# SurgMLLMBench: A Multimodal Large Language Model Benchmark Dataset for Surgical Scene Understanding

Tae-Min Choi[1]   Tae Kyeong Jeong[2,*]   Garam Kim[2,*]   Jaemin Lee[3]   Yeongyoon Koh[4]

In Cheul Choi[4]   Jae-Ho Chung[3]   Jong Woong Park[4]   Juyoun Park[2,†]

[1]Samsung Research    [2]Center for Humanoid Research, Korea Institute of Science and Technology
[3]Department of plastic surgery, College of medicine, Korea University
[4]Department of orthopedic surgery, College of medicine, Korea University
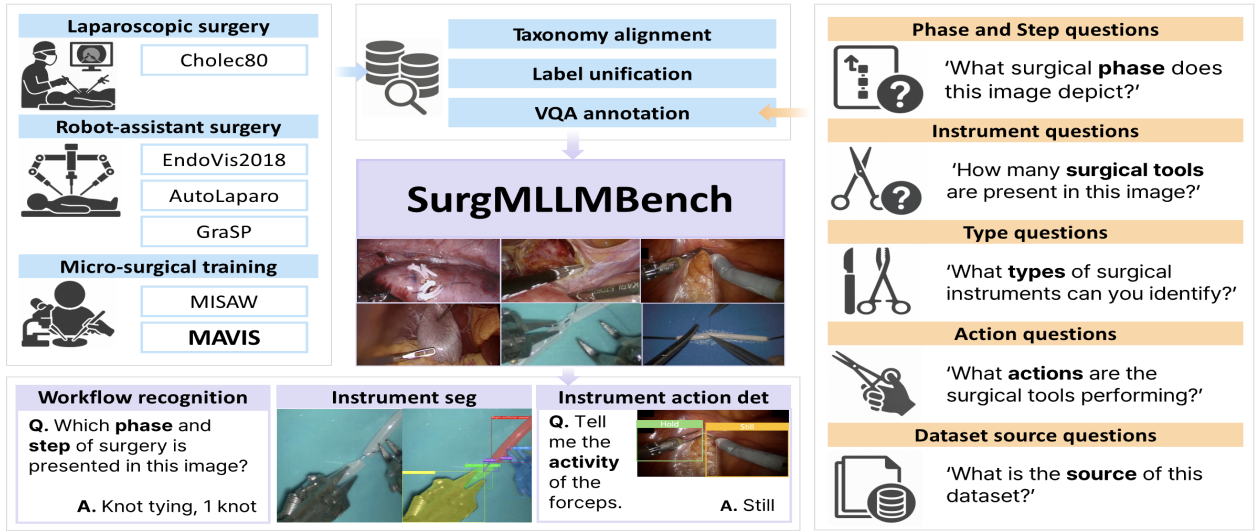
Figure 1. Overview of **SurgMLLMBench**. Multi-domain surgical datasets, including the newly collected **MAVIS**, are unified into a multimodal benchmark through taxonomy alignment, label unification, and VQA annotation. The template-based question generator (orange) produces structured VQA pairs using five query types. SurgMLLMBench supports interactive multimodal surgical scene understanding.

## Abstract

*Recent advances in multimodal large language models (LLMs) have highlighted their potential for medical and surgical applications. However, existing surgical datasets predominantly adopt a Visual Question Answering (VQA) format with heterogeneous taxonomies and lack support for pixel-level segmentation, limiting consistent evaluation and applicability. We present SurgMLLMBench, a unified multimodal benchmark explicitly designed for developing and evaluating interactive multimodal LLMs for surgical scene understanding, including the newly collected Microsurgical Artificial Vascular anastomosIS (MAVIS) dataset. It integrates pixel-level instrument segmentation masks and structured VQA annotations across laparoscopic, robot-assisted, and micro-surgical domains under a unified taxonomy, enabling comprehensive evaluation beyond traditional VQA tasks and richer visual–conversational interactions. Extensive baseline experiments show that a single model trained on SurgMLLMBench achieves consistent performance across domains and generalizes effectively to unseen datasets. SurgMLLMBench will be publicly released as a robust resource to advance multimodal surgical AI research, supporting reproducible evaluation and development of interactive surgical reasoning models.[1]*

## 1. Introduction

In minimally invasive procedures such as laparoscopic and robot-assisted surgery, surgeons depend almost entirely on the restricted, two-dimensional view provided by an endoscopic camera. Because the camera can rotate, zoom, or be

---

[1]http://surgmllmbench.github.io/

1

obscured by instruments, maintaining real-time awareness of the procedural stage and of each tool's precise action is challenging even for experienced surgeons [19, 23]. Multimodal large language models (LLMs) [18] offer a promising route to surgical assistance: by jointly interpreting images and text, a model could describe the ongoing step, highlight relevant instruments on the display, and answer intra-operative queries as they arise. Achieving this capability, however, requires training data that links rich workflow context with fine-grained visual localization [15, 28, 31].

Existing publicly available datasets have been developed primarily for Visual Question Answering (VQA) [5, 24, 29] tasks and therefore do not capture the full range of information needed for comprehensive scene understanding. EndoChat [26] and LLaVA-Surg [14] provide hundreds of thousands of question–answer pairs, yet omit annotations for the surgical phase and procedural steps, leaving the temporal structure of an operation unrepresented. Surgical-LLaVA [12] supplies phase labels, but its spatial supervision is limited to bounding-boxes rather than pixel-level delineations. The absence of detailed workflow annotations and fine-grained localization restricts how multimodal LLMs can be trained or evaluated as effective intra-operative assistants.

To address these limitations, we present **SurgMLLM-Bench**, a benchmark expressly designed around the full spectrum of surgical scene-understanding tasks, from global workflow recognition to fine-grained visual grounding. SurgMLLMBench integrates complementary annotations for the surgical phase, procedural step, instrument-centered actions, and pixel-level instrument segmentation, thereby covering an operation's temporal and spatial aspects. The dataset further embraces domain diversity by combining laparoscopic surgery, robot-assisted surgery, and micro-surgical training procedures with sequences captured in simulated environments, improving domain shift robustness. Also, unlike prior VQA-centric corpora that offer only bounding-box supervision, SurgMLLMBench provides dense segmentation masks, enabling multimodal LLMs to ground their textual outputs in pixel-accurate visual evidence. An overview of the proposed SurgMLLMBench is shown in Fig. 1.

The SurgMLLMBench benchmark unifies six surgical datasets—Cholec80 [25], EndoVis2018 [2], AutoLaparo [27], MISAW [10], GraSP [4], and a newly collected **Micro-surgical Artificial Vascular anastomosIS (MAVIS)** dataset. MAVIS captures complete surgical sequences in micro-surgical training sessions and introduces a hierarchical workflow taxonomy with a novel attribute, providing detailed representations of surgical workflows. Each frame across all datasets is organized around four complementary tasks: (1) phase recognition, assigning coarse workflow phases to capture global context; (2) step classification, refining this context by identifying the specific procedural step underway; (3) instrument-centered action detection, labeling the

functional primitive executed by each active tool (e.g., grasp, cut, suture) for fine-grained temporal reasoning; and (4) instrument segmentation, providing pixel-accurate masks for every visible instrument instance and offering the spatial detail required for precise visual grounding. By aligning all source videos to this standard set of tasks, SurgMLLMBench supplies a single supervision signal that spans the temporal structure and spatial layout of surgical scenes.

Because most current multimodal LLMs are limited to textual responses and coarse bounding-box outputs, we select OMG-LLaVA [30] as our baseline model. OMG-LLaVA combines a general-purpose segmentation encoder with a vision–language decoder, enabling the system to generate pixel-level masks while simultaneously producing natural-language answers to surgical queries. When trained on SurgMLLMBench, the model exhibits stable performance across datasets from multiple domains and demonstrates generalization ability to adapt even to datasets that were not included in the training. These findings indicate that the dense masks supplied by SurgMLLMBench not only encourage multimodal LLMs to link textual explanations with pixel-accurate visual evidence but also contribute to their robust understanding of surgical workflows across diverse scenarios. This paper contributes the SurgMLLMBench dataset, a reproducible integration pipeline that maps heterogeneous surgical resources into a coherent annotation scheme, and quantitative baselines that highlight both the benefits and the remaining challenges of applying interactive multimodal LLMs in surgical settings.

## 2. Related Work

### 2.1. Benchmarks for Surgical Scene Understanding

Recent advancements in surgical AI have accelerated the development of various benchmark datasets aimed at evaluating surgical scene understanding models comprehensively [13]. The GraSP [4] dataset addresses surgical scene understanding by providing hierarchical tasks including surgical phase recognition, procedural steps identification, and fine-grained tasks such as surgical instrument segmentation and atomic visual action detection specifically within robot-assisted prostatectomy scenarios. Additionally, the recently proposed MM-OR dataset [21] captures multimodal surgical environments by incorporating RGB-D video, audio data, speech transcripts, robot log data, and semantic scene graphs, facilitating tasks such as panoptic segmentation and holistic operating room (OR) scene understanding. Although these datasets significantly contribute to the surgical AI community by providing diverse annotations and tasks, they predominantly focus on specific surgical procedures or modalities, limiting their generalizability and application in broader surgical contexts. Moreover, many current benchmarks are confined to static question-answering frameworks, restrict-

ing their effectiveness in evaluating dynamic and interactive multimodal models for real-world surgical environments.

## 2.2. Multimodal LLMs in Surgical Applications

Multimodal LLMs have recently demonstrated significant promise across various medical and surgical applications by integrating visual, textual, and contextual data. EndoChat [26] utilizes the large-scale multimodal Surg-396K dataset to enable interactive surgical education, supporting complex interactions through combined visual and textual queries. LLaVA-Surg [14] introduces the Surg-QA dataset, generated from surgical lecture videos, to train a conversational vision-language assistant capable of answering open-ended queries regarding surgical scenarios. Similarly, Surgical-LLaVA [12] enhances multimodal interactions by integrating language models with visual encoders specifically trained for surgical spatiotemporal contexts, demonstrating improved performance in surgical video understanding tasks. Despite these advances, existing multimodal surgical LLMs largely remain restricted to predefined static interaction paradigms, lacking the ability to adapt dynamically to evolving surgical conditions. Particularly, tasks involving detailed pixel-level segmentation and real-time recognition of procedural contexts have been inadequately addressed due to the absence of richly annotated interactive datasets.

## 3. SurgMLLMBench Dataset

In this section, we introduce the SurgMLLMBench, a new benchmark recognizing the overall surgical procedure, scene information, and tool segmentation abilities for multimodal LLMs. We utilize five existing open-source surgical datasets and a newly generated dataset. An overview of the overall dataset composition and benchmark structure is illustrated in Fig. 1. To ensure compatibility across data sources, all six datasets were standardized under a unified annotation schema, and supplementary prompt-level annotations were introduced to support multimodal LLM training. The SurgMLLMBench dataset is publicly available at: huggingface.co/datasets/KIST-HARILAB/SurgMLLMBench.

### 3.1. Data Collection

**Public Datasets.** We construct an integrated benchmark dataset for training multimodal LLMs including five widely used surgical video datasets: Cholec80 [25], EndoVis2018 [2], AutoLaparo [27], GraSP [4], and MISAW [10]. Since the original MISAW dataset does not include instrument segmentation annotations, we incorporate supplementary segmentation data containing approximately 3,000 annotated images provided by [11]. To leverage the cross-domain characteristics in SurgMLLMBench, we included surgical video datasets covering diverse procedures and clinical domains.

We categorize the datasets according to the surgical environment, distinguishing them as Laparoscopic Surgery (LS), Robot-Assisted Surgery (RAS), or Micro-Surgical Training procedures (MST) following [4]. A detailed summary of this categorization is provided in Tab. 1.

Since the existing datasets differ in resolution, frame rate, annotation format, and task definitions, a standardization process was required to integrate all videos and annotations into a unified structure. First, all videos were converted into frame-level representations, and a COCO-style [16] metadata schema was designed to harmonize annotation information across datasets. Each frame includes the fields: video ID, frame ID, stage, phase, step, instrument action, and segmentation. Missing entries were left blank to maintain structural consistency. This integration process was carefully designed to harmonize annotation standards across different sources while maintaining spatiotemporal consistency and semantic interpretability. Based on this unified structure, all datasets were reorganized to enable training and evaluation under a single, consistent format. The range of annotations, video durations, and surgical domains provided by each dataset are summarized in Tab. 1, which illustrates the diversity and complexity of the tasks encompassed within SurgMLLMBench. The table visualizes the annotation coverage of each dataset under the unified SurgMLLMBench schema, demonstrating how heterogeneous datasets were standardized into a cohesive multimodal benchmark.

**MAVIS Dataset.** The MAVIS dataset is introduced as a comprehensive resource for micro-surgical workflow understanding. The proposed dataset targets the micro-surgery domain, where high-precision procedures demand prolonged focus from a few expert surgeons, emphasizing the role of AI and robotic assistance. Unlike MISAW [10], which includes only a single-step suturing task, our dataset captures the complete sequence of both anterior and posterior anastomosis of an artificial vessel within a single video. To enable a comprehensive understanding of the surgical workflow, we introduce a new attribute, *Stage*, allowing each procedure to be hierarchically annotated following a stage–phase–step structure. The MAVIS dataset is the first to incorporate the *Stage* attribute across surgical domains, providing a more detailed representation of real surgical workflows. The MAVIS dataset is available as a standalone dataset at: huggingface.co/datasets/KIST-HARILAB/MAVIS.

It comprises 19 videos of artificial vascular anastomosis procedures performed by three expert micro-surgeons from the College of Medicine, Korea University, who are also co-authors of this study. The dataset includes recordings of 1 mm artificial vessel anastomosis procedures, performed by pairs of surgeons (a lead surgeon and an assistant). Each pair executed the full anastomosis sequence in varying orders according to their preferred surgical techniques. Specifically, the dataset comprises seven videos from Surgeon 1, seven

| Dataset | Stage Recog. | Phase Recog. | Step Recog. | Instrument Action Recog. | Instrument Segmentation | Video Hour | Total Frames | Data Domain |
|---|---|---|---|---|---|---|---|---|
| Cholec80 [25] | | ✓ | | | | 51.25 | 184,498 | LS |
| EndoVis2018 [2] | | | ✓ | | ✓ | 1.58 | 2,235 | RAS |
| AutoLaparo [27] | | ✓ | | | ✓ | 23.13 | 83,243 | RAS |
| GraSP [4] | | ✓ | ✓ | ✓ | ✓ | 32.37 | 116,515 | RAS |
| MISAW‡ [10] | | ✓ | ✓ | ✓ | ✓ | 1.52 | 164,275 | MST |
| **MAVIS** | ✓ | ✓ | ✓ | | ✓ | 2.95 | 10,652 | MST |
| Total Annotations | 10,652 | 559,182 | 291,442 | 31,872 | 25,221 | 112.80 | 561,418 | |

Table 1. Comparison of datasets. A checkmark (✓) indicates the presence of corresponding annotations. (Recog. denotes recognition. MISAW‡ comprises the original MISAW dataset and supplementary segmentation annotations from [11].)
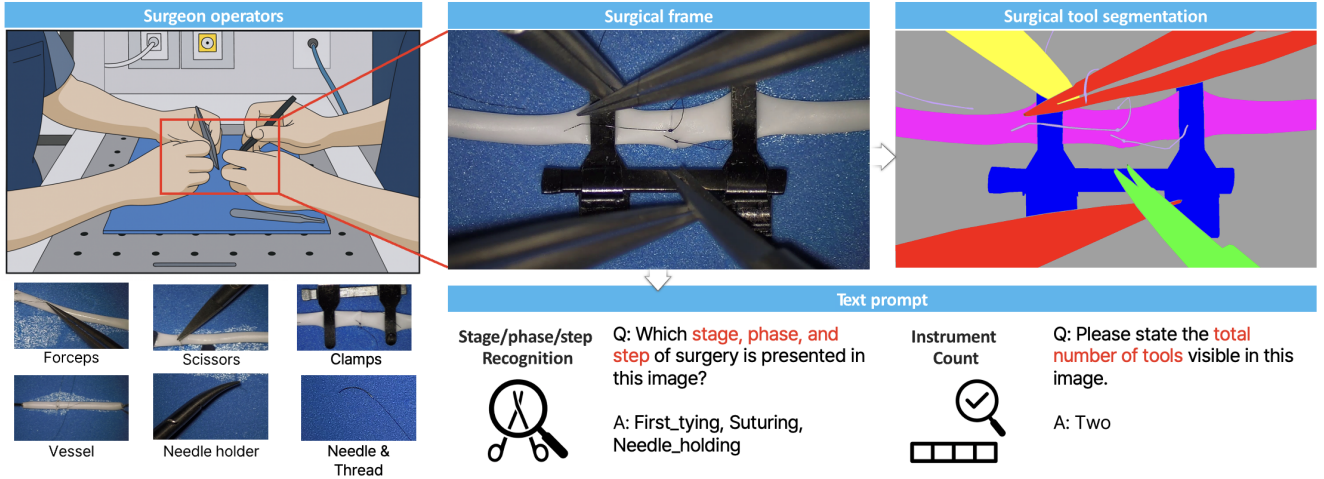


Figure 2. Overview of the MAVIS dataset collection process.

from Surgeon 2, and five from Surgeon 3. All videos were recorded at a resolution of 1920 × 1080 and temporally sampled at 1 FPS. For each frame, MAVIS provides pixel-level segmentation annotations for seven categories of surgical instruments and frame-level workflow annotations encompassing surgical stages, phases, and steps. The overall data collection setup and recording environment are illustrated in Fig. 2.

### 3.2. Data Annotation

**MAVIS Dataset.** The MAVIS dataset provides frame-level multi-level annotations designed for the precise analysis of micro-vascular anastomosis procedures. These annotations were constructed in collaboration with expert micro-surgeons (co-authors of this paper), reflecting real surgical workflows to ensure practical applicability in both clinical and research settings. Each frame contains workflow annotations that capture the temporal progression of surgery through a stage–phase–step hierarchy, as well as instrument segmentation annotations that represent the spatial components of the surgical scene. The surgical workflow is organized into six major stages: FIRST TYING, SECOND 180° TYING, SECOND 120° TYING, FRONT-SIDE TYING,

BACK-SIDE TYING, and FLIP. Each stage is subdivided into four phases—SUTURING, KNOT TYING, CUTTING, and FLIP—and each phase is further decomposed into fine-grained operational steps. The step annotations include eight categories: NEEDLE HOLDING, NEEDLE PASSING, NEEDLE DROPPING, KNOT TYING (1ST–3RD KNOT), CUTTING, and FLIPPING. This hierarchical annotation structure (see Fig. 3(d)) enables quantitative analysis of the temporal progression of surgical workflows and facilitates research on phase and step anticipation.

At the visual level, each frame includes pixel-wise segmentation annotations for all surgical instruments appearing in the scene. The annotated instrument categories comprise eight classes: BACKGROUND MATERIAL, FORCEPS, SCISSORS, VASCULAR CLAMPS, NEEDLE HOLDER, VESSEL, NEEDLE, and THREAD. Each object is delineated with polygonal coordinates following the COCO format, accurately capturing its spatial contour and extent. Multiple instruments may appear simultaneously within the same frame, and such pixel-level annotations can be utilized for various vision-based analyses, including instrument detection, segmentation, tracking, and action recognition [20].
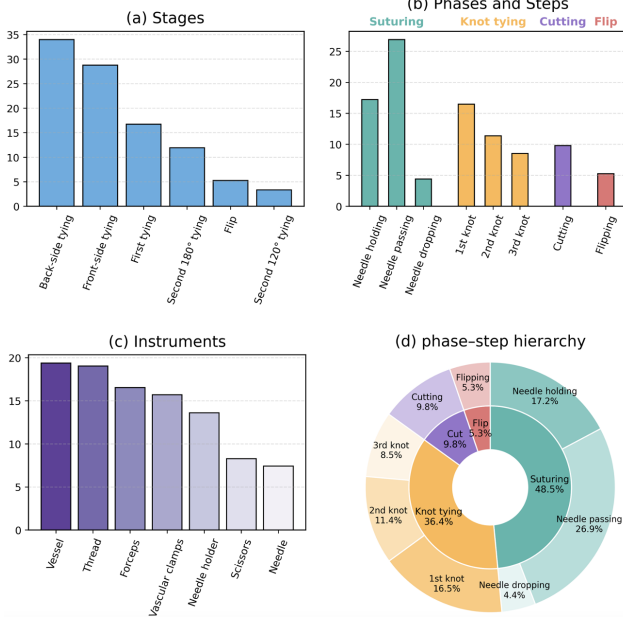
4

Figure 3. MAVIS dataset distributions.



(a) Dataset distribution



(b) Domain group distribution

Figure 4. Comparison of dataset composition across tasks.

**VQA Prompt Generation.** To construct the VQA annotations, we combined MAVIS frame-level workflow annotations (stage, phase, step) with short-term instrument metadata and, where available from other sources (e.g., EndoVis2018, GraSP, MISAW), instrument action labels. Each frame was paired with one of several fixed prompt templates, uniformly sampled across five categories: (1) *workflow queries* (e.g., "Which stage, phase, and step are shown?"), (2) *instrument count queries* (e.g., "How many surgical tools are visible?"), (3) *instrument type queries* (e.g., "Which instruments are present?"), (4) *instrument action queries* (e.g., "What action is the needle holder performing?"), and (5) *dataset source queries* (e.g., "What is the source of this dataset?"). Workflow and count prompts were always instantiated; type and action prompts were included when corresponding labels existed. The corresponding VQA prompts include dataset identifiers to help the model distinguish domain-specific contexts during training. A template-based VQA generation approach was adopted instead of using generative language models such as GPT [1]. Surgical scene understanding tasks, particularly workflow recognition and instrument counting, require high precision and deterministic supervision [17]. Free-form question generation can introduce linguistic variability and semantic inconsistency [6], whereas fixed prompt templates ensure consistent phrasing, reproducible dataset construction, and higher annotation accuracy, enabling stable and efficient multimodal training and evaluation [3] within SurgMLLMBench.

**MAVIS Dataset Statistics.** The MAVIS dataset comprises 10,652 frames, each annotated with stage, phase, and step
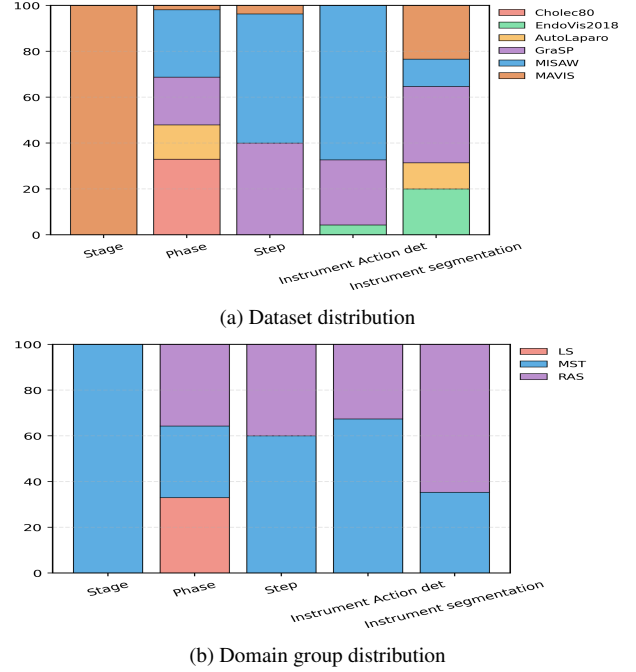
labels. Among the six surgical stages, BACK-SIDE TYING and FRONT-SIDE TYING occupy the largest proportions, followed by FIRST TYING and SECOND 180° TYING (see Fig. 3(a)). At the phase and step levels, SUTURING and KNOT TYING dominate the distribution, with NEEDLE PASSING and NEEDLE HOLDING as the most frequent fine-grained actions (see Fig. 3(b)). For instrument segmentation, VESSEL and THREAD appear most frequently, followed by FORCEPS, VASCULAR CLAMPS, and NEEDLE HOLDER, reflecting that suturing and knot tying constitute the core operations of micro-surgical anastomosis (see Fig. 3(c)). Overall, MAVIS presents a balanced composition of temporal workflow and visual annotations, quantitatively representing the procedural characteristics of micro-surgery. The dataset distributions are illustrated in Fig. 3.

**Dataset Comparison.** SurgMLLMBench integrates six surgical vision datasets: five existing datasets—Cholec80 [25], EndoVis2018 [2], AutoLaparo [27], GraSP [4], and MISAW [10]—plus a newly collected dataset, MAVIS. Fig. 4 visually summarizes the composition and task coverage of each dataset within the SurgMLLMBench framework. In (a), the proportion of data dedicated to each task (stage, phase, and step recognition; instrument action detection; instrument segmentation) is shown for each dataset, while (b) illustrates the distribution of data across surgical domains (LS, MST, RAS). Stage recognition was introduced in the proposed MAVIS dataset to represent hierarchical workflow information and is only available in this dataset among the six included. For the

5

remaining tasks, since phase recognition is only provided in Cholec80—the sole dataset in the LS domain—excluding this case, the datasets contribute roughly similar proportions of data across the different surgical domains.

## 4. Performance Evaluation

### 4.1. Evaluation Methods

**OMG-LLaVA.** We adopt OMG-LLaVA [30] as the primary interactive multimodal baseline because it can both answer open-ended surgical queries (e.g., phase, step, action, and count) and generate segmentation masks corresponding to these queries. This dual capability aligns closely with the design objectives of SurgMLLMBench, which provides workflow-level supervision alongside pixel-accurate instrument annotations across domains. Such integration enables joint learning of surgical reasoning and visual grounding within a unified multimodal framework.

Following OMG-LLaVA [30], we adopt a two-stage training framework comprising pre-training for broad image–text alignment and instruction tuning for multimodal task adaptation. In the pre-training stage, the perception model and LLM are frozen while only the visual and text projectors are optimized to establish vision–language alignment. Since OMG-LLaVA's pre-training stage already provides comprehensive multimodal alignment, we adopt the same pre-training dataset. During instruction tuning, we employ LoRA [9] to fine-tune the LLM and unfreeze the OMG-decoder to enable adaptation to surgical-specific visual features and pixel-level reasoning. The experiments use a batch size of 4 and one training epoch, while all other hyperparameters follow the original OMG-LLaVA configuration [30].

OMG-LLaVA is deliberately trained on SurgMLLM-Bench while excluding the MAVIS dataset. This experimental setup enables a clear assessment of cross-dataset generalization by evaluating how well the model adapts to MAVIS through additional fine-tuning, in comparison to models trained without instruction tuning.

**LLaVA.** We include LLaVA [18] as a text–vision baseline to quantify the effect of explicit pixel-level grounding on multimodal reasoning. Unlike OMG-LLaVA, LLaVA does not produce segmentation masks; instead, it processes image–text pairs to generate natural-language responses without a dedicated segmentation head. This makes it well-suited for recognition-oriented tasks such as phase, step, action, and count prediction. Similar to OMG-LLaVA, we follow the two-stage training paradigm of pre-training for general image–text alignment and instruction tuning for task adaptation. We initialize our models from the official LLaVA checkpoint and perform instruction tuning. Both experiments are conducted with a batch size of 16, one training epoch, and LoRA [9] for LLM fine-tuning. The learning rate is set to $2 \times 10^{-4}$ for per-dataset fine-tuning and $2 \times 10^{-5}$ for instruc-

tion tuning on SurgMLLMBench. All other hyperparameters follow the original LLaVA implementation [18]. LLaVA is trained on SurgMLLMBench with MAVIS withheld, in order to assess the generalization ability, after which additional evaluation experiments on MAVIS are conducted.

### 4.2. Evaluation Metrics

For tasks where both ground truth and predicted outputs are textual (phase, step, action, and count), we compute accuracy as the proportion of predictions whose text exactly matches the reference label [7, 22]. Formally, for $N$ samples with predictions $\hat{y}_i$ and ground truths $y_i$, accuracy is defined as $\text{Accuracy} = \frac{1}{N} \sum_{i=1}^{N} \mathbb{1}[\hat{y}_i = y_i]$, where $\mathbb{1}[\cdot]$ denotes the indicator function returning 1 when the prediction and ground truth are identical. This strict text-matching criterion ensures that the model's reasoning output aligns precisely with the ground truth workflow labels without relying on semantic similarity or partial credit.

For the segmentation task, we adopt the mean Intersection over Union (mIoU) metric [8], computed as the average IoU across all instrument classes. For a given class $c$, IoU is defined as $\text{IoU}_c = \frac{TP_c}{TP_c + FP_c + FN_c}$, where $TP_c$, $FP_c$, and $FN_c$ denote the number of true positive, false positive, and false negative pixels for class $c$, respectively. The mIoU is then calculated as $\text{mIoU} = \frac{1}{C} \sum_{c=1}^{C} \text{IoU}_c$, where $C$ is the total number of instrument classes in the dataset. This metric quantifies the pixel-level overlap between predicted and ground truth masks and allows for direct comparison among models with segmentation capability.

### 4.3. Evaluation Results

**SurgMLLMBench Instruction Tuning Performance.** Tab. 2 summarizes the performance of multimodal LLMs (LLaVA [18] and OMG-LLaVA [30]) on five datasets within the SurgMLLMBench. Each model is evaluated under two training configurations: (1) *fine-tuning on each dataset*, where the model is trained independently on each dataset, and (2) *instruction tuning on SurgMLLMBench*, where the model is trained once on the full SurgMLLMBench corpus excluding the MAVIS dataset and evaluated directly on all datasets without additional fine-tuning. This comparison highlights the generalization and scalability of a unified multimodal LLM trained across multiple surgical domains.

Tab. 2 demonstrates that a single model instruction-tuned on the full SurgMLLMBench corpus achieves competitive performance across diverse surgical datasets, despite not being optimized for any specific domain. In particular, for EndoVis2018 and AutoLaparo, the model trained on SurgM-LLMBench performs comparably to, or even surpasses, models fine-tuned individually on each dataset. Although some dataset-specific models exhibit slightly higher scores in their native domains, this is likely attributable to overfitting to dataset-specific label semantics or visual styles, and the per-

| Dataset | Method | Instruction Tuning on SurgMLLMBench | Fine-tuning on each dataset | Phase | Step | Action | Count | Instrument Segmentation |
|---|---|---|---|---|---|---|---|---|
| Cholec80 | LLaVA | ✗ | ✓ | 81.98 | | | 81.54 | |
| | LLaVA[§] | ✓ | ✗ | 77.32 | | | 78.73 | |
| | OMG-LLaVA | ✗ | ✓ | 76.51 | | | 83.08 | |
| | OMG-LLaVA[§] | ✓ | ✗ | 55.78 | | | 82.53 | |
| EndoVis2018 | LLaVA | ✗ | ✓ | | | 34.31 | 58.88 | - |
| | LLaVA[§] | ✓ | ✗ | | | 43.65 | 78.38 | - |
| | OMG-LLaVA | ✗ | ✓ | | | 45.47 | 60.63 | 26.04 |
| | OMG-LLaVA[§] | ✓ | ✗ | | | 44.53 | 63.09 | 27.23 |
| AutoLaparo | LLaVA | ✗ | ✓ | 23.00 | | | 57.21 | - |
| | LLaVA[§] | ✓ | ✗ | 24.67 | | | 82.56 | - |
| | OMG-LLaVA | ✗ | ✓ | 22.64 | | | 67.91 | 59.57 |
| | OMG-LLaVA[§] | ✓ | ✗ | 26.79 | | | 57.91 | 44.97 |
| MISAW | LLaVA | ✗ | ✓ | 87.16 | 64.03 | 48.58 | 94.17 | - |
| | LLaVA[§] | ✓ | ✗ | 88.20 | 58.69 | 19.73 | 94.70 | - |
| | OMG-LLaVA | ✗ | ✓ | 86.44 | 62.85 | 58.49 | 93.46 | 61.51 |
| | OMG-LLaVA[§] | ✓ | ✗ | 86.33 | 56.56 | 48.68 | 83.57 | 59.06 |
| GraSP | LLaVA | ✗ | ✓ | 68.83 | 52.75 | 39.95 | 54.40 | - |
| | LLaVA[§] | ✓ | ✗ | 66.78 | 50.83 | 2.95 | 47.56 | - |
| | OMG-LLaVA | ✗ | ✓ | 47.85 | 37.26 | 33.37 | 55.02 | 66.65 |
| | OMG-LLaVA[§] | ✓ | ✗ | 36.55 | 31.84 | 43.69 | 44.80 | 53.06 |

Table 2. Comparison of multimodal LLM performance across surgical tasks and datasets. Gray cells indicate the absence of corresponding annotations in each dataset. [§]Single instruction-tuned model per method (no additional fine-tuning).

| Method | Instruction Tuning on SurgMLLMBench | Fine-tuning on MAVIS | Stage | Phase | Step | Count | Instrument Segmentation |
|---|---|---|---|---|---|---|---|
| LLaVA | ✗ | ✓ | 67.67 | 65.09 | 37.59 | 54.67 | - |
| | ✓ | ✓ | 75.70 | 69.43 | 44.80 | 68.47 | - |
| OMG-LLaVA | ✗ | ✓ | 62.84 | 63.63 | 37.28 | 55.46 | 53.96 |
| | ✓ | ✓ | 43.90 | 55.72 | 29.48 | 47.51 | 36.89 |

Table 3. Results of multimodal LLMs on the MAVIS dataset, which was excluded from SurgMLLMBench instruction tuning to assess model adaptability to unseen surgical dataset.

formance differences remain modest. These results suggest that performing large-scale, cross-domain instruction tuning on SurgMLLMBench enables consistently strong performance across heterogeneous surgical environments, without incurring significant accuracy degradation. Notably, the findings highlight the effectiveness of SurgMLLMBench in facilitating the development of a single model that can perform interactive VQA across multiple surgical domains.

For GraSP and MISAW, LLaVa trained on SurgMLLM-Bench exhibits lower action prediction accuracy. This degradation stems from the imbalanced task distribution within the benchmark—although the number of action labels is large, the actual sample count per label is limited. As a result, the model tends to conflate semantically related actions across datasets (e.g., IDLE vs. STILL, RETRACTION vs. PULL), where functionally similar gestures are expressed under different naming conventions (see Fig. 5(c)). We argue that this imbalance leads to reduced numerical accuracy, even when the predicted actions remain contextually valid.

**Generalization Ability.** To further assess the generalization ability of instruction tuning on the entire SurgMLLM-Bench, we evaluate all models on the MAVIS dataset, which is intentionally excluded from the instruction tuning corpus (Tab. 3). This setup allows us to test whether models trained on SurgMLLMBench can adapt to unseen datasets and novel tasks, such as stage recognition. In this evaluation, the instruction-tuned model[§] from Tab. 2 serves as the initialization for additional fine-tuning on MAVIS and is compared against models trained solely on the MAVIS dataset. The checkpoint trained on SurgMLLMBench is further fine-tuned on MAVIS (3 epochs for LLaVA and 1 epoch for OMG-LLaVA), enabling evaluation of feature transfer to unseen datasets.

LLaVA benefits from SurgMLLMBench initialization, achieving consistent improvements across tasks and demonstrating strong adaptability even when transferred to a new dataset that introduces a novel task, stage recognition. These results indicate that performing large-scale cross-domain in-

Original Image | GT | OMG-LLaVA | OMG-LLaVA§
(a) Instrument segmentation comparison

| | | |
|---|---|---|
| User | | Could you specify the surgical phase and step shown here? |
| OMG-LLaVA | | Knot_Tying, 1_Knot |
| OMG−LLaVA§ | | Knot_Tying, 3_Knot |
| User | | What is the action of the right needle holder? |
| OMG-LLaVA | | Pull |
| OMG−LLaVA§ | | Make a Loop |

| | | |
|---|---|---|
| User | | Can you describe what the bipolar forceps is doing? |
| OMG-LLaVA | | Retraction |
| OMG−LLaVA§ | | Idle |
| User | | Tell me the activity of the monopolar curved scissors. |
| OMG-LLaVA | | Idle |
| OMG−LLaVA§ | | Cutting |

(b) Interactive workflow recognition by OMG-LLaVA

| | | |
|---|---|---|
| User | | Identify the current action of the bipolar forceps. |
| Ground Truth | | Hold, Pull |
| LLaVA | | Hold, Still |
| LLaVA§ | | Retraction |

| | | |
|---|---|---|
| User | | Tell me the activity of the right needle holder. |
| Ground Truth | | Insert |
| LLaVA | | Hold |
| LLaVA§ | | Sew |

(c) Interactive action detection by LLaVA
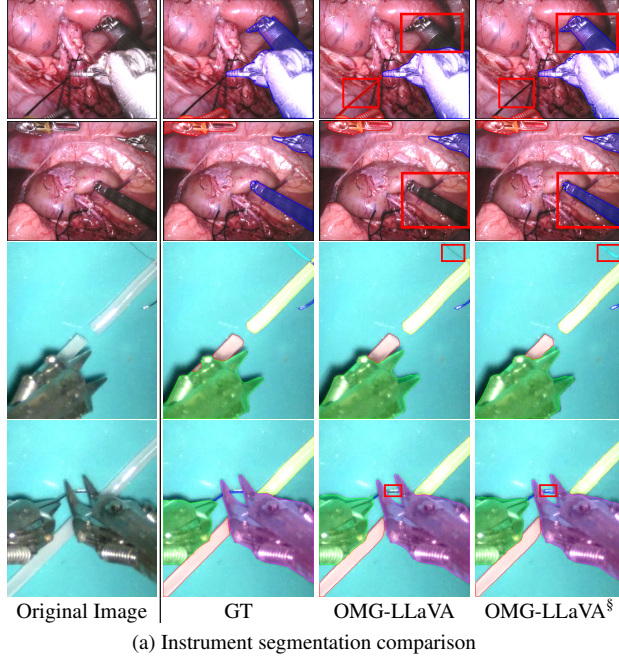(LLaVA§ provides an incorrect but contextually valid answer)

Figure 5. Qualitative visualization results of (a) instrument segmentation and (b, c) workflow recognition via VQA (green: correct, red: incorrect). OMG-LLaVA and LLaVA denote models trained individually on each dataset, whereas OMG-LLaVA§ and LLaVA§ represent a single model trained on SurgMLLMBench without additional per-dataset fine-tuning.

struction tuning on SurgMLLMBench enhances the ability of text–vision reasoning models to generalize to unseen tasks and datasets. For OMG-LLaVA, we observe that initializing from the model trained on SurgMLLMBench can result in performance degradation, likely due to visual gaps between datasets limiting the transferability of its pixel-level decoder.

**Visualization.** Fig. 5(a) presents qualitative results of OMG-LLaVA predictions under two training strategies: dataset-specific fine-tuning and SurgMLLMBench instruction tuning. The top two rows show results from the EndoVis2018 dataset [2], while the bottom two rows correspond to the MISAW dataset [10]. In the EndoVis2018 results, the model fine-tuned only on EndoVis2018 often misses existing instrument masks or produces spurious background predictions, suggesting overfitting to dataset-specific visual patterns. Conversely, the model trained on SurgMLLM-Bench—exposed to more diverse surgical scenes—shows more stable and conservative segmentation, reducing false detections and improving mask consistency. In MISAW, the SurgMLLMBench-trained model more accurately captures thin and small structures, such as needles and wires, compared to the MISAW-only fine-tuned model. This indicates that cross-domain instruction tuning encourages more generalized visual representations, improving recognition of delicate, low-contrast surgical tools.

Fig. 5(b) shows that OMG-LLaVA instruction-tuned on

SurgMLLMBench performs more accurate workflow recognition on MISAW and EndoVis2018 through VQA than the dataset-specific model, evidencing stronger workflow-aware visual reasoning. Conversely, Fig. 5(c) illustrates the LLaVA's action recognition result on MISAW and GraSP: the SurgMLLMBench-tuned model predicts the label from the Endovis2018 dataset (RETRACTION) instead of PULL, revealing cross-dataset semantic alignment that is contextually valid yet penalized by dataset-specific label names. Overall, these visual comparisons demonstrate that instruction tuning on the entire SurgMLLMBench improves both visual grounding and workflow reasoning, and achieves more consistent VQA predictions across tasks while maintaining semantic alignment even under cross-dataset label variations.

## 5. Conclusion

We introduced SurgMLLMBench, a unified benchmark for interactive multimodal LLMs that integrates workflow annotations (stage, phase, and step) with pixel-level instrument masks across multiple surgical domains, including the newly proposed MAVIS dataset. Baseline experiments demonstrate that training on SurgMLLMBench yields stable multimodal reasoning, enables pixel-accurate visual grounding, and improves robustness under domain shifts. Beyond serving as a benchmark, SurgMLLMBench provides a practical foundation for training multimodal surgical assistants that

can answer textual queries and provide visual evidence via segmentation masks. It also enables cross-domain generalization analysis and supports the development of agentic systems that maintain conversational context, leverage vision tools, and offer interpretable visual explanations—advancing applications in surgical education, intraoperative assistance, and robotic autonomy.

Future work will extend temporal and multimodal coverage (e.g., kinematics, audio, depth), broaden action taxonomies and instance tracking, introduce topology-aware metrics, and assess real-time reliability, uncertainty, and safety for clinical deployment.

# References

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 5

[2] Max Allan, Satoshi Kondo, Sebastian Bodenstedt, Stefan Leger, Rahim Kadkhodamohammadi, Imanol Luengo, Felix Fuentes, Evangello Flouty, Ahmed Mohammed, Marius Pedersen, et al. 2018 robotic scene segmentation challenge. *arXiv preprint arXiv:2001.11190*, 2020. 2, 3, 4, 5, 8

[3] Rabiul Awal, Le Zhang, and Aishwarya Agrawal. Investigating prompting techniques for zero-and few-shot visual question answering. *arXiv preprint arXiv:2306.09996*, 2023. 5

[4] Nicolás Ayobi, Santiago Rodríguez, Alejandra Pérez, Isabela Hernández, Nicolás Aparicio, Eugénie Dessevres, Sebastián Peña, Jessica Santander, Juan Ignacio Caicedo, Nicolás Fernández, et al. Pixel-wise recognition for holistic surgical scene understanding. *arXiv preprint arXiv:2401.11174*, 2024. 2, 3, 4, 5

[5] Long Bai, Guankun Wang, Mobarakol Islam, Lalithkumar Seenivasan, An Wang, and Hongliang Ren. Surgical-vqla++: Adversarial contrastive learning for calibrated robust visual question-localized answering in robotic surgery. *Information Fusion*, 113:102602, 2025. 2

[6] Wenjie Dong, Shuhao Shen, Yuqiang Han, Tao Tan, Jian Wu, and Hongxia Xu. Generative models in medical visual question answering: A survey. *Applied Sciences*, 15(6):2983, 2025. 5

[7] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010. 6

[8] Cristina González, Laura Bravo-Sánchez, and Pablo Arbelaez. Isinet: An instance-based approach for surgical instrument segmentation. In *Proceedings of the International conference on medical image computing and computer-assisted intervention*, pages 595–605. Springer, 2020. 6

[9] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. In *Proceedings of the International Conference on Learning Representations*, 2022. 6

[10] Arnaud Huaulmé, Duygu Sarikaya, Kévin Le Mut, Fabien Despinoy, Yonghao Long, Qi Dou, Chin-Boon Chng, Wenjun Lin, Satoshi Kondo, Laura Bravo-Sánchez, et al. Micro-surgical anastomose workflow recognition challenge report. *Computer Methods and Programs in Biomedicine*, 212:106452, 2021. 2, 3, 4, 5, 8

[11] Tae Kyeong Jeong, Garam Kim, and Juyoun Park. Micro-surgical instrument segmentation for robot-assisted surgery. *arXiv preprint arXiv:2509.11727*, 2025. 3, 4

[12] Juseong Jin and Chang Wook Jeong. Surgical-llava: Toward surgical scenario understanding via large language and vision models. *arXiv preprint arXiv:2410.09750*, 2024. 2, 3

[13] Ufaq Khan, Umair Nawaz, Adnan Qayyum, Shazad Ashraf, Muhammad Bilal, and Junaid Qadir. Surgical scene understanding in the era of foundation ai models: A comprehensive review. *arXiv preprint arXiv:2502.14886*, 2025. 2

[14] Jiajie Li, Garrett Skinner, Gene Yang, Brian R Quaranto, Steven D Schwaitzberg, Peter CW Kim, and Jinjun Xiong. Llava-surg: Towards multimodal surgical assistant via structured surgical video learning. *arXiv preprint arXiv:2408.07981*, 2024. 2, 3

[15] Pengpeng Li, Xiangbo Shu, Chun-Mei Feng, Yifei Feng, Wangmeng Zuo, and Jinhui Tang. Surgical video workflow analysis via visual-language learning. *npj Health Systems*, 2 (1):5, 2025. 2

[16] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Proceedings of the European conference on computer vision*, pages 740–755. Springer, 2014. 3

[17] Zhihong Lin, Donghao Zhang, Qingyi Tao, Danli Shi, Gholamreza Haffari, Qi Wu, Mingguang He, and Zongyuan Ge. Medical visual question answering: A survey. *Artificial Intelligence in Medicine*, 143:102611, 2023. 5

[18] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *Proceedings of the Advances in neural information processing systems*, pages 34892–34916, 2023. 2, 6

[19] Pietro Mascagni, Deepak Alapatt, Luca Sestini, Maria S Altieri, Amin Madani, Yusuke Watanabe, Adnan Alseidi, Jay A Redan, Sergio Alfieri, Guido Costamagna, et al. Computer vision in surgery: from potential to clinical value. *npj Digital Medicine*, 5(1):163, 2022. 2

[20] Shervin Minaee, Yuri Boykov, Fatih Porikli, Antonio Plaza, Nasser Kehtarnavaz, and Demetri Terzopoulos. Image segmentation using deep learning: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 44(7):3523–3542, 2021. 4

[21] Ege Özsoy, Chantal Pellegrini, Tobias Czempiel, Felix Tristram, Kun Yuan, David Bani-Harouni, Ulrich Eck, Benjamin Busam, Matthias Keicher, and Nassir Navab. Mm-or: A large multimodal operating room dataset for semantic understanding of high-intensity surgical environments. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 19378–19389, 2025. 2

[22] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*, 2016. 6

[23] Sanat Ramesh, Vinkle Srivastav, Deepak Alapatt, Tong Yu, Aditya Murali, Luca Sestini, Chinedu Innocent Nwoye, Idris Hamoud, Saurav Sharma, Antoine Fleurentin, et al. Dissecting self-supervised learning methods for surgical computer vision. *Medical Image Analysis*, 88:102844, 2023. 2

[24] Lalithkumar Seenivasan, Mobarakol Islam, Adithya K Krishna, and Hongliang Ren. Surgical-vqa: Visual question answering in surgical scenes using transformer. In *Proceedings of the International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 33–43. Springer, 2022. 2

[25] Andru P Twinanda, Sherif Shehata, Didier Mutter, Jacques Marescaux, Michel De Mathelin, and Nicolas Padoy. Endonet: A deep architecture for recognition tasks on laparoscopic videos. *IEEE transactions on medical imaging*, 36(1):86–97, 2016. 2, 3, 4, 5

[26] Guankun Wang, Long Bai, Junyi Wang, Kun Yuan, Zhen Li, Tianxu Jiang, Xiting He, Jinlin Wu, Zhen Chen, Zhen Lei, et al. Endochat: Grounded multimodal large language model for endoscopic surgery. *arXiv preprint arXiv:2501.11347*, 2025. 2, 3

[27] Ziyi Wang, Bo Lu, Yonghao Long, Fangxun Zhong, Tak-Hong Cheung, Qi Dou, and Yunhui Liu. Autolaparo: A new dataset of integrated multi-tasks for image-guided surgical automation in laparoscopic hysterectomy. In *Proceedings of the International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 486–496. Springer, 2022. 2, 3, 4, 5

[28] Mengya Xu, Zhongzhen Huang, Jie Zhang, Xiaofan Zhang, and Qi Dou. Surgical action planning with large language models. In *Proceedings of the International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 563–572. Springer, 2025. 2

[29] Kun Yuan, Manasi Kattel, Joël L Lavanchy, Nassir Navab, Vinkle Srivastav, and Nicolas Padoy. Advancing surgical vqa with scene graph knowledge. *International Journal of Computer Assisted Radiology and Surgery*, 19(7):1409–1417, 2024. 2

[30] Tao Zhang, Xiangtai Li, Hao Fei, Haobo Yuan, Shengqiong Wu, Shunping Ji, Chen Change Loy, and Shuicheng Yan. Omg-llava: Bridging image-level, object-level, pixel-level reasoning and understanding. In *Proceedings of the Advances in Neural Information Processing Systems*, pages 71737–71767, 2024. 2, 6

[31] Zhihong Zhu, Yunyan Zhang, Xuxin Cheng, Zhiqi Huang, Derong Xu, Xian Wu, and Yefeng Zheng. Alignment before awareness: Towards visual question localized-answering in robotic surgery via optimal transport and answer semantics. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 711–721, 2024. 2