIS, which often follows predefined procedural steps. These surgeries typically adhere to a fixed sequence of steps, and under normal circumstances, camera movements should exhibit similar patterns.

Li et al. (2022b) introduce a method for the proactive adjustment of the camera's field of view, achieved by modifying the camera's position in the *x*, *y*, and *z* coordinates. To replicate the motion patterns of a laparoscope camera, the authors employ a ConvLSTM for sequence modelling, predicting feature motions on a per-frame basis. A unique aspect of this research is the generation of ground truth laparoscope motion from laparoscopic videos. This process involves dynamic camera motion estima-
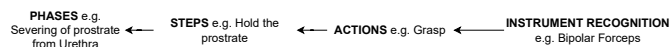
tion to infer camera pose. The authors utilize the Neural-Guided Random Sample Consensus (NG-RANSAC) method (Brachmann and Rother, 2019) to match stereo-images under dynamic conditions. Additionally, they apply the remote centre of motion constraint to optimize the pose estimation. The ground truth motion in the *x*, *y*, and *z* directions is derived from different sequential poses. The inputs to the ConvLSTM model comprise estimated segmentation and optical flow obtained from off-the-shelf models. The model optimizes for correct outputs of the subsequent *N* optical flow and segmentation results. The optical flow outputs are then passed through a laparoscopic action head to predict the camera's movements in the *x*, *y*, and *z* coordinates.

### 3.3. MTL for surgical video workflow analysis

Surgical video workflow analysis is the systematic examination of videos recorded during surgical procedures, aiming to extract valuable insights into various facets of the surgery. This analysis serves multiple critical purposes, including providing context-aware intraoperative assistance, enhancing surgeon training, facilitating procedure planning, supporting research endeavours, and enabling retrospective analysis (Lalys and Jannin, 2014).

The analysis process involves breaking down a surgical procedure video into distinct segments, which are categorized based on the surgeon's various activities. Different surgeries can be divided into specific activities, which can be examined at multiple levels of granularity.

We adopt a representation for surgical video workflow analysis from Valderrama et al. (2022), where the holistic analysis of a surgical video and its workflow is divided into four granularities. The initial granularity level involves identifying *instruments* present in the surgical scene, which are features of the surgical scene directly responsible for the surgical workflow. This stage entails understanding the instruments within the field of view as these instruments movements and interactions with different tissues give rise to specific *actions*, such as grasping or cutting, which represents the second granularity level. The third granularity stage involves the sequence of actions performed in pursuit of a particular surgical objective, which is called a *step*. Multiple steps are executed together to accomplish a portion of the surgery, termed a surgical *phase*, which is the fourth and final granularity of a surgical procedure. Fig. 8 provides a visual representation of the different levels of granularity.

Addressing the problem of surgical video workflow analysis is a critical step towards enhancing computer-aided interventions in surgical procedures. However, recognizing and distinguishing between different phases, steps, or actions within surgical videos remains a formidable challenge. This challenge is attributed to several factors, including the limited availabil-
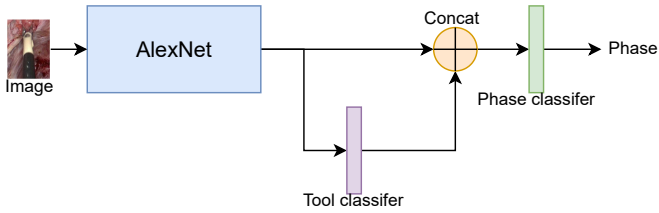
**Fig. 9:** One of the earliest works that attempts to solve the phase recognition problem is the EndoNet (Twinanda et al., 2017). It uses an Alexnet as a feature extractor with two task heads - A tool classifier head and a phase classifier head. Logits from the tool classifier head are concatenated with features from Alexnet before predicting phase.

ity of publicly accessible data, the high resemblance between long-range sequences belonging to different phases or steps, the substantial variability within sequences associated with a single phase or step, and the extended duration of surgical procedures necessitating thorough analysis.

In this section, we explore the literature which applies MTL to solve the surgical activity recognition problem and divide it into three subsections, namely: auxiliary tasks for surgical activity recognition, surgical action triplet learning, and multi-level activity recognition.

### 3.3.1. Auxiliary tasks for surgical activity recognition

As discussed in Section 2.2.3, auxiliary tasks are designed to aid in training primary tasks and enhance overall performance. Early research in surgical activity recognition with MTL primarily centred on phase recognition (Twinanda et al., 2017; Czempiel et al., 2020; Sanchez-Matilla et al., 2022; Jin et al., 2020). Notably, the prevalent auxiliary tasks for activity detection revolve around instrument recognition, specifically tool presence detection. The rationale for focusing on instrument recognition is rooted in the understanding that surgical instruments are the means by which surgeons interact with the surgical scene. Assuming standard surgical practices are followed, a substantial correlation exists between the choice of specific tools and the corresponding activities undertaken during a surgical procedure.

One of the pioneering papers that introduces the application of MTL and CNNs to surgical phase recognition is EndoNet (Twinanda et al., 2017). This method follows a two-stage prediction approach. The first stage jointly learns tool presence and phase prediction for each frame. The second stage refines the phase predictions generated in the initial stage. To implement this, EndoNet employs a pretrained AlexNet (Krizhevsky et al., 2012) as the encoder for its first stage, which is fine-tuned by simultaneously predicting tool presence and phase detection using dedicated task heads. In the second stage, EndoNet combines the features for phase detection in the current frame with information from the previous frames and feeds this combined data into a Hidden Markov Model (HMM) (Padoy et al., 2009). An architectural overview of EndoNet is illustrated in Fig. 9. In another approach, Twinanda et al. (2016) replaces the HMM in EndoNet with an LSTM (Hochreiter and Schmidhuber, 1997) to improve temporal modelling.

In a manner similar to the approach employed by Twinanda

et al. (2017), Czempiel et al. (2020) also leverage the auxiliary tool presence detection task along with a multistage network for phase recognition with different temporal modelling and optimization. Instead of utilizing the LSTM/HMM models as seen in previous works, they opt for the Temporal Convolutional Network (TCN), as described by Lea et al. (2017), for the refinement phase. The multitask optimization method used involves implementing the median frequency balancing technique (Eigen and Fergus, 2015). As for the input to the TCN, it comprises a concatenation of features extracted from the current frame and the $N$ preceding frames. TCN is preferred to LSTM/HMM approaches because it is notably faster, less resource-intensive and gives competitive results.

In a similar vein, Mondal et al. (2019) delve into the intricate relationship between tool presence detection and phase detection. The authors introduce a two-stage network identical to previous approaches, but Mondal et al. (2019) differ in how they link the primary with the auxiliary task. Unlike the methodologies found in (Czempiel et al., 2020; Twinanda et al., 2016, 2017), where refinement primarily concerns phase detection, this approach extends refinement to both tool and phase detection. Key to their methodology is the introduction of a joint probability loss function, which serves as the binding agent between the two tasks at each stage. The joint probability loss function is a product of the probabilities associated with tool and phase detection, weighted by the inverse of the tool's appearance frequency in a given phase. This approach seeks to establish a more integrated and interdependent relationship between the tasks, offering a fresh perspective on the problem.

Sanchez-Matilla et al. (2022) take a distinct approach by strongly emphasising enhanced temporal modelling by adopting a better auxiliary task for their two-stage network. In this paper, the authors opt for semantic segmentation as the auxiliary task, deviating from the conventional choice of tool presence detection. The first stage is a multitask network incorporating two distinct decoders, one for segmentation and another for phase prediction. The second stage of their network implementation employs a Temporal Convolutional Network (TCN) for the refinement phase. A notable finding of their research is the potential for performance enhancement with as little as 5% of the data labelled for segmentation.

In their research, Jin et al. (2020) offer an end-to-end trained network that capitalizes on the strong correlation between tool presence and surgical phases, much like previous works. However, it distinguishes itself by adopting a single-stage network design. The authors introduce a novel network known as MTRCNet-CL, featuring a shared encoder with two branches. One branch is dedicated to tool presence detection, while the other is focused on phase recognition. The phase recognition decoder branch incorporates an LSTM for temporal modelling. An innovative aspect of their approach is formulating a correlation loss that models the intricate relationship between tool presence and the predicted phase in the decoders. Their hypothesis posits that given the correlation between these two tasks, the logit values of the two tasks should also exhibit a certain degree of correlation. To measure this correlation, they calculate the Kullback-Leibler (KL) divergence between the phase

logit values and the tool presence logit values.

As evident from the studies presented in this section, early research in surgical activity detection primarily concentrated on the surgical phase granularity. Researchers explored various architectural and loss function strategies to leverage the relationship between the auxiliary task of tool recognition, with tool presence detection being the predominant focus, and the core task of phase detection. However, it is important to note that tool presence detection has inherent limitations. It can introduce noise into the analysis and lacks the precision required to pinpoint where the surgical action is unfolding within the video. Moreover, beyond the tools, the specific tissues that surgical tools interact with during a procedure also offer valuable insights into the ongoing surgical activity. Additionally, the identification of the instrument, tissue, and the actions taking place in a scene can provide valuable clues about the ongoing surgical steps and phases (Valderrama et al., 2022). Each step and phase requires very specific actions performed on specific tissues, in a specific order, and with specific instruments.

### 3.3.2. Surgical action triplet learning

More recent literature has compellingly argued against the exclusive prediction of surgical workflow either directly or solely based on tool presence. It has been contended that such an approach proves inadequate for a comprehensive understanding of surgical scenes. Instead, these works shift their focus towards the more granular action recognition task, striving to establish connections between three crucial elements before attempting to predict activities with longer sequences: the specific instrument in use, the action being performed, and the target anatomy undergoing the procedure, often collectively represented as three distinct labels, <instrument, action, target>. This predictive task illustrated in Fig. 10 and involving these three distinct labels is often referred to as the 'surgical triplet prediction task' (Katić et al., 2014). The surgical action triplet task fundamentally represents a multilabel (instrument, action, target) multiclass classification problem, with the understanding of the relationship between the labels being paramount. Typically, the problem is framed with the input as a single frame (Nwoye et al., 2020, 2022b), or multiple consecutive frames (Sharma et al., 2023a), with the objective of predicting the <instrument, action, target> triplet.

In their study, Nwoye et al. (2020) train a multitask network tailored specifically for the surgical action triplet task. They proposed the Tripnet architecture, which employs an encoder to extract joint features and three dedicated decoders for each of the task components. The first decoder operates as a convolutional-based unit, with a dual purpose: it predicts instrument classes in the image and generates class activation maps. These class activation maps are essential for the functioning of the other two decoders. The second and third decoders are *class-activation-map-guided* convolutional units that extract features relevant to action recognition and target recognition, respectively. They use the class activation maps from the instrument decoder as a guide.

The class activation maps are concatenated with the features in the action and tissue target decoders. This approach is un-
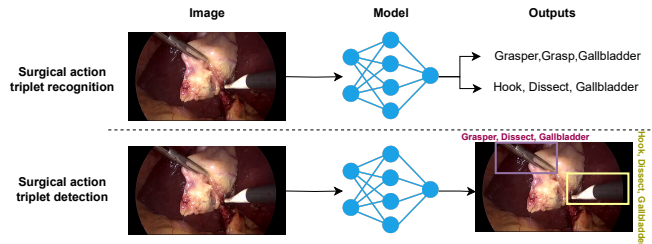


**Fig. 10:** The triplet task of recognizing the <instrument, action and target> is traditionally cast as surgical action triplet recognition, which is a multi-label multi-classification recognition problem (Nwoye et al., 2020, 2022b; Sharma et al., 2023a) as shown in the top image. Recently, there has been a progression in the task difficulty that cast this problem as surgical action triplet detection, which involves localizing all instruments and associating the corresponding instrument triplet to all localized instruments (Sharma et al., 2023b).
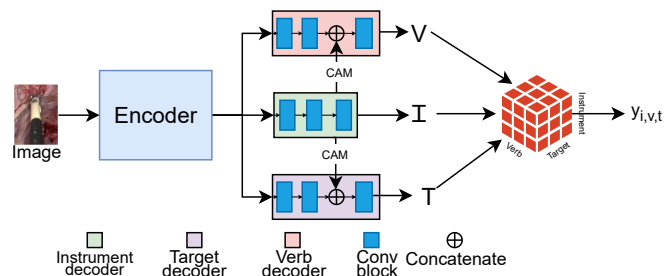


**Fig. 11:** The proposed Tripnet architecture for surgical action triplet recognition (Nwoye et al., 2020). The model contains an encoder and three separate branches for predicting instrument verb and targets. The class activation maps from the instrument branch are used as weak localization guides for action and target prediction.

derpinned by the hypothesis that the instrument plays a pivotal role in interacting with the surgical environment to initiate actions on a desired target. By incorporating a weak localization guide indicating the instrument's current location, the authors expect to enhance the performance of both the action recognition and tissue target recognition tasks. A visual representation of the Tripnet structure can be found in Fig. 11.

Furthermore, the authors observe that not all triplet combinations are feasible, and a data association problem arises as these components (predicted instruments, actions, and tissues) are interconnected. Instead of predicting each label separately, the authors opt for a joint prediction approach. The logits produced from each encoder (*I* for Instrument, *V* for action verb, and *T* for tissue target) are used to create a 3D interaction space volume (*Y*) through an outer product operation:

$$Y = \alpha I \otimes \beta V \otimes \gamma T$$

Here, $\alpha$, $\beta$, and $\gamma$ represent learnable weights. The 3D volume is quantized, with values above a chosen threshold accepted as valid triplets, while spaces in the 3D volume that can never form triplets are masked out. The loss function for this approach comprises the standard cross-entropy loss for all tasks, combined into a linear combination of losses.

Nwoye et al. (2022b), also known as RDV, is a sequel to the earlier work (Nwoye et al., 2020). This paper introduces

a noteworthy improvement by incorporating repeated attention mechanisms to facilitate inter-task knowledge transfer, drawing inspiration from the transformer architecture. Similar to Nwoye et al. (2020), RDV employs a single encoder with three decoders. However, it distinguishes itself by employing attention-based decoders, referred to as the *Class Activation Guided Attention Mechanism* (CAGAM).

The CAGAM mechanism replaces class-activation-map-guided convolutional units that extract features relevant to action recognition and target recognition used in the Tripnet architecture (Nwoye et al., 2020). Similar to the Tripnet architecture, It leverages class activation maps (CAM) from the Instrument decoder as a source of weakly supervised spatial guidance for the action and tissue target recognition tasks. But instead of just concatenating, the CAM is used to produce Key and Query vectors to form an attention matrix, which is used to scale the features from the Targets and Actions. Following the CAGAM Mechanism, RDV yields instrument, action, and target features, which are further processed through another attention mechanism. This mechanism is a variant of the multi-head attention models used in transformers. This process is integral to combining the features and is notably the second application of the attention mechanism. This dual-attention mechanism strategy is what the authors aptly refer to as a *rendezvous*. The rendezvous attention mechanism employs self-attention and cross-attention to capture complex semantic features in the instrument, verb, and target features. In contrast to the volume-based predictions in Nwoye et al. (2020), RDV opts for a simpler classifier for prediction tasks. To fine-tune the model and optimize its performance, RDV employs uncertainty weighting (Cipolla et al., 2018) to automatically determine the hyperparameters for each loss function.

The third instalment in this series of papers, Sharma et al. (2023a) called Rendezvous in Time (RIT), follows the preceding works. While RDV focuses solely on single-frame features for triplet recognition, RIT incorporates temporal modelling into the RDV model. In particular, Sharma et al. (2023a) introduce the Class Activation Guided Temporal Attention Module (CAGATAM), designed to enhance verb prediction and build upon the foundation of CAGAM. The Temporal Attention Module (TAM) plays a role in the temporal fusion of verb features extracted from the current and past frames, weighted by attention scores.

Yamlahi et al. (2023) apply the principle of self-distillation to solve the problem of class imbalance and label ambiguity in surgical action triplet recognition task. The method utilizes MTL and ensemble models as regularization to improve performance. A single teacher model with a Swin Transformer (Liu et al., 2021) backbone and a classifier are trained on hard labels with binary cross entropy loss for the three tasks of instrument, action and target detection separately. Then, the Sigmoid probabilities from the teacher model are used to train three self-distilled student models as an ensemble to minimize a distillation loss.

In the latest addition to the RDV series of papers, Sharma et al. (2023b) introduce the surgical action triplet detection task: the localization of surgical instruments along with the recogni-

tion of surgical action triplets. The authors address a challenge in datasets designed for binary triplet label recognition when used for detection tasks. Specifically, there may be multiple instruments in an image from the CholecT50 dataset (Nwoye et al., 2022b), but only the primary instrument that is most prominent is annotated. This leads to a data association problem. To tackle this issue, the authors propose a two-stage network named the Multi-class Instrument-aware Transformer - Interaction Graphs (MCIT-IG). The first stage is referred to as a multi-class instrument-aware Transformer (MCIT) that performs target prediction sub-task by using the knowledge of the current instruments and their classes. It does this by first detecting all the instruments in an image using a deformable DETR model (Zhu et al., 2020) trained on an annotated Cholec80 dataset, then extracting global features from the same image by using a feature extractor chained with a small transformer to provide global features, and finally combining the instrument detection information (from DETR) and image global feature information (from feature extractor) with a lightweight transformer to predict targets for each instrument. The second stage (IG) utilizes the instrument and target features from the first stage to construct a bipartite graph with action relationships as interaction edges. Although the supervision for the second stage is weak, a heuristic is employed to address the data association problem. The authors note that the better the instrument localization prediction, the better the accuracy score for surgical action triplet prediction is.

In their work, Chen et al. (2023) point out the challenges of jointly optimizing three distinct classification problems as a single multiclass multilabel problem. They highlight that this problem is unbalanced, as positive results are only achieved when all three components of the triplet are predicted correctly. Moreover, the presence of multiple instruments in a single image and the lack of annotations for some instruments further introduce ambiguities to the triplet recognition task. To address the complexities of triplet recognition, the authors propose a solution involving five smaller subnetworks. The first subnetwork is responsible for counting the number of tools and predicting the presence of key triplets or irrelevant triplets, including null actions and null targets. The second subnetwork predicts the tool classes in the image, and similar to the Tripnet approach (Nwoye et al., 2020), it utilizes class activation maps. However, in this case, the Inflated 3D Convolutional Network (I3D) (Carreira and Zisserman, 2017) is used, which allows for the generation of class activation maps, which the authors found to be more accurate. The third and fourth subnetworks jointly predict verbs and targets. Finally, the fifth subnetwork serves as both a fine-tuning and masking network, removing impossible triplets. It takes the logits predicted by the previous subnetworks and the current video clip as input to predict fine-tuned triplet logits and perform classification. The authors emphasize that they train each subnetwork in different stages, as attempting to address multiple auxiliary tasks simultaneously can lead to task distraction and negative transfer.
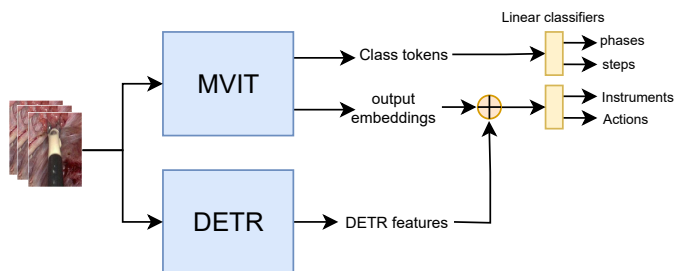
**Fig. 12:** The TAPIR architecture from Valderrama et al. (2022) for the prediction of phases, steps, instruments, and actions. It utilizes two separate encoders, a video feature extractor (MViT), and an instrument detector (DETR), along with various classification heads.

### 3.3.3. Multi-granularity activity detection

A comprehensive understanding of temporal relationships in MIS hinges on a network's ability to grasp the concept that multiple actions are required to complete a step and that multiple steps collectively constitute phases. This interconnection of granularities in surgical activities can serve as valuable signals for training deep neural networks in surgical activity recognition. Despite the limited available datasets, there is a growing interest in multi-granularity surgical activity learning, spurred by the recent introduction of multiple datasets providing labels for various granularities (Wang et al., 2022b; Huaulmé et al., 2021; Valderrama et al., 2022).

In an early endeavour to simultaneously predict phase and step activities, Ramesh et al. (2021) offer an innovative approach. Their proposed network adopts a two-stage design similar to the architecture of EndoNet (Twinanda et al., 2017). The initial stage comprises a feature extractor with a classifier, where two distinct linear classifier heads are responsible for predicting phase and step for $N$ frames. In the subsequent stage, the vectors extracted in the first stage for multiple consecutive frames are concatenated and injected into a temporal convolutional neural network for temporal modelling, which leads to improved predictions for both phase and step activities.

The authors of Valderrama et al. (2022) aim to completely address the temporal scene understanding problem in MIS. They introduce a multi-level activity detection model illustrated in Fig. 12 named Transformer for Action, Phase, Instrument, and Steps Recognition (TAPIR) to achieve this goal. TAPIR leverages two distinct backbones to capture various types of information. The first, the Multiscale Vision Transformer (MViT) (Fan et al., 2021), for extracting global and temporal features from sequential video frames. The second, Deformable Transformers for End-to-End Object Detection (Deformable DETR), focuses on capturing spatial features relevant to instrument detection and box proposals. The authors use a concatenation approach to combine the insights from these two backbones. The concatenated features from the MVIT and DETR are passed to linear classifiers to predict the action and instrument detection tasks. As for the prediction of phases and steps, this is solely carried out using the class tokens from the MViT backbone and linear classifiers.

### 3.4. MTL for surgical skill assessment

Traditionally, surgical skill assessment involves senior surgeons observing and evaluating less experienced counterparts, utilizing standardized rating checklists such as the Objective Structure Assessment of Technical Skills (OSATS) (Martin et al., 1997). However, the demand for training more doctors exceeds the number of experienced surgeons. Consequently, an automated system capable of scoring a surgeon's skill is invaluable for trainees.

In Jian et al. (2020), the authors introduce a multitask network designed to tackle the dual objectives of overall surgical skill level classification (categorizing surgeons as expert, intermediate, or novice) and attribute score regression (modified OSATS attributes from datasets like JIGSAWS (Gao et al., 2014)). This network architecture revolves around a shared encoder featuring two heads: one for classification and the other for the regression task. The core of this encoder is a variation of I3D (Carreira and Zisserman, 2017), a proven feature extractor in the domain of action recognition. To enhance the I3D capabilities, the encoder incorporates attention modules that enable it to focus on crucial segments of the video clips. To facilitate efficient computation, the authors split the input videos into $K$ equal parts, and clip snippets are sampled within these segments. Subsequently, a classification head is employed to predict the overall surgical skill assessment score, while a regression head predicts individual attribute scores. The classification and regression scores for each snippet are then aggregated, with the average of output features forming the final solution.

In Wang et al. (2020), the focus is on achieving interpretable skill assessment inspired by how senior clinicians assess other surgeons' performance. The authors observe that senior clinicians assess how well each surgical gesture is executed. The authors propose a two-stage network to replicate this assessment method. The first stage primarily aims to detect surgical gestures in different clips along with predicting the skill level classification and skill level score as auxiliary tasks. The surgical video input is initially split into $C$ clips, designed to streamline the computational process. These clips are then processed by a shared C3D encoder adapted from Tran et al. (2015). The features obtained from all $C$ clips are concatenated and used as the input to the decoder. The surgical gesture recognition decoder adopts a multi-stage temporal convolutional network. In parallel, the skill level classification and skill level prediction are executed through a shared LSTM, with classification and regression heads. The second stage focuses on the task of predicting skill assessment for the gestures identified in the first stage. To achieve this, the predicted gestures from the initial stage serve as a basis for segmenting the original surgical video into distinct gesture clips. These clips are then utilized as input for a C3D network, initialized with the weights from the C3D encoder employed in the first stage. Similarly, the decoder for predicting gesture level immediate skill score is implemented using an LSTM, featuring a dedicated head for gesture level immediate skill score prediction.

### 3.5. MTL for surgical report generation

The automatic generation of surgical reports can free surgeons and nurses from the tedious task of document entry, al-
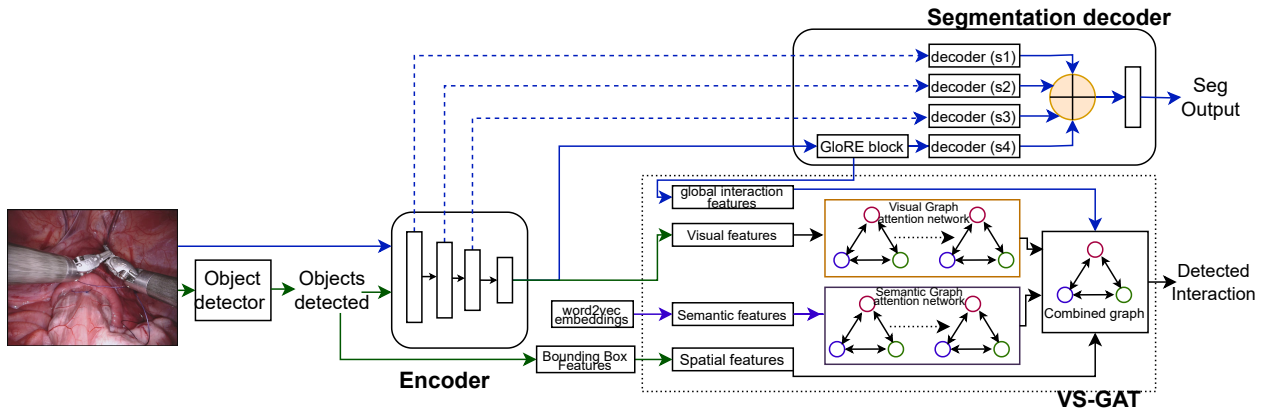
**Fig. 13:** The proposed architecture for globally-reasoned multi-task surgical scene understanding in Seenivasan et al. (2022b) for performing instrument segmentation and tool-tissue interaction. The model features a feature encoder, global and local reasoning for instrument segmentation and a VS-GAT (Liang et al., 2021) model for interaction detection.

lowing them to focus more on patients and post-operative interventions (Xu et al., 2021; Lin et al., 2022). Currently, in the existing literature, the approach to the surgical report generation task primarily revolves around frame-by-frame scene captioning (Xu et al., 2021; Lin et al., 2022; Seenivasan et al., 2022b, 2023b). The increasing interest in surgical report generation is closely tied to advancements in scene captioning and scene graph generation within the broader context of scene captioning research (Liang et al., 2021; Cornia et al., 2020) Notably, MTL has been applied to this field, incorporating aspects like image captioning and scene graph generation alongside other tasks to further enhance the overall capabilities of the system.

In Seenivasan et al. (2022b), a new approach is introduced for optimizing scene graphs to facilitate object-to-object interactions, focusing on the surgical report generation automation task. The authors frame this challenge as a multitask problem involving two core tasks: instrument segmentation and tool-tissue interaction. The authors propose a network with a single encoder and two decoders, as visualized in Fig. 13. The process commences with a preprocessing step, which generates bounding boxes and their corresponding classifications. These bounding boxes serve as a foundation for relationship modelling. The encoder, based on ResNet18, is responsible for shared feature extraction. The first decoder concentrates on local and global reasoning for instrument segmentation, configured as a U-net-like decoder with two noteworthy differences. The first is that the outputs of each lower-resolution decoder block do not feed into the next higher-resolution. Instead, each decoder block produces an output directly, which is concatenated together and passed through a conv block to produce the final segmentation prediction. Additionally, a GloRE block (Chen et al., 2019) is integrated into the output features of the encoder, enriching the global reasoning capabilities in the feature space. The features obtained from the GloRE block, known as Global Reasoned features, are also leveraged in the tool-tissue interaction decoder. The tool-tissue interaction decoder is structured as a scene graph from which the surgical report is generated. The authors employ the visual-semantic graph attention network (VS-GAT) (Liang et al., 2021). This network essentially comprises two components: a visual graph attention network and a semantic graph attention network. Both networks are harmoniously combined to create a unified graph. Notably, the authors introduce global interaction features from the GloRE block into this combined graph before instrument-tissue interaction prediction.

Seenivasan et al. (2023b) explore surgical report generation as a two-task problem of scene graph optimization for object-to-object relationships in a scene and frame-by-frame image captioning. The network architecture developed for this purpose is built around a shared encoder and two decoders. Based on the standard ResNet-18, the shared encoder forms the foundational feature extraction component. The object-to-object scene graph optimization decoder leverages the visual-semantic graph attention network from Liang et al. (2021). This aspect of the network focuses on enhancing the understanding of complex relationships between objects within the surgical scene. In parallel, the frame-by-frame image captioning task is addressed through the utilization of a meshed-memory transformer from Cornia et al. (2020). The optimization process is facilitated through ATO (Islam et al., 2020a). Recognizing the challenge of model generalization to target domains and the need to accommodate factors like new instruments, the authors introduce a loss function inspired by continual learning, termed class incremental contrastive loss.

### 3.6. Large models for solving multiple tasks

Large models pretrained on extensive datasets have demonstrated remarkable generalization capabilities (Radford et al., 2021, 2019). They exhibit emergent properties, such as zero-shot learning (Radford et al., 2021) and the capability to solve multiple tasks (Radford et al., 2019; Kirillov et al., 2023). In the current literature, there exist large models which exhibit this capability of solving multiple tasks in surgical scene understanding.

Seenivasan et al. (2022a) introduce the problem of visual question answering (VQA) in the context of surgery. They develop a large model capable of providing textual answers to a

wide range of questions, including tasks such as tissue presence recognition, instrument localization, tool presence recognition, and action recognition. To achieve this, they expand the Endovis2018 and Cholec80 datasets by adding sentence-based and classification-based answers to a predefined set of questions. The authors propose two models based on the VisualBERT architecture (Li et al., 2020) for classification-based and sentence-based answering in surgical VQA. They improve the VisualBERT encoder by introducing cross-token and cross-channel submodules to enhance the interaction between visual and text tokens. For the classification-based model, they utilize linear classifiers, with initial sentence classification determining the appropriate linear layer. The sentence-based model employs a standard transformer decoder.

In subsequent work, SurgicalGPT (Seenivasan et al., 2023a) for the surgical VQA task, the authors replace the VisualBERT architecture with GPT-2 (Radford et al., 2019), enhancing it with a vision encoder. They concatenate tokens from the text and vision encoders before feeding them to the GPT decoder, which they refer to as LV-GPT. Additionally, they change the unidirectional attention to a bidirectional attention model in the GPT decoder. Notably, they report improved performance compared to Seenivasan et al. (2022a).

Another application of the concept of generalizable tasks is *promptable segmentation*, which primarily focuses on binary segmentation for ideas represented in various forms, including text (Zhou et al., 2023), points (Kirillov et al., 2023), bounding boxes (Kirillov et al., 2023), and reference images (Lüddecke and Ecker, 2022). This approach allows for the execution of binary segmentation, instance segmentation, part segmentation, and instrument-type segmentation, provided that suitable prompts are provided.

An example of promptable segmentation in surgical scene understanding can be found in the work of Zhou et al. (2023). Their work addresses the challenge of distinguishing and segmenting diverse surgical instruments using textual prompts. To achieve this, they leverage pretrained CLIP text and vision encoder foundational models and introduce a novel text promptable mask decoder. The authors report capabilities to generalize text promptable segmentation to tissue segmentation, instrument-part, and instrument-type segmentation. Notably, their method demonstrates strong generalization capabilities, as evidenced by robust performance in cross-dataset evaluations.

Recent advancements in promotable segmentation have seen remarkable progress through methods based on the Segment Anything Model (SAM) (Kirillov et al., 2023). SAM serves as a foundational model that accommodates different prompting strategies, namely: points, bounding box, segmentation mask, and text.

SAM has demonstrated exceptional generalization abilities and has found applications in surgical segmentation tasks. An empirical study conducted in Wang et al. (2023) examines the robustness and generalizability of SAM when applied to the Endovis2017 and Endovis2018 datasets. Their findings reveal that pretrained SAM excels in terms of generalizability, especially when used with bounding box prompts, achieving state-of-the-art results. However, it is important to acknowledge that comparing bounding box prompts to class-based segmentation techniques might not be entirely fair.

Wang et al. (2023) also observe that SAM does not perform well for surgical instrument segmentation with point-based prompts, and its generalizability can be compromised under conditions of data corruption commonly encountered in surgical segmentation tasks. In summary, while SAM exhibits strong generalization capabilities, it relies on users to supply real-time bounding boxes for each image. In addition, SAM may not be robust in the presence of noise and data corruption. To address these challenges, several innovative approaches have been proposed.

MEDSAM (Ma and Wang, 2023) is a foundational model designed for universal medical image segmentation, curated from a dataset comprising over one million images, including CT scans, histology images, and surgical scenes. This model aims to bridge the gap between SAM for natural images and SAM for medical images. It adopts bounding box prompts as input, with a key distinction being the freezing of the SAM encoders, while only training the SAM decoder.

SurgicalSAM (Yue et al., 2023) proposes using the concepts of a prototype image, a real image of a particular instrument that captures the interested image, as a prompt instead of bounding boxes per image as prompts or directly using class names. More specifically, their method builds a memory bank of prototype images then they use the prototype-based class prompt encoder to exploit similarities between images in the dataset and class prototype images to create prompts, which they call prompt embeddings. In addition, the authors also propose a contrastive prototype learning loss to ensure that during training, the feature space for each different prototype is far from each other.

AdaptiveSAM (Paranjape et al., 2023) addresses the bounding box requirement by utilizing text as prompts. It employs text embeddings from pretrained CLIP and passes them through a trainable affine layer before applying them to the SAM prompt decoder. A notable feature of AdaptiveSAM is the introduction of *bias tuning* as a more memory-efficient method to adapt the SAM encoders. This approach trains the bias of the multi-head attention layers in the SAM image encoder and the normalization layers, achieving adaptation with higher efficiency.

All these methods (Ma and Wang, 2023; Yue et al., 2023; Paranjape et al., 2023) report significant improvement over the vanilla SAM approach when applied to surgical segmentation datasets over various tasks.

## 4. Public datasets for MTL in MIS

In this section, we explore the public datasets curated to support MTL for MIS. These datasets offer a valuable foundation for researchers to experiment, innovate, and address real-world MIS challenges.

Due to the inherent difficulties in producing surgical video datasets, the availability of public datasets suitable for MTL in this domain is currently limited. While there are some datasets designed specifically for MTL, such as the *multi-granularity surgical activity recognition datasets*, the number of datasets

**Table 1:** A summary of publically available multitask learning datasets for minimally invasive surgeries with information on the dataset characteristics and annotation characteristics.

| Dataset | Brief description | Procedure | Task | Input | Annotations | Paper |
|---|---|---|---|---|---|---|
| AutoLaparo | A dataset for image-guided surgical automation in laparoscopic hysterectomy comprising three sub-datasets: surgical workflow recognition, laparoscope motion prediction, instrument and anatomy segmentation. | 21 | Phase recognition task | 1388 minutes of surgical videos | 1388 minutes of phase labels | (Wang et al., 2022b) |
|  |  |  | Laparoscope motion prediction | 300 video clips | 300 next motion mode labels |  |
|  |  |  | Instrument part segmentation | 1800 keyframes | 1800 part segmentation maps |  |
|  |  |  | Key anatomy segmentation | 1800 keyframes | 1800 tissue segmentation maps |  |
| ART-Net | A non-robotic laparoscopic hysterectomy dataset designed for 3D graphics applications like 3D measurement and Augmented Reality (AR). | 29 | Tool presence prediction | 1500 frames | 1500 frames with tool presence labels | (Hasan et al., 2021) |
|  |  |  | Binary instrument segmentation | 635 keyframes | 635 binary instrument segmentation maps |  |
|  |  |  | Geometric map prediction | 635 keyframes | 635 geometric map labels |  |
| HeiSURF | A laparoscopic cholecystectomy designed for surgical activity recognition, full scene segmentation, and skill assessment. It is a parent dataset to the HeiChole dataset | 33 | Phase recognition | 23.3 hours of surgical videos | 23.3 hours of phase labels | (Bodenstedt et al., 2021) |
|  |  |  | Full scene segmentation | 827 keyframes | 827 full scene segmentation maps |  |
|  |  |  | Action recognition | 5514 instances in frames | 5514 action labels |  |
|  |  |  | Tool presence prediction | 6980 instances in frames | 6980 tool presence labels |  |
|  |  |  | Surgical skill-score prediction | 99 video clips | 495 skill scores |  |
| PSI-AVA | A dataset of robot-assisted radical prostatectomy designed for research into the complementary nature of surgical activity recognition tasks. | 8 | Phase recognition | 20.45 hours of surgical videos | 20.45 hours of phase labels | (Valderrama et al., 2022) |
|  |  |  | Step recognition | 20.45 hours of surgical videos | 20.45 hours of step labels |  |
|  |  |  | Action recognition | 5804 instances in frames | 5804 action labels |  |
|  |  |  | Instrument detection | 5804 instances in frames | 5804 bounding boxes with labels |  |
| HeiCo | A dataset of three different colorectal procedures (proctocolectomy, rectal resection, and sigmoid resection procedures) with emphasis on dataset generalization and diversity. | 30 | Phase recognition | 9.45 hours of surgical videos | 20.45 hours of phase labels | (Maier-Hein et al., 2021) |
|  |  |  | Instrument instance segmentation | 10,040 keyframes | 10,040 instance segmentation maps |  |
| JIGSAWS | An automatic gesture recognition and surgical skill assessment dataset containing videos of surgeons performing suturing, knot-tying and needle-passing on a bench-top model with a da Vinci. | 104 | Surgical gesture classification | 208 minutes of surgical videos | 208 minutes of surgical gesture labels | (Gao et al., 2014) |
|  |  |  | Surgical skill-score prediction | 104 video clips | 104 global surgical skill scores |  |
| MISAW | A micro-surgical anastomosis dataset with a focus on evaluating the impact of learning multiple surgical activity recognition tasks together. | 27 | Phase recognition | 1.5 hours of surgical videos | 1.5 hours of phase labels | (Huaulmé et al., 2021) |
|  |  |  | Step recognition | 1.5 hours of surgical videos | 1.5 hours of step labels |  |
|  |  |  | Action recognition | 1.5 hours of surgical videos | 1.5 hours of action labels |  |
|  |  |  | Tool presence prediction | 1.5 hours of surgical videos | 1.5 hours of Tool presence labels |  |
|  |  |  | target tissue prediction | 1.5 hours of surgical videos | 1.5 hours of tool target labels |  |
| PETRAW | A dataset containing the peg transfer task in laparoscopic surgery training. It is designed for multiple surgical activity recognition tasks. This dataset is collected from a virtual reality simulator | 150 | Phase recognition | 5.86 hours of surgical videos | 5.86 hours of phase labels | (Huaulmé et al., 2023) |
|  |  |  | Step recognition | 5.86 hours of surgical videos | 5.86 hours of step labels |  |
|  |  |  | Action recognition | 5.86 hours of surgical videos | 5.86 hours of action labels |  |
|  |  |  | target tissue prediction | 5.86 hours of surgical videos | 5.86 hours of tool target labels |  |
|  |  |  | Full scene segmentation | 5.86 hours of surgical videos | 5.86 hours of scene segmentation |  |
| SAR-RARP50 | A dataset of suturing segments of robotic assisted radical prostatectomy. The dataset focuses on tool segmentation and surgical action recognition. | 50 | Phase recognition | 3 hours of surgical videos | 3 hours of action labels | (Psychogyios et al., 2023) |
|  |  |  | Instrument segmentation | 1000 keyframes | 1000 Instrument segmentation maps |  |
| CholecT50 | The CholecT50 is designed for fine-grained action recognition and tool-tissue interaction in laparoscopic cholecystectomy surgeries. CholecT50 is the latest in the Cholec series of datasets. | 50 | Tool presence recognition | 100.9k frames | 100.9k frames with tool presence labels | (Nwoye et al., 2022b) |
|  |  |  | Action recognition | 100.9k frames | 100.9k frames with action labels |  |
|  |  |  | Tissue target recognition | 100.9k frames | 100.9k frames with target labels |  |
|  |  |  | Instrument detection | 13k instances in frames | 13k bounding boxes with labels |  |
| SARAS-MESAD | A dataset of both real and virtual human prostatectomy procedures containing. It is a dataset for research into surgical activity recognition and instrument detection with a focus on cross-domain learning | 9 | Instrument detection | 59k frames | 59k frames with bounding boxes | (Cuzzolin and Bawa, 2021) |
|  |  |  | Action recognition | 59k frames | 59k frames with action labels |  |

catering to other applications of MTL in surgical vision remains scarce. Nevertheless, it is worth noting that some authors have successfully explored approaches to leverage single-task datasets for MTL by generating self-supervised auxiliary tasks from the existing data. Moreover, there are other multitask datasets that have received less attention in the literature, which we aim to highlight. We present a comprehensive list of datasets designed to address multiple tasks in surgical vision. Additionally, we include two datasets initially intended for single-task applications but have been utilized in literature for MTL. For a quick overview of datasets designed for multiple tasks, along with information such as the amount of of images and labels, and other relevant information, readers can refer to Table 1. We do not include private datasets (e.g. ByPass40 (Ramesh et al., 2021)). For more information about surgical tool datasets, including non-multitask learning datasets, we recommend the following online resource: https://github.com/luiscarlosgph/list-of-surgical-tool-datasets.

### 4.1. Integrated Multi-tasks for Image-guided Surgical Automation in Laparoscopic Hysterectomy Dataset (AutoLaparo)

The AutoLaparo dataset (Wang et al., 2022b) supports image-guided surgical automation in laparoscopic hysterectomy, offering three sub-datasets: phase recognition, laparoscope motion prediction, and instrument segmentation with key anatomy annotation. The 21 procedures are recorded at 25 fps with a resolution of 1920x1080 pixels. These procedures are annotated with phase labels (7 phases). The laparoscope motion prediction sub-dataset consists of 300 clips, extracted from phases 2-4 of recorded procedures, each annotated with one of seven motion modes (up, down, left, right, zoom-in and zoom-out). The segmentation sub-dataset includes instruments (4 instruments) and key anatomy (1 anatomy) segmentation for keyframes in the motion prediction clips.

### 4.2. Augmented Reality Tool Network dataset (ART-Net)

The ART-Net dataset (Hasan et al., 2021) is tailored for non-robotic laparoscopic hysterectomy, emphasizing 3D graphics applications. Annotations cover tool presence detection, binary tool segmentation, and 2D pose estimation. Extracted from 29 procedures, the dataset provides frames with and without instruments, annotating these frames for tool presence. Keyframes are annotated for binary tool segmentation and 2D pose in the form of geometric primitives (tool-tip, midline and instrument shaft heatmaps).

### 4.3. HeiChole Surgical Workflow Analysis and Full Scene Segmentation dataset (HeiSURF)

The HeiSURF dataset (Bodenstedt et al., 2021) for day-to-day laparoscopic cholecystectomy includes annotations for phase recognition, full scene segmentation, action recognition, instrument presence detection, and skill assessment. Each procedure in the dataset is annotated for phase (7 phases), and keyframes are annotated for full scene segmentation and actions. Three videos per procedure are extracted and annotated with 5 skill annotations. The HeiSURF dataset is a parent dataset of an earlier dataset - the HeiChole dataset (This did not contain segmentation information) (Wagner et al., 2023).

### 4.4. Phase, Step, Instrument, and Atomic Visual Action recognition dataset (PSI-AVA)

The PSI-AVA dataset (Valderrama et al., 2022) focuses on robot-assisted radial prostatectomy for phase recognition, step recognition, instrument presence, instrument bounding boxes, and action recognition. The whole dataset is labelled for phase (10 phases) and step (20 steps) recognition. Keyframes annotated with bounding boxes provide detailed instrument detection (7 instruments) and corresponding actions (16 actions).

### 4.5. The Heidelberg Colorectal Dataset for Surgical Data Science in the Sensor Operating Room dataset (HeiCo)

The HeiCo (Maier-Hein et al., 2021) dataset features colorectal procedures for instrument segmentation and phase recognition. The dataset, derived from 30 procedures, covers proctocolectomy, rectal resection, and sigmoid resection. Phase recognition (14 phases) annotations are provided for the whole dataset. Keyframes with instance segmentation labels are also provided. A 10-second video snippet preceding each annotated keyframe is provided for temporal context. HeiCo contains data on surgical workflow analysis from the sensorOR challenge (Maier-Hein et al., 2021) and data from the Robust-MIS segmentation dataset (Ross et al., 2020).

### 4.6. The John Hopkins University - Intuitive Surgical Inc. Gesture and Skill Assessment Working Set dataset (JIGSAWS)

JIGSAWS dataset (Gao et al., 2014) contains annotations for manipulator gestures and surgical skills collected from recorded trials of eight medical doctors with varying skill levels performing suturing, knot-tying, and needle-passing task with the da Vinci Robot Surgical System (DSS). The dataset includes kinematic data (76D vector capturing the position, orientation, velocity and angle of the manipulator and camera), synchronized stereo video recordings (each at 30fps for approximately 2 minutes), and annotations for automatic gesture recognition (15 gestures) and skill assessment (a modified form of OSATS (Martin et al., 1997)).

### 4.7. The MIcro-Surgical Anastomose Workflow recognition on training sessions dataset (MISAW)

MISAW dataset (Huaulmé et al., 2021) focuses on multi-granularity surgical activity recognition in micro-surgical anastomosis. The dataset consists of 27 sequences recorded using a stereo-microscope, along with annotations for phase (2 phases), step (6 steps), action (17 actions), instrument presence (1 instrument), and instrument tissue targets (9 targets) for each frame. Synchronized kinematic data are also provided.

### 4.8. PEg TRAnsfer Workflow recognition by different modalities dataset (PETRAW)

The PETRAW dataset (Huaulmé et al., 2023) is designed for workflow recognition in peg transfer training sessions. The dataset comprises 150 sequences of peg transfer sessions recorded on a virtual reality simulator, kinematic data, videos, and annotations for semantic segmentation (2 targets and 1 instrument), phase (2 phases), step (12 steps), and action annotations (6 actions).

### 4.9. Surgical Instrumentation Segmentation and Action Recognition on Robot-Assisted Radical Prostatectomy dataset (SAR-RARP50)

The SAR-RARP50 dataset (Psychogyios et al., 2023) is designed to address data-scarcity challenges in surgical action recognition and tool segmentation for in vivo Robotic Assisted Radical Prostatectomy. The dataset comprises of 50 suturing segments acquired with a DaVinci Si Robot that features a stereo endoscope. The dataset is annotated for instrument segmentation (9 instruments) at keyframes and action (8 actions) for the whole dataset.

### 4.10. Cholecystectomy Action Triplet Datasets (Cholec series)

Cholec is a series of datasets containing laparoscopic cholecystectomy surgeries introduced to enable research on surgical workflow and tool-tissue interaction in the form of surgical action triplets. CholecT50 (Nwoye et al., 2022b), the most recent iteration of the Cholec series containing 50 sequences, offers annotations for six tasks: tool presence, action workflow, tool detection, target recognition, phase recognition, and surgical action triplets. The Cholec series of datasets includes Cholec80 (Twinanda et al., 2017), CholecSeg8k (Hong et al., 2020), CholecT40 (Nwoye et al., 2020), CholecT45 (Nwoye et al., 2022a), and CholecT50(Nwoye et al., 2022a,b). (Nwoye and Padoy, 2022) gives in-depth information about the Cholec series of datasets, its benchmarks, metrics, and other relevant information.

### 4.11. SARAS challenge on Multi-domain Endoscopic Surgeon Action Detection dataset (SARAS-MESAD)

SARAS-MESAD dataset (Cuzzolin and Bawa, 2021) facilitates surgical activity recognition research and cross-domain learning. It includes MESAD-Real and MESAD-Phantom sub-datasets, offering annotated frames for instrument detection and action recognition for human prostatectomy and phantom surgeries. This dataset builds on the previous SARAS-ESAD dataset (Bawa et al., 2021).

### 4.12. The Robotic Instrument Segmentation 2017 Sub-challenge Dataset (EndoVis2017)

EndoVis2017 (Allan et al., 2019) is a popular dataset that provides instrument segmentation data for endoscopic abdominal porcine procedures. The dataset is annotated for binary, parts, and type segmentation. To extend the EndoVis2017 dataset for MTL, researchers have employed various approaches to create auxiliary tasks, such as creating bounding boxes or scanpath heuristics from the existing segmentation masks (Laina et al., 2017; Islam et al., 2020b).

### 4.13. 2018 Robotic Scene Segmentation Challenge Dataset (EndoVis2018)

EndoVis2018 (Allan et al., 2020) is a full scene semantic segmentation dataset, a follow-up dataset to Endovis2017. The dataset comprises porcine robotic nephrectomy procedures annotated with full scene segmentation information. To extend the dataset to MTL, researchers followed a similar procedure

to Endovis2017 datasets. In addition, some authors separately annotated the Endovis2018 dataset to include extra information such as Tool-Tissue interaction (Seenivasan et al., 2022b), and Surgical VQA (Seenivasan et al., 2022a).

## 5. Discussion and conclusion

In this section, we consolidate our insights and observations from the MIS-related papers in Section 3 and connect them with generic MTL techniques for natural images introduced in Section 2. We also draw parallels with recent trends in the deep learning community.

### 5.1. Learning tasks for minimally invasive surgeries together

A prominent observation from the reviewed papers is that there are successful applications of MTL in minimally invasive surgeries (MIS). MTL not only leverages the inherent relationships between tasks but also allows for the extraction of richer information from surgical data. We have observed several methodologies for learning perceptual tasks together in this context. While some reported positive outcomes, others indicate negative transfer effects or negligible improvements, demonstrating the need for in-depth investigations on which tasks assist each other or regularly cause negative transfer, similar to the work of Standley et al. (2020). Furthermore, we found a lack of studies focusing on learning perceptual tasks like motion flow, and normal estimation, suggesting an unexplored research direction or that these tasks may not work well in the multitask learning framework.

Another interesting avenue for future research would be to systematically compare the efficiency of various multitask optimization techniques for learning multiple tasks in MIS compared to standard linear scalarization with grid search similar to the works of Kurin et al. (2022) and Xin et al. (2022). There is a noted rarity of multitask optimization methods in MIS, which should be investigated.

The architecture of how multiple tasks are learned together in the context of MIS presents an interesting observation. In many MIS papers, the standard approach involves the use of hard parameter sharing, often with a single-stage or two-stage architecture. However, one noticeable absence is the use of soft parameter sharing in MTL for MIS. Soft parameter sharing can allow tasks to share information at different levels and determine how to share information. This could be because of the complicated nature of soft parameter sharing networks or that hard parameter sharing is just good enough in most MIS scenarios. An area of exploration in future research could be to determine if soft parameter sharing can provide advantages over hard parameter sharing in the MIS context by comparing and contrasting the performance of both architectural styles on minimally invasive surgeries.

Another observation pertains to the impact of auxiliary tasks on learning accuracy. While many studies have shown the advantages of auxiliary tasks in enhancing task performance, it is essential to remain critical and consider possible negative transfers. The choice between directly incorporating domain-specific information into the primary network or predicting this information as an auxiliary task should be further explored. For instance, the prediction of contour images as an auxiliary task or a boundary loss should be evaluated for efficiency and informativeness.

Another noteworthy point is that only a few works incorporate information flow between decoders in the MIS field. This additional connectivity between decoders can facilitate inter-task relationship learning, which is a unique advantage of MTL. By allowing tasks to influence and inform each other, models can develop a deeper understanding of the underlying relationships between tasks and potentially improve overall performance. While learning better representations through MTL is valuable, it is essential to acknowledge that there are other learning paradigms, such as self-supervised learning, which excel in representation learning. However, the specific advantages of MTL in MIS, including inter-task relationship learning, make it a compelling choice for solving complex surgical tasks.

### 5.2. Multitask learning for automatic camera control

The discussion on automatic camera control in the context of MIS raises some interesting points. While the scanpath prediction task has been a significant development, it may not provide a comprehensive solution for automating camera movements in these surgeries. The approach of focusing primarily on surgical instruments (Gruijthuijsen et al., 2022; Islam et al., 2019c, 2020b), while valuable, may overlook the importance of capturing the surgical field comprehensively, including tissues and organs. Surgical procedures often involve dynamic movements where instruments enter and exit the field of view, and the camera's focus may need to adapt to these changes, as detailed by Brigham and Hospital (2019). This highlights the complexity of the task, where understanding the surgical context and deciding what to focus on is crucial.

An emerging approach to automating camera control is camera imitation and surgical intent learning (Huber et al., 2023; Li et al., 2022c). While this field is still relatively new, and there are only a handful of papers exploring it, it holds promise for improving camera control in minimally invasive surgeries. The release of datasets like AutoLaparo is expected to drive more interest and research in this area. Incorporating MTL into camera control, where tasks include predicting future segmentation and motion for the camera, is great. The future development of this field may involve adding more tasks, such as predicting ongoing and future actions, surgical steps and phases, to provide a holistic solution for camera control.

As the field of automatic camera control in MIS matures, it will be fascinating to see how researchers tackle the challenges of understanding and target domain generalization of surgical scenes to enhance camera movements and, ultimately, improve the surgical experience and outcomes.

### 5.3. Multitask learning for surgical activity detection

The evolution of surgical activity recognition from phase detection to action triplets and multi-granular activity detection is a notable development. This progression underscores

the importance of understanding the intricate relationships between different aspects of surgical activities, leading to improved recognition of complex surgical procedures and their contextual understanding.

In the early stages, much of the focus was on phase detection, which was partially influenced by the availability of datasets like Cholec80. Early models were primarily two-stage systems, incorporating temporal modelling to capture the dynamics of surgical activities. A significant shift has been observed in recent research, where surgical activity detection has been framed as the prediction of surgical action triplets, which include instrument, verb, and target. This approach provides a more detailed and informative way to describe surgical activities. It also highlights the need for understanding complex interactions between these components during surgical procedures.

The exploration of multi-granularity research, as exemplified by the PSI-AVA and MISAW datasets, offers valuable insights into recognizing different levels of surgical activities. It would be interesting to see how these methods that are surgery-specific and models that are trained specifically for multiple granularities of this specific procedure are generalized to other surgical procedures, surgeon styles, and surgical contexts.

As the field of surgical activity recognition continues to evolve, addressing the challenge of adapting models to different surgeries and achieving broader generalization is essential. Future research may focus on creating more versatile models that can understand and recognize surgical activities across various surgical procedures, improving the practical utility of these techniques in clinical settings.

### 5.4. Multitask learning for report generation

The application of MTL to surgical report generation is a promising avenue for improving the documentation and record-keeping aspects of minimally invasive surgeries. This complex task requires the generation of detailed and organized information about the events and procedures that occur in the intra- or per-operative period.

Two approaches have been explored in recent research: one involving training scene graphs with segmentation models to obtain reports and the other incorporating scene graphs and frame captioning techniques to generate reports. These approaches are initial steps towards addressing the challenge of surgical report generation, and they have shown potential for producing structured information from surgical data.

Future research in this area may focus on developing more advanced techniques that can capture and organize detailed surgical information comprehensively. This could involve more sophisticated approaches for scene graph generation with MTL, information aggregation, and natural language generation with foundation models, or multisystem approaches that utilize information from multiple specific task models with the ultimate goal of automating the process of generating surgical reports accurately and efficiently.

### 5.5. Large models in surgical scene understanding

The trend towards large models with the capacity to solve multiple problems by predicting a universal task is an intriguing development in the field of computer vision. Language models like the GPT series (Radford et al., 2019) have demonstrated the potential of understanding the nature of language through a single pretext task (next-token prediction) and subsequently applying that understanding to various language-related tasks. This notion of having a universal pretext task that can be leveraged for solving diverse problems is both promising and efficient.

While the concept of universal pretext tasks has gained traction in natural language processing, it is noteworthy that the vision domain is yet to have a similarly unified and versatile framework. Recent efforts, such as the SAM (Kirillov et al., 2023) and models derived from SAM, bridge this gap for segmentation, but this is just a single visual task. Looking forward, it is possible that developing a universal vision task and devising efficient methods for querying this universal task for specific applications could become a foundational approach for visual models. This paradigm would allow for a more streamlined and versatile way of addressing multiple vision-related challenges and potentially lead to breakthroughs in the field and, by extension, the computer vision for the MIS.

The generalizability of large models for multiple tasks is great. However, a main issue with large models in the surgical scenario is that for the networks to be utilized in Operating rooms (ORs), they need to perform inference in real-time and fit into devices and computers in the OR, which would be a cost-inefficient solution.

### 5.6. Conclusion

In conclusion, MTL has firmly established itself as an important paradigm within the domain of minimally invasive surgeries. Its influence extends across various facets of this specialized field,

This review analysed the applications of the MTL paradigm to minimally invasive surgeries. Firstly, the review gave an introduction to MTL and its objectives. Secondly, a detailed exploration of six distinct areas where MTL is applied in MIS was provided. Thirdly, the datasets that support MTL for MIS were presented. Lastly, we discussed some of the inferences and interesting observations on MTL and minimally invasive surgeries.

# References

Al Hajj, H., Lamard, M., Conze, P.H., Cochener, B., Quellec, G., 2018. Monitoring tool usage in surgery videos using boosted convolutional and recurrent neural networks. Medical image analysis 47, 203–218.

Allan, M., Kondo, S., Bodenstedt, S., Leger, S., Kadkhodamohammadi, R., Luengo, I., Fuentes, F., Flouty, E., Mohammed, A., Pedersen, M., Kori, A., Alex, V., Krishnamurthi, G., Rauber, D., Mendel, R., Palm, C., Bano, S., Saibro, G., Shih, C.S., Chiang, H.A., Zhuang, J., Yang, J., Iglovikov, V., Dobrenkii, A., Reddiboina, M., Reddy, A., Liu, X., Gao, C., Unberath, M., Kim, M., Kim, C., Kim, C., Kim, H., Lee, G., Ullah, I., Luna, M., Park, S.H., Azizian, M., Stoyanov, D., Maier-Hein, L., Speidel, S., 2020. 2018 robotic scene segmentation challenge. Eprint arXiv:2001.11190. arXiv:2001.11190.

Allan, M., Shvets, A., Kurmann, T., Zhang, Z., Duggal, R., Su, Y.H., Rieke, N., Laina, I., Kalavakonda, N., Bodenstedt, S., Herrera, L., Li, W., Iglovikov, V., Luo, H., Yang, J., Stoyanov, D., Maier-Hein, L., Speidel, S., Azizian, M., 2019. 2017 robotic instrument segmentation challenge. Eprint arXiv:1902.06426. arXiv:1902.06426.

Baby, B., Thapar, D., Chasmai, M., Banerjee, T., Dargan, K., Suri, A., Banerjee, S., Arora, C., 2023. From forks to forceps: A new framework for instance segmentation of surgical instruments, in: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, pp. 6191–6201.

Bawa, V.S., Singh, G., KapingA, F., Skarga-Bandurova, I., Oleari, E., Leporini, A., Landolfo, C., Zhao, P., Xiang, X., Luo, G., et al., 2021. The SARAS endoscopic surgeon action detection (ESAD) dataset: Challenges and methods. arXiv preprint arXiv:2104.03178 .

Bengio, Y., Courville, A., Vincent, P., 2013. Representation learning: A review and new perspectives. IEEE transactions on pattern analysis and machine intelligence 35, 1798–1828.

Besl, P.J., McKay, N.D., 1992. Method for registration of 3D shapes, in: Sensor fusion IV: control paradigms and data structures, Spie. pp. 586–606.

Bhattarai, B., Subedi, R., Gaire, R.R., Vazquez, E., Stoyanov, D., 2023. Histogram of oriented gradients meet deep learning: A novel multi-task deep network for 2D surgical image semantic segmentation. Medical Image Analysis 85, 102747.

Bodenstedt, S., Speidel, S., Wagner, M., Chen, J., Kisilenko, A., Müller, B., Maier-Hein, L., Oliveira, B., Hong, S., Zamora-Anaya, J., et al., 2021. Heichole surgical workflow analysis and full scene segmentation (HeiSurF).

Brachmann, E., Rother, C., 2019. Neural-guided RANSAC: Learning where to sample model hypotheses, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 4322–4331.

Brigham, Hospital, W., 2019. Robotic assisted laparoscopic radical prostatectomy — Brigham and Women's hospital.

Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S., 2020. End-to-end object detection with transformers, in: European conference on computer vision, Springer. pp. 213–229.

Carreira, J., Zisserman, A., 2017. Quo vadis, action recognition? a new model and the kinetics dataset, in: proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 6299–6308.

Caruana, R., 1997. Multitask learning. Machine learning 28, 41–75.

Chen, Y., He, S., Jin, Y., Qin, J., 2023. Surgical activity triplet recognition via triplet disentanglement, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer. pp. 451–461.

Chen, Y., Rohrbach, M., Yan, Z., Shuicheng, Y., Feng, J., Kalantidis, Y., 2019. Graph-based global reasoning networks, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).

Chen, Z., Badrinarayanan, V., Lee, C.Y., Rabinovich, A., 2018. GradNorm: Gradient normalization for adaptive loss balancing in deep multitask networks, in: Dy, J., Krause, A. (Eds.), Proceedings of the 35th International Conference on Machine Learning, PMLR. pp. 794–803.

Cheng, B., Schwing, A.G., Kirillov, A., 2021. Per-pixel classification is not all you need for semantic segmentation, in: NeurIPS.

Cipolla, R., Gal, Y., Kendall, A., 2018. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics, in: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 7482–7491.

Collobert, R., Weston, J., 2008. A unified architecture for natural language processing: Deep neural networks with multitask learning, in: ICML '08: Proceedings of the 25th international conference on Machine learning, Association for Computing Machinery, New York, NY, USA. p. 160–167.

Cornia, M., Stefanini, M., Baraldi, L., Cucchiara, R., 2020. Meshed-memory transformer for image captioning, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).

Cuzzolin, F., Bawa, V., 2021. Saras-mesad - benchmark and challenge. doi:10.13140/RG.2.2.10022.04160.

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: Martel, A.L., Abolmaesumi, P., Stoyanov, D., Mateus, D., Zuluaga, M.A., Zhou, S.K., Racoceanu, D., Joskowicz, L. (Eds.), Medical Image Computing and Computer Assisted Intervention – MICCAI 2020, Springer International Publishing, Cham. pp. 343–352.

Das, A., Khan, D.Z., Williams, S.C., Hanrahan, J.G., Borg, A., Dorward, N.L., Bano, S., Marcus, H.J., Stoyanov, D., 2023. A multi-task network for anatomy identification in endoscopic pituitary surgery, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer. pp. 472–482.

Doersch, C., Zisserman, A., 2017. Multi-task self-supervised visual learning, in: Proceedings of the IEEE international conference on computer vision, pp. 2051–2060.

Dorent, R., Booth, T., Li, W., Sudre, C.H., Kafiabadi, S., Cardoso, J., Ourselin, S., Vercauteren, T., 2021. Learning joint segmentation of tissues and brain lesions from task-specific hetero-modal domain-shifted datasets. Medical image analysis 67, 101862.

Eigen, D., Fergus, R., 2015. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture, in: Proceedings of the IEEE international conference on computer vision, pp. 2650–2658.

Fan, H., Xiong, B., Mangalam, K., Li, Y., Yan, Z., Malik, J., Feichtenhofer, C., 2021. Multiscale vision transformers, in: Proceedings of the IEEE/CVF international conference on computer vision, pp. 6824–6835.

Foulsham, T., 2019. Scenes, saliency maps and scanpaths, in: Klein, C., Ettinger, U. (Eds.), Eye Movement Research: An Introduction to its Scientific Foundations and Applications. Springer International Publishing, Cham, pp. 197–238.

Gao, Y., Ma, J., Zhao, M., Liu, W., Yuille, A.L., 2019. Nddr-cnn: Layerwise feature fusing in multi-task cnns by neural discriminative dimensionality reduction, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 3205–3214.

Gao, Y., Vedula, S.S., Reiley, C.E., Ahmidi, N., Varadarajan, B., Lin, H.C., Tao, L., Zappella, L., Béjar, B., Yuh, D.D., et al., 2014. JHU-IsI gesture and skill assessment working set (JIGSAWS): A surgical activity dataset for human motion modeling, in: MICCAI workshop: M2cai.

Garcia-Peraza-Herrera, L.C., Li, W., Fidon, L., Gruijthuijsen, C., Devreker, A., Attilakos, G., Deprest, J., Vander Poorten, E., Stoyanov, D., Vercauteren, T., et al., 2017. Toolnet: holistically-nested real-time segmentation of robotic surgical tools, in: 2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), IEEE. pp. 5717–5722.

Ghiasi, G., Zoph, B., Cubuk, E.D., Le, Q.V., Lin, T.Y., 2021. Multi-task self-training for learning general representations, in: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), pp. 8856–8865.

Girshick, R., 2015. Fast R-CNN, in: 2015 IEEE International Conference on Computer Vision (ICCV), pp. 1440–1448.

Godard, C., Mac Aodha, O., Brostow, G.J., 2017. Unsupervised monocular depth estimation with left-right consistency, in: CVPR.

Gruijthuijsen, C., Garcia-Peraza-Herrera, L.C., Borghesan, G., Reynaerts, D., Deprest, J., Ourselin, S., Vercauteren, T., Vander Poorten, E., 2022. Robotic endoscope control via autonomous instrument tracking. Frontiers in Robotics and AI 9, 832208.

Guo, M., Haque, A., Huang, D.A., Yeung, S., Fei-Fei, L., 2018. Dynamic task prioritization for multitask learning, in: Proceedings of the European Conference on Computer Vision (ECCV).

Hasan, M.K., Calvet, L., Rabbani, N., Bartoli, A., 2021. Detection, segmentation, and 3D pose estimation of surgical tools using convolutional neural networks and algebraic geometry. Medical Image Analysis 70, 101994.

He, K., Gkioxari, G., Dollár, P., Girshick, R., 2017. Mask R-CNN, in: Proceedings of the IEEE international conference on computer vision, pp. 2961–2969.

Heuer, F., Mantowsky, S., Bukhari, S.S., Schneider, G., 2021. Multitask-

centernet (mcn): Efficient and diverse multitask learning using an anchor free approach. 2021 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW) , 997–1005.

Hochreiter, S., Schmidhuber, J., 1997. Long short-term memory. Neural computation 9, 1735–1780.

Hong, W.Y., Kao, C.L., Kuo, Y.H., Wang, J.R., Chang, W.L., Shih, C.S., 2020. Cholecseg8k: A semantic segmentation dataset for laparoscopic cholecystectomy based on Cholec80. Eprint arXiv:2012.12453. arXiv:2012.12453.

Hu, J., Shen, L., Sun, G., 2018. Squeeze-and-excitation networks, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).

Hu, Y., Xian, R., Wu, Q., Fan, Q., Yin, L., Zhao, H., 2023. Revisiting scalarization in multi-task learning: A theoretical perspective. arXiv preprint arXiv:2308.13985 .

Huang, B., Nguyen, A., Wang, S., Wang, Z., Mayer, E., Tuch, D., Vyas, K., Giannarou, S., Elson, D.S., 2022a. Simultaneous depth estimation and surgical tool segmentation in laparoscopic images. IEEE Transactions on Medical Robotics and Bionics 4, 335–338.

Huang, B., Zheng, J.Q., Nguyen, A., Xu, C., Gkouzionis, I., Vyas, K., Tuch, D., Giannarou, S., Elson, D.S., 2022b. Self-supervised depth estimation in laparoscopic image using 3D geometric consistency, in: Wang, L., Dou, Q., Fletcher, P.T., Speidel, S., Li, S. (Eds.), Medical Image Computing and Computer Assisted Intervention – MICCAI 2022, Springer Nature Switzerland, Cham. pp. 13–22.

Huaulmé, A., Sarikaya, D., Le Mut, K., Despinoy, F., Long, Y., Dou, Q., Chng, C.B., Lin, W., Kondo, S., Bravo-Sánchez, L., et al., 2021. Micro-surgical anastomose workflow recognition challenge report. Computer Methods and Programs in Biomedicine 212, 106452.

Huaulmé, A., Harada, K., Nguyen, Q.M., Park, B., Hong, S., Choi, M.K., Peven, M., Li, Y., Long, Y., Dou, Q., Kumar, S., Lalithkumar, S., Hongliang, R., Matsuzaki, H., Ishikawa, Y., Harai, Y., Kondo, S., Mitsuishi, M., Jannin, P., 2023. Peg transfer workflow recognition challenge report: Do multimodal data improve recognition? Computer Methods and Programs in Biomedicine 236, 107561. URL: https://www.sciencedirect.com/science/article/pii/S0169260723002262, doi:https://doi.org/10.1016/j.cmpb.2023.107561.

Huber, M., Ourselin, S., Bergeles, C., Vercauteren, T., 2023. Deep homography prediction for endoscopic camera motion imitation learning, in: Greenspan, H., Madabhushi, A., Mousavi, P., Salcudean, S., Duncan, J., Syeda-Mahmood, T., Taylor, R. (Eds.), Medical Image Computing and Computer Assisted Intervention – MICCAI 2023, Springer Nature Switzerland, Cham. pp. 217–226.

Islam, M., Atputharuban, D.A., Ramesh, R., Ren, H., 2019a. Real-time instrument segmentation in robotic surgery using auxiliary supervised deep adversarial learning. IEEE Robotics and Automation Letters 4, 2188–2195.

Islam, M., Atputharuban, D.A., Ramesh, R., Ren, H., 2019b. Real-time instrument segmentation in robotic surgery using auxiliary supervised deep adversarial learning. IEEE Robotics and Automation Letters 4, 2188–2195.

Islam, M., Li, Y., Ren, H., 2019c. Learning where to look while tracking instruments in robot-assisted surgery, in: Shen, D., Liu, T., Peters, T.M., Staib, L.H., Essert, C., Zhou, S., Yap, P.T., Khan, A. (Eds.), Medical Image Computing and Computer Assisted Intervention – MICCAI 2019, Springer International Publishing, Cham. pp. 412–420.

Islam, M., Vibashan, V., Ren, H., 2020a. AP-MTL: Attention pruned multitask learning model for real-time instrument detection and segmentation in robot-assisted surgery, in: 2020 IEEE International Conference on Robotics and Automation (ICRA), IEEE. pp. 8433–8439.

Islam, M., Vibashan, V.S., Lim, C.M., Ren, H., 2020b. ST-MTL: Spatiotemporal multitask learning model to predict scanpath while tracking instruments in robotic surgery. Medical image analysis 67, 101837.

Jaffray, B., 2005. Minimally invasive surgery. Archives of disease in childhood 90, 537.

Jian, Z., Yue, W., Wu, Q., Li, W., Wang, Z., Lam, V., 2020. Multitask learning for video-based surgical skill assessment, in: 2020 Digital Image Computing: Techniques and Applications (DICTA), IEEE. pp. 1–8.

Jin, Y., Li, H., Dou, Q., Chen, H., Qin, J., Fu, C.W., Heng, P.A., 2020. Multitask recurrent convolutional network with correlation loss for surgical video analysis. Medical Image Analysis 59, 101572.

Kalashnikov, D., Varley, J., Chebotar, Y., Swanson, B., Jonschkowski, R., Finn, C., Levine, S., Hausman, K., 2021. Scaling up multi-task robotic reinforcement learning, in: 5th Annual Conference on Robot Learning.

Katić, D., Wekerle, A.L., 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: Stoyanov, D., Collins, D.L., Sakuma, I., Abolmaesumi, P., Jannin, P. (Eds.), Information Processing in Computer-Assisted Interventions, Springer International Publishing, Cham. pp. 158–167.

Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A.C., Lo, W.Y., Dollar, P., Girshick, R., 2023. Segment anything, in: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), pp. 4015–4026.

Kokkinos, I., 2017. Ubernet: Training a universal convolutional neural network for low-, mid-, and high-level vision using diverse datasets and limited memory, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).

Koonce, B., Koonce, B., 2021. Efficientnet. Convolutional Neural Networks with Swift for Tensorflow: Image Recognition and Dataset Categorization , 109–123.

Krizhevsky, A., Sutskever, I., Hinton, G.E., 2012. Imagenet classification with deep convolutional neural networks, in: Pereira, F., Burges, C., Bottou, L., Weinberger, K. (Eds.), Advances in Neural Information Processing Systems, Curran Associates, Inc.

Kurin, V., De Palma, A., Kostrikov, I., Whiteson, S., Kumar, M.P., 2022. In defense of the unitary scalarization for deep multi-task learning, in: Neural Information Processing Systems.

Laina, I., Rieke, N., Rupprecht, C., Vizcaíno, J.P., Eslami, A., Tombari, F., Navab, N., 2017. Concurrent segmentation and localization for tracking of surgical instruments, in: Descoteaux, M., Maier-Hein, L., Franz, A., Jannin, P., Collins, D.L., Duchesne, S. (Eds.), Medical Image Computing and Computer-Assisted Intervention - MICCAI 2017, Springer International Publishing, Cham. pp. 664–672.

Lalys, F., Jannin, P., 2014. Surgical process modelling: a review. International journal of computer assisted radiology and surgery 9, 495–511.

Lea, C., Flynn, M.D., Vidal, R., Reiter, A., Hager, G.D., 2017. Temporal convolutional networks for action segmentation and detection, in: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), IEEE Computer Society, Los Alamitos, CA, USA. pp. 1003–1012.

Lee, C.Y., Xie, S., Gallagher, P., Zhang, Z., Tu, Z., 2015. Deeply-Supervised Nets, in: Lebanon, G., Vishwanathan, S.V.N. (Eds.), Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics, PMLR, San Diego, California, USA. pp. 562–570. URL: https://proceedings.mlr.press/v38/lee15a.html.

Li, B., Li, S., Yang, J., 2022a. Multi-task semi-supervised learning framework for surgical instrument pose estimation, in: Proceedings of the 8th International Conference on Computing and Artificial Intelligence, Association for Computing Machinery, New York, NY, USA. p. 698–704.

Li, B., Lu, B., Wang, Z., Zhong, F., Dou, Q., Liu, Y.H., 2022b. Learning laparoscope actions via video features for proactive robotic field-of-view control. IEEE Robotics and Automation Letters 7, 6653–6660.

Li, B., Lu, B., Wang, Z., Zhong, F., Dou, Q., Liu, Y.H., 2022c. Learning laparoscope actions via video features for proactive robotic field-of-view control. IEEE Robotics and Automation Letters 7, 6653–6660. doi:10.1109/LRA.2022.3173442.

Li, L.H., Yatskar, M., Yin, D., Hsieh, C.J., Chang, K.W., 2020. What does BERT with vision look at?, in: Jurafsky, D., Chai, J., Schluter, N., Tetreault, J. (Eds.), Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, Association for Computational Linguistics, Online. pp. 5265–5275. URL: https://aclanthology.org/2020.acl-main.469, doi:10.18653/v1/2020.acl-main.469.

Li, W.H., Liu, X., Bilen, H., 2022d. Learning multiple dense prediction tasks from partially annotated data, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 18879–18889.

Liang, Z., Liu, J., Guan, Y., Rojas, J., 2021. Visual-semantic graph attention networks for human-object interaction detection, in: 2021 IEEE International Conference on Robotics and Biomimetics (ROBIO), pp. 1441–1447.

Liao, Y.H., Kar, A., Fidler, S., 2021. Towards good practices for efficiently annotating large-scale image classification datasets, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 4350–4359.

Liebel, L., Körner, M., 2018. Auxiliary tasks in multi-task learning. arXiv preprint arXiv:1805.06334 .

Lin, B., Zhang, Y., 2023. LibMTL: A Python library for multi-task learning. Journal of Machine Learning Research 24, 1–7.

Lin, C., Zheng, S., Liu, Z., Li, Y., Zhu, Z., Zhao, Y., 2022. SGT: Scene graph-guided transformer for surgical report generation, in: International conference on medical image computing and computer-assisted intervention, Springer. pp. 507–518.

Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S., 2017. Feature pyramid networks for object detection, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 2117–2125.

Lin, X., Zhen, H.L., Li, Z., Zhang, Q.F., Kwong, S., 2019. Pareto multi-task learning. Advances in neural information processing systems 32.

Liu, S., Johns, E., Davison, A.J., 2019. End-to-end multi-task learning with attention, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 1871–1880.

Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.Y., Berg, A.C., 2016. SSD: Single shot multibox detector, in: Leibe, B., Matas, J., Sebe, N., Welling, M. (Eds.), Computer Vision – ECCV 2016, Springer International Publishing, Cham. pp. 21–37.

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 (ICCV).

Lüddecke, T., Ecker, A., 2022. Image segmentation using text and image prompts, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 7086–7096.

Luo, H., Hu, Q., Jia, F., 2019. Details preserved unsupervised depth estimation by fusing traditional stereo knowledge from laparoscopic images. Healthcare technology letters 6, 154–158.

Luo, H., Wang, C., Duan, X., Liu, H., Wang, P., Hu, Q., Jia, F., 2022. Unsupervised learning of depth estimation from imperfect rectified stereo laparoscopic images. Computers in biology and medicine 140, 105109.

Ma, J., Wang, B., 2023. Segment anything in medical images. arXiv preprint arXiv:2304.12306 .

Maier-Hein, L., Wagner, M., Ross, T., Reinke, A., Bodenstedt, S., Full, P.M., Hempe, H., Mindroc-Filimon, D., Scholz, P., Tran, T.N., Bruno, P., Kisilenko, A., Müller, B., Davitashvili, T., Capek, M., Tizabi, M.D., Eisenmann, M., Adler, T.J., Gröhl, J., Schellenberg, M., Seidlitz, S., Lai, T.Y.E., Pekdemir, B., Roethlingshoefer, V., Both, F., Bittel, S., Mengler, M., Mündermann, L., Apitz, M., Kopp-Schneider, A., Speidel, S., Nickel, F., Probst, P., Kenngott, H.G., Müller-Stich, B.P., 2021. Heidelberg colorectal data set for surgical data science in the sensor operating room. Scientific Data 8, 101.

Martin, J.A., Regehr, G., Reznick, R., Macrae, H., Murnaghan, J., Hutchison, C., Brown, M., 1997. Objective structured assessment of technical skill (OSATS) for surgical residents. BJS (British Journal of Surgery) 84, 273–278.

Mayer, N., Ilg, E., Hausser, P., Fischer, P., Cremers, D., Dosovitskiy, A., Brox, T., 2016. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 4040–4048.

Milletari, F., Rieke, N., Baust, M., Esposito, M., Navab, N., 2018. CFCM: segmentation via coarse to fine context memory, in: Medical Image Computing and Computer Assisted Intervention–MICCAI 2018: 21st International Conference, Granada, Spain, September 16-20, 2018, Proceedings, Part IV 11, Springer. pp. 667–674.

Mishra, K., Sathish, R., Sheet, D., 2017. Learning latent temporal connectionism of deep residual visual abstractions for identifying surgical tools in laparoscopy procedures, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, pp. 58–65.

Misra, I., Shrivastava, A., Gupta, A., Hebert, M., 2016. Cross-stitch networks for multi-task learning, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).

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 abs/1905.08315.

Nguyen, A.T., Lim, S.N., Torr, P., 2022. Task-agnostic robust representation learning. arXiv preprint arXiv:2203.07596 .

Nwoye, C.I., Alapatt, D., Vardazaryan, Armine ... Gonzalez, C., Padoy, N., 2022a. Cholectriplet2021: a benchmark challenge for surgical action triplet recognition. arXiv preprint arXiv:2204.04746 .

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: Martel, A.L., Abolmaesumi, P., Stoyanov,

D., Mateus, D., Zuluaga, M.A., Zhou, S.K., Racoceanu, D., Joskowicz, L. (Eds.), Medical Image Computing and Computer Assisted Intervention – MICCAI 2020, Springer International Publishing, Cham. pp. 364–374.

Nwoye, C.I., Padoy, N., 2022. Data splits and metrics for benchmarking methods on surgical action triplet datasets. arXiv preprint arXiv:2204.05235 .

Nwoye, C.I., Yu, T., Gonzalez, C., Seeliger, B., Mascagni, P., Mutter, D., Marescaux, J., Padoy, N., 2022b. Rendezvous: Attention mechanisms for the recognition of surgical action triplets in endoscopic videos. Medical Image Analysis 78, 102433.

Ono, H., Yao, K., Fujishiro, M., Oda, I., Uedo, N., Nimura, S., Yahagi, N., Iishi, H., Oka, M., Ajioka, Y., et al., 2021. Guidelines for endoscopic submucosal dissection and endoscopic mucosal resection for early gastric cancer. Digestive Endoscopy 33, 4–20.

Padoy, N., Mateus, D., Weinland, D., Berger, M.O., Navab, N., 2009. Workflow monitoring based on 3D motion features, in: 2009 IEEE 12th International Conference on Computer Vision Workshops, ICCV Workshops, pp. 585–592.

Pakhomov, D., Navab, N., 2020. Searching for efficient architecture for instrument segmentation in robotic surgery, in: Medical Image Computing and Computer Assisted Intervention–MICCAI 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part III 23, Springer. pp. 648–656.

Paranjape, J.N., Nair, N.G., Sikder, S., Vedula, S.S., Patel, V.M., 2023. Adaptivesam: Towards efficient tuning of sam for surgical scene segmentation. arXiv preprint arXiv:2308.03726 .

Psychogyios, D., Colleoni, E., Van Amsterdam, B., Li, C.Y., Huang, S.Y., Li, Y., Jia, F., Zou, B., Wang, G., Liu, Y., et al., 2023. SAR-RARP50: Segmentation of surgical instrumentation and action recognition on robot-assisted radical prostatectomy challenge. arXiv preprint arXiv:2401.00496 .

Psychogyios, D., Mazomenos, E., Vasconcelos, F., Stoyanov, D., 2022. Msdesis: Multi-task stereo disparity estimation and surgical instrument segmentation. IEEE Transactions on Medical Imaging , 1–1.

Qin, F., Lin, S., Li, Y., Bly, R.A., Moe, K.S., Hannaford, B., 2020a. Towards better surgical instrument segmentation in endoscopic vision: Multi-angle feature aggregation and contour supervision. IEEE Robotics and Automation Letters 5, 6639–6646.

Qin, X., Zhang, Z., Huang, C., Dehghan, M., Zaiane, O., Jagersand, M., 2020b. U2-net: Going deeper with nested U-structure for salient object detection. Pattern Recognition 106, 107404.

Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al., 2021. Learning transferable visual models from natural language supervision, in: International conference on machine learning, PMLR. pp. 8748–8763.

Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al., 2019. Language models are unsupervised multitask learners. OpenAI blog 1, 9.

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 16, 1111–1119.

Rivas-Blanco, I., Pérez-Del-Pulgar, C.J., García-Morales, I., Muñoz, V.F., 2021. A review on deep learning in minimally invasive surgery. IEEE Access 9, 48658–48678. doi:10.1109/ACCESS.2021.3068852.

Ronneberger, O., Fischer, P., Brox, T., 2015. U-net: Convolutional networks for biomedical image segmentation, in: Navab, N., Hornegger, J., Wells, W.M., Frangi, A.F. (Eds.), Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015, Springer International Publishing, Cham. pp. 234–241.

Ross, T., Reinke, A., Full, P.M., Wagner, M., Kenngott, H., Apitz, M., Hempe, H., Filimon, D.M., Scholz, P., Tran, T.N., Bruno, P., Arbeláez, P., Bian, G.B., Bodenstedt, S., Bolmgren, J.L., Bravo-Sánchez, L., Chen, H.B., González, C., Guo, D., Halvorsen, P., Heng, P.A., Hosgor, E., Hou, Z.G., Isensee, F., Jha, D., Jiang, T., Jin, Y., Kirtac, K., Kletz, S., Leger, S., Li, Z., Maier-Hein, K.H., Ni, Z.L., Riegler, M.A., Schoeffmann, K., Shi, R., Speidel, S., Stenzel, M., Twick, I., Wang, G., Wang, J., Wang, L., Wang, L., Zhang, Y., Zhou, Y.J., Zhu, L., Wiesenfarth, M., Kopp-Schneider, A., Müller-Stich, B.P., Maier-Hein, L., 2020. Robust medical instrument segmentation challenge 2019. Eprint arXiv:2003.10299. arXiv:2003.10299.

Ruder, S., Bingel, J., Augenstein, I., Søgaard, A., 2019. Latent multi-task architecture learning, in: Proceedings of the AAAI Conference on Artificial Intelligence, pp. 4822–4829.

Sanchez-Matilla, R., Robu, M., Grammatikopoulou, M., Luengo, I., Stoyanov,

D., 2022. Data-centric multi-task surgical phase estimation with sparse scene segmentation. International Journal of Computer Assisted Radiology and Surgery 17, 953–960.

Sanchez-Matilla, R., Robu, M., Luengo, I., Stoyanov, D., 2021. Scalable joint detection and segmentation of surgical instruments with weak supervision, in: de Bruijne, M., Cattin, P.C., Cotin, S., Padoy, N., Speidel, S., Zheng, Y., Essert, C. (Eds.), Medical Image Computing and Computer Assisted Intervention – MICCAI 2021, Springer International Publishing, Cham. pp. 501–511.

Seenivasan, L., Islam, M., Kannan, G., Ren, H., 2023a. Surgicalgpt: End-to-end language-vision gpt for visual question answering in surgery, in: Greenspan, H., Madabhushi, A., Mousavi, P., Salcudean, S., Duncan, J., Syeda-Mahmood, T., Taylor, R. (Eds.), Medical Image Computing and Computer Assisted Intervention – MICCAI 2023, Springer Nature Switzerland, Cham. pp. 281–290.

Seenivasan, L., Islam, M., Krishna, A.K., Ren, H., 2022a. Surgical-VQA: Visual question answering in surgical scenes using transformer, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer. pp. 33–43.

Seenivasan, L., Islam, M., Xu, M., Lim, C.M., Ren, H., 2023b. Task-aware asynchronous multi-task model with class incremental contrastive learning for surgical scene understanding. International Journal of Computer Assisted Radiology and Surgery , 1–8.

Seenivasan, L., Mitheran, S., Islam, M., Ren, H., 2022b. Global-reasoned multi-task learning model for surgical scene understanding. IEEE Robotics and Automation Letters .

Sener, O., Koltun, V., 2018. Multi-task learning as multi-objective optimization, in: Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., Garnett, R. (Eds.), Advances in Neural Information Processing Systems 31, Curran Associates, Inc.. pp. 525–536.

Sharma, S., Nwoye, C.I., Mutter, D., Padoy, N., 2023a. Rendezvous in time: an attention-based temporal fusion approach for surgical triplet recognition. International Journal of Computer Assisted Radiology and Surgery , 1–7.

Sharma, S., Nwoye, C.I., Mutter, D., Padoy, N., 2023b. Surgical action triplet detection by mixed supervised learning of instrument-tissue interactions, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer. pp. 505–514.

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.

Simonyan, K., Zisserman, A., 2015. Very deep convolutional networks for large-scale image recognition, in: International Conference on Learning Representations.

Standley, T., Zamir, A., Chen, D., Guibas, L., Malik, J., Savarese, S., 2020. Which tasks should be learned together in multi-task learning?, in: International Conference on Machine Learning, PMLR. pp. 9120–9132.

Tarvainen, A., Valpola, H., 2017. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. Advances in neural information processing systems 30.

Teichmann, M., Weber, M., Zöllner, M., Cipolla, R., Urtasun, R., 2018. Multinet: Real-time joint semantic reasoning for autonomous driving, in: 2018 IEEE Intelligent Vehicles Symposium (IV), pp. 1013–1020.

Tran, D., Bourdev, L., Fergus, R., Torresani, L., Paluri, M., 2015. Learning spatiotemporal features with 3D convolutional networks, in: Proceedings of the IEEE international conference on computer vision, pp. 4489–4497.

Twinanda, A.P., Mutter, D., Marescaux, J., de Mathelin, M., Padoy, N., 2016. Single- and multi-task architectures for surgical workflow challenge at m2cai 2016. Eprint arXiv:1610.08844. `arXiv:1610.08844`.

Twinanda, A.P., Shehata, S., Mutter, D., Marescaux, J., de Mathelin, M., Padoy, N., 2017. Endonet: A deep architecture for recognition tasks on laparoscopic videos. IEEE transactions on medical imaging 36, 86–97.

Valderrama, N., Ruiz Puentes, P., Hernández, I., Ayobi, N., Verlyck, M., Santander, J., Caicedo, J., Fernández, N., Arbeláez, P., 2022. Towards holistic surgical scene understanding, in: International conference on medical image computing and computer-assisted intervention, Springer. pp. 442–452.

Vandenhende, S., Georgoulis, S., Gansbeke, W.V., Proesmans, M., Dai, D., Gool, L.V., 2022. Multi-task learning for dense prediction tasks: A survey. IEEE Transactions on Pattern Analysis & Machine Intelligence 44, 3614–3633.

Vercauteren, T., Unberath, M., Padoy, N., Navab, N., 2020. 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., Reid, C., 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., Zhao, L., Ge, Z., Sun, H., Xie, D., Guo, M., Liu, D., Kenngott, H.G., Nickel, F., von Frankenberg, M., Mathis-Ullrich, F., Kopp-Schneider, A., Maier-Hein, L., Speidel, S., Bodenstedt, S., 2023. Comparative validation of machine learning algorithms for surgical workflow and skill analysis with the heichole benchmark. Medical Image Analysis 86, 102770.

Wang, A., Islam, M., Xu, M., Zhang, Y., Ren, H., 2023. Sam meets robotic surgery: An empirical study in robustness perspective. arXiv preprint arXiv:2304.14674 .

Wang, J., Jin, Y., Cai, S., Xu, H., Heng, P.A., Qin, J., Wang, L., 2022a. Real-time landmark detection for precise endoscopic submucosal dissection via shape-aware relation network. Medical Image Analysis 75, 102291.

Wang, S., Raju, A., Huang, J., 2017. Deep learning based multi-label classification for surgical tool presence detection in laparoscopic videos, in: 2017 IEEE 14th international symposium on biomedical imaging (ISBI 2017), IEEE. pp. 620–623.

Wang, T., Wang, Y., Li, M., 2020. Towards accurate and interpretable surgical skill assessment: A video-based method incorporating recognized surgical gestures and skill levels, in: Martel, A.L., Abolmaesumi, P., Stoyanov, D., Mateus, D., Zuluaga, M.A., Zhou, S.K., Racoceanu, D., Joskowicz, L. (Eds.), Medical Image Computing and Computer Assisted Intervention – MICCAI 2020, Springer International Publishing, Cham. pp. 668–678.

Wang, Z., Lu, B., Long, Y., Zhong, F., Cheung, T.H., Dou, Q., Liu, Y., 2022b. Autolaparo: A new dataset of integrated multi-tasks for image-guided surgical automation in laparoscopic hysterectomy, in: Medical Image Computing and Computer Assisted Intervention–MICCAI 2022: 25th International Conference, Singapore, September 18–22, 2022, Proceedings, Part VII, Springer. pp. 486–496.

Wang, Z., Tsvetkov, Y., Firat, O., Cao, Y., 2021. Gradient vaccine: Investigating and improving multi-task optimization in massively multilingual models, in: International Conference on Learning Representations.

Ward, T.M., Mascagni, P., Ban, Y., Rosman, G., Padoy, N., Meireles, O., Hashimoto, D.A., 2021. Computer vision in surgery. Surgery 169, 1253–1256.

Wei, R., Li, B., Mo, H., Lu, B., Long, Y., Yang, B., Dou, Q., Liu, Y., Sun, D., 2022. Stereo dense scene reconstruction and accurate localization for learning-based navigation of laparoscope in minimally invasive surgery. IEEE Transactions on Biomedical Engineering 70, 488–500.

Xiao, B., Xu, W., Guo, J., Lam, H.K., Jia, G., Hong, W., Ren, H., 2020. Depth estimation of hard inclusions in soft tissue by autonomous robotic palpation using deep recurrent neural network. IEEE Transactions on Automation Science and Engineering 17, 1791–1799.

Xin, D., Ghorbani, B., Gilmer, J., Garg, A., Firat, O., 2022. Do current multi-task optimization methods in deep learning even help? Advances in Neural Information Processing Systems 35, 13597–13609.

Xu, D., Ouyang, W., Wang, X., Sebe, N., 2018. Pad-net: Multi-tasks guided prediction-and-distillation network for simultaneous depth estimation and scene parsing, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).

Xu, M., Islam, M., Lim, C.M., Ren, H., 2021. Class-incremental domain adaptation with smoothing and calibration for surgical report generation, in: Medical Image Computing and Computer Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part IV 24, Springer. pp. 269–278.

Yamlahi, A., Tran, T.N., Godau, P., Schellenberg, M., Michael, D., Smidt, F.H., Nölke, J.H., Adler, T.J., Tizabi, M.D., Nwoye, C.I., et al., 2023. Self-distillation for surgical action recognition, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer. pp. 637–646.

Yu, T., Kumar, S., Gupta, A., Levine, S., Hausman, K., Finn, C., 2020a. Gradient surgery for multi-task learning. Advances in Neural Information Processing Systems 33, 5824–5836.

Yu, T., Quillen, D., He, Z., Julian, R., Hausman, K., Finn, C., Levine, S., 2020b. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning, in: Conference on robot learning, PMLR. pp. 1094–1100.

Yue, W., Zhang, J., Hu, K., Xia, Y., Luo, J., Wang, Z., 2023. Surgical-sam: Efficient class promptable surgical instrument segmentation. Eprint

arXiv:2308.08746v1. `arXiv:2308.08746`.

Zamir, A.R., Sax, A., , Shen, W.B., Guibas, L., Malik, J., Savarese, S., 2018. Taskonomy: Disentangling task transfer learning, in: 2018 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), IEEE.

Zhang, Y., Yang, Q., 2022. A survey on multi-task learning. IEEE Transactions on Knowledge and Data Engineering 34, 5586–5609.

Zhang, Z., Cui, Z., Xu, C., Jie, Z., Li, X., Yang, J., 2018. Joint task-recursive learning for semantic segmentation and depth estimation, in: Proceedings of the European Conference on Computer Vision (ECCV).

Zhang, Z., Luo, P., Loy, C.C., Tang, X., 2014. Facial landmark detection by deep multi-task learning, in: Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part VI 13, Springer. pp. 94–108.

Zhao, Z., Jin, Y., Heng, P.A., 2022. TraSeTR: Track-to-segment transformer with contrastive query for instance-level instrument segmentation in robotic surgery, in: 2022 International Conference on Robotics and Automation (ICRA), IEEE. pp. 11186–11193.

Zhou, Z., Alabi, O., Wei, M., Vercauteren, T., Shi, M., 2023. Text promptable surgical instrument segmentation with vision-language models. arXiv preprint arXiv:2306.09244 .

Zhu, X., Su, W., Lu, L., Li, B., Wang, X., Dai, J., 2020. Deformable DETR: Deformable transformers for end-to-end object detection, in: International Conference on Learning Representations.