

Supplementary Material

Kun Yuan^{1,3}, Vinkle Srivastav^{1,2}, Nassir Navab³, and Nicolas Padoy^{1,2}

¹ University of Strasbourg, CNRS, INSERM, ICube, UMR7357, Strasbourg, France

² IHU Strasbourg, Strasbourg, France

³ CAMP, Technische Universität München, Munich, Germany

$$\mathcal{L}_{clip}^{aws} = -\frac{1}{B} \sum_{i=1}^B \log \left(\frac{\exp(\mathcal{F}_v(v_i)^T \cdot \mathcal{F}_t(a_i)/\tau)}{\sum_{j=1}^B \exp(\mathcal{F}_v(v_j)^T \cdot \mathcal{F}_t(a_j)/\tau)} \right)$$

$$\mathcal{L}_{clip}^{whisper} = -\frac{1}{B} \sum_{i=1}^B \log \left(\frac{\sum_{m=1}^M \exp(\mathcal{F}_v(v_i)^T \cdot \mathcal{F}_t(w_i^m)/\tau)}{\sum_{j=1}^B \sum_{m=1}^M \exp(\mathcal{F}_v(v_j)^T \cdot \mathcal{F}_t(w_j^m)/\tau)} \right)$$

$$L_{clip} = L_{clip}^{aws} + L_{clip}^{whisper}$$

Fig. 1: Loss function of clip-level video-text pairs in Sec. 2.3. a : narration text (n_{ij} in the main paper) transcribed from AWS Medical Transcribe [1]. w : narration text transcribed from Whisper [2] system. v : short-term video clips, which correspond to 1 aws text and M whisper texts.

Table 1: Manually designed prompts for the class names to recognize the surgical phase in Cholec80 dataset.

Phase Labels	Prompts
<i>Preparation</i>	In preparation phase I insert trocars to patient abdomen cavity
<i>CalotTriangleDissection</i>	In calot triangle dissection phase I use grasper to hold gallbladder and use hook to expose the hepatic triangle area and cystic duct and cystic artery
<i>ClippingCutting</i>	In clip and cut phase I use clipper to clip the cystic duct and artery then use scissor to cut them
<i>GallbladderDissection</i>	In dissection phase I use the hook to dissect the connective tissue between gallbladder and liver
<i>GallbladderPacking</i>	In packaging phase I put the gallbladder into the specimen bag
<i>CleaningCoagulation</i>	In clean and coagulation phase I use suction and irrigation to clear the surgical field and coagulate bleeding vessels
<i>GallbladderRetraction</i>	In retraction phase I grasp the specimen bag and remove it from trocar

Table 2: Manually designed prompts for the class names to recognize the surgical phase in AutoLaparo dataset.

Phase Labels	Prompts
<i>Preparation</i>	I use grasper to grasp and explore the field
<i>Dividing Ligament and Peritoneum</i>	I divide ligament and peritoneum
<i>Dividing Uterine Vessels and Ligament</i>	I divide uterine vessels and ligament
<i>Transecting the Vagina</i>	I use the dissecting hook to transect the vagina
<i>Specimen Removal</i>	I remove the specimen bag and uterus
<i>Suturing</i>	I suture the tissue
<i>Washing</i>	Washing

Table 3: Manually designed prompts for the class names to recognize the surgical phase in gastric bypass dataset. We use the same prompts for both StrasBypass70 and BernBypass70.

Phase Labels	Prompts
<i>Preparation</i>	In preparation phase I insert trocars to the abdominal cavity and expose of the operating field
<i>Gastric pouch creation</i>	I cut the fat tissue and open retrogastric window at stomach
<i>Omentum division</i>	I grasp and lift the omentum and divide it
<i>Gastrojejunal anastomosis</i>	I see the proximal jejunum and determine the length of the biliary limb. I open the distal jejunum and create the gastrojejunostomy using a stapler. I reinforcement of the gastrojejunostomy with an additional suture.
<i>Anastomosis test</i>	I place the retractor and move the gastric tube and detect any leakage of the gastrojejunostomy
<i>Jejunal separation</i>	I open the mesentery to facilitate the introduction of the stapler and transect the jejunum proximal
<i>Petersen space closure</i>	I expose between the alimentary limb and the transverse colon and close it with sutures
<i>Jejunojejunal anastomosis</i>	I expose between the alimentary limb and the transverse colon and close it with sutures
<i>Mesenteric defect closure</i>	I expose the mesenteric defect and then close it by stitches
<i>Cleaning and coagulation</i>	In clean and coagulation phase I use suction and irrigation to clear the surgical field and coagulate bleeding vessels
<i>Disassembling</i>	I remove the instruments, retractor, ports, and camera
<i>Other</i>	Other

References

1. AWS: Amazon transcribe medical (2023), <https://aws.amazon.com/transcribe/medical/>
2. Radford, A., Kim, J.W., Xu, T., Brockman, G., McLeavey, C., Sutskever, I.: Robust speech recognition via large-scale weak supervision. In: International Conference on Machine Learning. pp. 28492–28518. PMLR (2023)