# arXiv:2311.12603v2 [cs.CV] 22 Nov 2023

# Surgical Temporal Action-aware Network with Sequence Regularization for Phase Recognition

Zhen Chen<sup>1\*</sup>, Yuhao Zhai<sup>2\*</sup>, Jun Zhang<sup>2⊠</sup>, Jinqiao Wang<sup>1,3,4,5⊠</sup>

<sup>1</sup>Centre for Artificial Intelligence and Robotics, Hong Kong Institute of Science&Innovation, Chinese Academy of Sciences

<sup>2</sup>Beijing Friendship Hospital, Capital Medical University <sup>3</sup>Institute of Automation, Chinese Academy of Sciences <sup>4</sup>Wuhan AI Research <sup>5</sup>ObjectEye Inc.

zhen.chen@cair-cas.org.hk, zhaiyuhao@ccmu.edu.cn, zhangjun5986@ccmu.edu.cn, jqwang@nlpr.ia.ac.cn

Abstract—To assist surgeons in the operating theatre, surgical phase recognition is critical for developing computer-assisted surgical systems, which requires comprehensive understanding of surgical videos. Although existing studies made great progress, there are still two significant limitations worthy of improvement. First, due to the compromise of resource consumption, framewise visual features are extracted by 2D networks and disregard spatial and temporal knowledge of surgical actions, which hinders subsequent inter-frame modeling for phase prediction. Second, these works simply utilize ordinary classification loss with onehot phase labels to optimize the phase predictions, and cannot fully explore surgical videos under inadequate supervision. To overcome these two limitations, we propose a Surgical Temporal Action-aware Network with sequence Regularization, named STAR-Net, to recognize surgical phases more accurately from input videos. Specifically, we propose an efficient multi-scale surgical temporal action (MS-STA) module, which integrates visual features with spatial and temporal knowledge of surgical actions at the cost of 2D networks. Moreover, we devise the dualclassifier sequence regularization (DSR) to facilitate the training of STAR-Net by the sequence guidance of an auxiliary classifier with a smaller capacity. Our STAR-Net with MS-STA and DSR can exploit visual features of surgical actions with effective regularization, thereby leading to the superior performance of surgical phase recognition. Extensive experiments on a large-scale gastrectomy surgery dataset and the public Cholec80 benchmark prove that our STAR-Net significantly outperforms state-of-thearts of surgical phase recognition.

Index Terms-video analysis, surgery workflow, gastric cancer

### I. INTRODUCTION

The computer-assisted surgery can improve the quality of interventional healthcare, thereby facilitating patient safety [1]–[3]. In particular, surgical phase recognition [4] is significant for developing systems to monitor surgical procedures [5], schedule surgeons [6], promote surgical team coordination [7], and educate junior surgeons [8]. Compared with offline analysis of surgical videos, online recognition can support decision-making during surgery without using future frames, which is more practical in surgical applications.

Z. Chen and Y. Zhai contribute equally to this work.

Online phase recognition of surgical videos is challenging, and has received great research attention and progress [9]-[11]. Earlier works [12] formulated this task as the frameby-frame classification, and used auxiliary annotations of surgical tools for multi-task learning [13]. Meanwhile, some works [14]-[16] utilized 3D convolutions to capture temporal knowledge of surgical videos. To overcome the huge resource consumption of 3D convolution, mainstream methods [17]-[21] first used 2D convolutional neural networks (CNNs) to extract the feature vector of each surgical video frame, and then predicted the surgical phase with the inter-frame temporal relationship aggregated by the long short-term memory (LSTM) [17], temporal convolutions [18], [22], or transformers [19]. On this basis, recent works [20], [21] further improved this multi-stage paradigm of phase recognition by leveraging longrange temporal relation among frame-wise feature vectors.

However, existing works [17]–[21] on surgical phase recognition suffer from two major limitations, including the insufficient visual information of frame-wise feature vectors, and the inadequate supervision knowledge provided by surgical phase labels. First, most surgical workflow studies [17]-[19] first extracted frame-wise feature vectors with 2D networks, and then aggregated these feature vectors for surgical phase prediction. Note that the spatial and temporal information of surgical videos is discarded when 2D networks process frames into feature vectors, thus hindering the subsequent inter-frame modeling. To overcome this bottleneck, we aim to efficiently formulate the surgical actions during feature extraction and provide visual features with spatial and temporal information for sequence modeling and phase prediction. Second, existing works [17]-[21] formulated the phase prediction as a classification task of the current frame, and the supervision information provided by the ordinary loss with one-hot phase labels is inadequate, which makes the training susceptible to over-fitting. To guarantee that networks fully learn surgical knowledge as possible, it is beneficial to conduct reasonable regularization in training. Inspired by this idea, we introduce an auxiliary classifier with a smaller capacity to regularize the phase prediction of the input video sequence.

To address these two problems in surgical phase recognition, we propose a Surgical Temporal Action-aware Network with sequence Regularization, named STAR-Net, from the

This work is supported by National Key R&D Program of China (No. 2022ZD0160601), National Natural Science Foundation of China (No. 62276260, 61976210, 62076235, 62176254, 62006230, 62206283), Beijing Municipal Science & Technology Commission (No. D17100006517003), and InnoHK program.



Fig. 1. (a) The overview of the STAR-Net, (b) multi-scale surgical temporal action (MS-STA), (c) temporal difference (TDiff) operation, and (d) dual-classifier sequence regularization (DSR). The MS-STA module is inserted into the 2D visual backbone, which progressively conducts TDiff operations to efficiently capture multi-scale surgical action features. The DSR introduces the mutual regularization between the auxiliary classifier and the task classifier at the early and late sequence respectively.

perspective of feature extraction and surgical supervision. Specifically, we first devise an efficient Multi-Scale Surgical Temporal Action (MS-STA) module and insert it into the visual extraction network, which enables the visual features to perceive the surgical actions at the computational cost of 2D networks. In particular, we progressively conduct Temporal Difference (TDiff) operations to capture multi-scale surgical action features for MS-STA. Moreover, we devise the Dualclassifier Sequence Regularization (DSR) to regularize the training of STAR-Net by introducing an auxiliary classifier with a smaller capacity. As such, this auxiliary classifier regularizes the task classifier at the early sequence to prevent over-fitting, and the task classifier provided with spatial and temporal knowledge enhances the auxiliary classifier at the late sequence in turn. With the proposed MS-STA and DSR, our STAR-Net can exploit visual features with the knowledge of surgical actions and learn from abundant surgical supervision, thereby leading to the superior performance of surgical phase recognition. We perform extensive experiments on a largescale gastrectomy surgery dataset and the public Cholec80 benchmark to validate the effectiveness of our STAR-Net, which outperforms state-of-the-art surgical phase recognition methods by a large margin.

## II. METHODOLOGY

## A. Overview

As illustrated in Fig. 1 (a), our STAR-Net predicts the phase of each frame in surgical videos to achieve online phase

recognition. Following previous studies [20], our STAR-Net classifies the current frame  $x_n$  as one of C surgical phases by taking the current frame and T-1 preceding frames as sequence input  $\{x_{n-t}\}_{t=0}^{T-1}$ . By progressively shifting the sequence input over time, the STAR-Net can predict the surgical phase of each frame in the entire video.

Specifically, the STAR-Net first utilizes a 2D CNN with the MS-STA module as the backbone to extract visual features with spatial and temporal information of surgical actions. Then, a transformer with spatial and temporal attention blocks efficiently aggregates visual features by exploiting global relationships in spatial and temporal dimensions sequentially. Finally, we introduce the DSR with an auxiliary classifier to mutually regularize sequence predictions produced by the task classifier, thereby facilitating the training of the STAR-Net.

### B. Multi-Scale Surgical Temporal Action for Visual Features

Existing studies [18], [19] extracted frame-wise visual information into feature vectors, which lost spatial and temporal information of surgical videos. As a result, the surgical actions in surgical videos are not well represented, thereby leading to inaccurate modeling of the inter-frame relation. To address this problem, we propose the MS-STA module to efficiently model the multi-scale surgical temporal actions during visual extraction of the 2D backbone, which provides visual features with spatial and temporal knowledge for STAR-Net.

As shown in Fig. 1 (b), the MS-STA integrates visual features  $\boldsymbol{f} \in \mathbb{R}^{T \times H \times W \times D}$  of video sequences with multi-

scale temporal information of surgical actions to facilitate surgical phase recognition, where T is the length of the input sequence, and H, W and D are the numbers of height, width and channel dimensions of visual features. In particular, we devise the Temporal Difference (TDiff) operation to capture surgical actions between two adjacent frames, which can be used for longer range surgical actions based on previous operations progressively. In Fig. 1 (c), the input visual features f of TDiff operation are first shifted along the temporal dimension for one frame as delayed features  $\mathcal{D}(f, 1)$ , where the first and the last frame is performed with zero-padding and truncation, respectively. Then, we subtract the delayed features  $\mathcal{D}(f, 1)$  from the input visual features f elementwise to calculate the surgical action features of each frame relative to the previous adjacent frame, as follows:

$$\boldsymbol{a}_1 = \mathcal{M}(\boldsymbol{f} - \mathcal{D}(\boldsymbol{f}, 1)), \tag{1}$$

where action mask  $\mathcal{M}(\cdot)$  sets the first frame substraction to 0. Note that the TDiff operation efficiently captures surgical action features for each frame with only one shift operation and element-wise subtraction. In this way, we obtain surgical action features  $a_1$  and delayed features  $\mathcal{D}(f, 1)$  as the output of TDiff operation.

With the delayed features, the MS-STA can further perform the TDiff operation to progressively generate the action features with a longer temporal range, *e.g.*,  $\mathcal{D}(\boldsymbol{f}, 2) = \mathcal{D}(\mathcal{D}(\boldsymbol{f}, 1), 1)$ . By conducting multiple TDiff operations sequentially in Fig. 1 (b), we concatenate these surgical action features  $\{\boldsymbol{a}_k\}_{k=1}^{\tau}$  with multiple temporal scales, where  $[\boldsymbol{a}_1, \boldsymbol{a}_2, \cdots, \boldsymbol{a}_{\tau}] \in \mathbb{R}^{T \times \tau \times H \times W \times D}$  and  $\tau$  denotes the number of temporal scales, and then perform a 3D convolution to integrate the multi-scale temporal features of surgical actions, as follows:

$$\boldsymbol{a}_{\mathrm{ms}} = \boldsymbol{W} \circledast [\boldsymbol{a}_1, \boldsymbol{a}_2, \cdots, \boldsymbol{a}_{\tau}], \qquad (2)$$

where  $a_{ms} \in \mathbb{R}^{T \times H \times W \times D}$ , W is the parameters of a 3D convolutional layer and  $\circledast$  is the convolution operation. In contrast to the burdensome 3D convolutional networks, we only insert one 3D convolutional layer into the STAR-Net to integrate multi-scale temporal features of surgical actions, which perceive the surgical actions at the computational cost of 2D networks.

Finally, we add the multi-scale surgical action features  $a_{\rm ms}$  with the input features f as residual learning, which can provide each frame with the knowledge of surgical actions for the surgical phase recognition. Different from TSM [23] that shifted partial channels for temporal information at different layers, our MS-STA can efficiently capture multi-scale temporal information of surgical actions at once, while preserving the channel alignment of visual features, thereby providing surgical action features for phase recognition.

### C. Dual-Classifier Sequence Regularization

With multi-scale surgical action features provided by MS-STA, the STAR-Net can predict the surgical phase with discriminative spatial and temporal features. However, existing



Fig. 2. Typical examples of eight phases in gastrectomy phase dataset. Each surgical phase carries a distinct and specific clinical significance and serves as the necessary procedure of the gastrectomy.

works [12], [18], [20], [21] employed ordinary classification loss, *e.g.*, the cross-entropy loss and its variants, to train the network, which cannot provide sufficient supervision for the training. Since the phase label y is a one-hot vector to indicate the correct class, the cross-entropy loss  $L_{CE} = -\sum_{c=1}^{C} y_c \log p_c$  merely produces a single non-zero constraint among these C terms to supervise the network training. As a result, the lack of supervision makes the network prone to over-fitting, and thus restricts the performance of surgical phase recognition.

To address the lack of supervision, we devise the Dualclassifier Sequence Regularization (DSR) to regularize sequence predictions by introducing a frame-wise auxiliary classifier, as illustrated in Fig. 1 (d). With the tokens of each frame provided by the transformer in STAR-Net, the task classifier can generate frame-wise phase predictions of the input video sequence, where the predicted probabilities are denoted as  $p_{task}$ . Meanwhile, the auxiliary classifier uses the sequence features extracted by the 2D visual backbone, and performs spatial global average pooling to predict the phase probabilities  $p_{aux}$  of each frame.

Since MS-STA provides multi-scale temporal information of surgical actions, the auxiliary classifier can achieve relatively satisfactory prediction for each video frame. Considering that the small number of previous frames in the early sequences E cannot provide sufficient temporal knowledge for the task classifier after the transformer, we adopt the auxiliary classifier with a smaller capacity to regularize the predicted probabilities  $p_{\mathrm{task}}$  of the task classifier. This provides effective regularization for the training of STAR-Net, thereby avoiding over-fitting. On the other hand, due to the lack of long-range surgical video knowledge, the auxiliary classifier is inferior to the task classifier on the late sequences L, and thus we further improve the auxiliary classifier with the task classifier. In turn, this can promote the learning of the task classifier with an improved auxiliary classifier. Therefore, the objective of our DSR is summarized as follows:

$$L_{\text{DSR}} = \sum_{i \in \boldsymbol{E}} \text{KL}(\boldsymbol{p}_{\text{task}}^{(i)} || \hat{\boldsymbol{p}}_{\text{aux}}^{(i)}) + \sum_{j \in \boldsymbol{L}} \text{KL}(\boldsymbol{p}_{\text{aux}}^{(j)} || \hat{\boldsymbol{p}}_{\text{task}}^{(j)}),$$
(3)

where KL is the Kullback–Leibler divergence to measure the distance between two probabilities, and  $\hat{p}$  represents stopping the gradients from p by regarding it as constants. Therefore,

 TABLE I

 Comparison with state-of-the-arts on gastrectomy phase dataset. Best and second best results are highlighted and <u>underlined</u>.

Method	AC (%)	PR (%)	RE (%)	JA (%)	P-value
PhaseNet [12]	$72.9_{\pm 7.2}$	$66.5_{\pm 17.6}$	$70.4_{\pm 5.1}$	$52.2 \pm 13.6$	$2.8 \times 10^{-17}$
SV-RCNet [17]	$84.3_{\pm 7.6}$	$79.8_{\pm 9.4}$	$78.9_{\pm 7.9}$	$66.1_{\pm 10.6}$	$1.4 \times 10^{-12}$
TeCNO [18]	$85.4_{\pm 7.1}$	$80.9_{\pm 9.5}$	$80.3_{\pm 7.7}$	$68.0_{\pm 10.9}$	$1.3 \times 10^{-9}$
TMRNet [20]	$86.8_{\pm 6.2}$	$85.1_{\pm 7.1}$	$81.9_{\pm 8.3}$	$71.8_{\pm 9.3}$	$2.7 \times 10^{-7}$
Trans-SVNet [21]	$87.7_{\pm 6.0}$	$85.1 \pm 6.7$	$82.0_{\pm 8.4}$	$71.9_{\pm 9.4}$	$2.5 \times 10^{-6}$
STAR-Net w/o MS-STA, DSR	$85.1 \pm 6.3$	$80.5_{\pm 11.6}$	$81.5_{\pm 5.0}$	$68.6_{\pm 10.8}$	$1.6 \times 10^{-8}$
STAR-Net w/o MS-STA	$86.8_{\pm 6.2}$	$81.8_{\pm 10.3}$	$82.7 \pm 6.1$	$69.6_{\pm 10.2}$	$2.6 \times 10^{-8}$
STAR-Net w/o DSR	$87.9 \pm 6.9$	$85.5 \pm 7.1$	$82.3_{\pm 8.5}$	$72.0 \pm 9.1$	$3.0 \times 10^{-6}$
STAR-Net	<b>89.2</b> ±6.1	$86.6_{\pm 6.4}$	$83.7_{\pm 8.1}$	$73.5_{\pm 9.0}$	-



Fig. 3. The proportion of eight phases in gastrectomy phase dataset. The inherent imbalance of surgical phases makes online recognition challenging.

the first term in Eq. (3) optimizes  $p_{task}$  on the early sequences E, while the second term optimizes  $p_{aux}$  on the late sequences L. In this way, the DSR can facilitate the training of STAR-Net with the sequence regularization between the task classifier and the auxiliary classifier.

# D. Training and Inference

Following the efficient multi-stage training paradigm in previous works [20], [21], we first train the 2D visual backbone with MS-STA using the cross-entropy loss  $L_{CE}$ , and generate frame features with spatial and temporal knowledge. Then, we train the transformer with the task and auxiliary classifiers under DSR for surgical phase recognition, as follows:

$$L = L_{\rm CE} + \lambda L_{\rm DSR},\tag{4}$$

where the coefficient  $\lambda$  controls the trade-off between  $L_{\rm DSR}$ and the cross-entropy loss  $L_{\rm CE}$  of phase predictions. In the inference, the well-trained STAR-Net sequentially conducts the 2D visual backbone with MS-STA and the transformer with spatial and temporal attention blocks to extract visual features, and performs the online frame-wise prediction using the task classifier for the surgical video streaming in an endto-end manner.

### III. EXPERIMENT

### A. Dataset and Implementation Details

1) Gastrectomy Phase Dataset: To evaluate the online phase recognition of surgical videos, we collect a large-scale laparoscopic gastrectomy dataset consisting of 100 surgical videos from different gastric cancer patients, and its data size



Fig. 4. Color-coded ribbon comparison of PhaseNet, Trans-SVNet, STAR-Net and ground truth.

is 22.1 times<sup>1</sup> of the Cholec80 dataset [12]. The surgical videos are recorded with  $1,920 \times 1,080$  resolution and 25 frame-per-second (fps). The average length of surgical videos is 2.53 hours. All surgical videos are annotated by two surgeons with expertise in gastric cancer surgery. Each frame of surgical videos is assigned to one out of eight surgical phases, including the preparation, the greater curvature separation, the distal stomach separation, the lesser curvature separation, the gastrointestinal (GI) tract reconstruction, and the operation ending. We randomly split the dataset at the patient-level, as 70 videos for training and 30 videos for test.

To elaborate the collected gastrectomy phase dataset for surgical phase recognition, we show typical examples of eight phases in the gastrectomy surgery in Fig. 2. It is evident that each of these surgical phases carries distinct and specific clinical significance, and together these phases constitute the entire procedures of gastrectomy. Moreover, the proportion of eight phases is illustrated in Fig. 3. It is worth noting that the inherent imbalance of these eight phases makes it more difficult to accurately achieve the online phase recognition.

2) Cholec80 Dataset: We further perform the comparison on public Cholec80 dataset [12] of laparoscopic cholecystectomy procedures, which contains 80 surgical videos with a resolution of  $854 \times 480$  or  $1,920 \times 1,080$  at 25 fps. The surgery

<sup>&</sup>lt;sup>1</sup>The size of the dataset is measured in the number of pixels.

### TABLE II

COMPARISON WITH STATE-OF-THE-ARTS ON CHOLEC80 DATASET. BEST AND SECOND BEST RESULTS ARE highlighted and <u>underlined</u>.

Method	AC (%)	PR (%)	RE (%)	JA (%)	Param (107)	FLOPs (10 <sup>10</sup> )
PhaseNet [12]	$78.8 \pm 4.7$	$71.3_{\pm 15.6}$	$76.6_{\pm 16.6}$	-	4.23	0.07
SV-RCNet [17]	$85.3 \pm 7.3$	$80.7_{\pm 7.0}$	$83.5_{\pm 7.5}$	-	2.88	4.14
UATD [24]	$88.6 \pm 6.7$	$86.1 \pm 6.7$	$88.0 \pm 10.1$	$73.7_{\pm 10.2}$	2.80	5.72
TeCNO [18]	$88.6_{\pm 7.8}$	$86.5_{\pm 7.0}$	$87.6_{\pm 6.7}$	$75.1_{\pm 6.9}$	2.36	8.29
MTRCNet-CL [13]	$89.2_{\pm 7.6}$	$86.9_{\pm 4.3}$	$88.0_{\pm 6.9}$	-	2.98	4.14
TMRNet [20]	$89.2 \pm 9.4$	$89.7_{\pm 3.5}$	<b>89.5</b> +4.8	$78.9 \pm 5.8$	6.30	24.86
Trans-SVNet [21]	$90.3 \pm 7.1$	$90.7 \pm 5.0$	$88.8 {\pm} 7.4$	$79.3 \pm 6.6$	2.37	12.47
STAR-Net	$\overline{91.2}_{\pm 5.3}$	$\overline{\textbf{91.6}}_{\pm 3.4}$	$89.2 \pm 9.4$	$\overline{79.5}_{\pm 8.1}$	1.68	3.92



Fig. 5. Visualization of surgical action features  $a_{ms}$  of MS-STA in (a) gastrectomy and (b) Cholec80 datasets. The motion of the ultrasound knife, grasper and hook is captured in MS-STA, which provides spatial and temporal information for phase recognition.

procedures are divided into seven surgical phases, including the preparation, the calot triangle dissection, the clipping and cutting, the gallbladder dissection, the gallbladder packaging, the cleaning and coagulation, and the gallbladder retraction. We exactly follow the standard splits [12], [20], *i.e.*, the first 40 videos for training and the rest 40 videos for test.

3) Implementation Details: We compare STAR-Net with state-of-the-arts using PyTorch [25] on a single NVIDIA A100 GPU. In our STAR-Net, we adopt ResNet-18 [26] as the 2D visual backbone for feature extraction, and implement the temporal attention block with causal mask [19] to achieve online recognition without using future frames. For MS-STA, the temporal scale  $\tau$  is set as 5, and the sequence length T is 20. The coefficient  $\lambda$  of  $L_{\text{DSR}}$  is set as 1.0, and E and L are set as the 20% – 60% and 80% – 100% ranges of input video sequences, respectively. All models are optimized in SGD with the batch size of 32. The learning rate is initialized as  $1 \times 10^{-3}$  and halved after every 5 epochs.

4) Evaluation Metrics: We adopt four commonly-used metrics to comprehensively evaluate the performance of surgical phase recognition, including accuracy (AC), precision (PR), recall (RE) and Jaccard (JA). Following the evaluation protocol in previous works [12], [20], we calculate PR, RE and JA in the phase-wise manner, and report the average and standard deviation. The AC represents the percentage of frames correctly classified into ground truth. To perform fair comparisons, the selected state-of-the-art methods are evaluated with the same criteria as the STAR-Net. Note that all experiments are performed in the online mode, where future information is not accessible when estimating the current frame.

### B. Comparison on Gastrectomy Dataset

1) Comparison with state-of-the-arts: To verify the effectiveness of our STAR-Net, we perform a comprehensive comparison with the state-of-the-art methods [12], [17], [18], [20], [21]. As illustrated in Table I, our STAR-Net achieves the best performance among these methods, with the AC of 89.2% and JA of 73.5%. Noticeably, our STAR-Net outperforms the transformer-based method, Trans-SVNet [21], by a large margin, *e.g.*, 1.5% in AC and 1.6% in JA. In addition, we conduct the t-test of AC among paired test videos, which confirms a significant advantage of our STAR-Net over [12], [17], [18], [20], [21] with P-values  $< 1 \times 10^{-5}$ . These results demonstrate the performance advantage of our STAR-Net over state-of-the-arts on gastrectomy phase recognition.

2) Ablation Study: As elaborated in Table I, we perform the detailed ablation study to validate the effectiveness, by implementing three ablative baselines of STAR-Net without MS-STA or DSR. Compared with the baseline without both MS-STA and DSR, the MS-STA can bring an improvement of 2.8% in AC and 3.4% in JA, which reveals the impact of surgical actions on the task. Meanwhile, the DSR can also increase the baseline with 1.7% in AC, which validates the sequence regularization of the auxiliary classifier benefits the training of STAR-Net. The ablation experiments indicate that the proposed MS-STA and DSR are crucial to improving the performance of STAR-Net on surgical phase recognition.

*3) Qualitative Results of Phase Recognition:* We further qualitatively compare our STAR-Net with Trans-SVNet [21] and PhaseNet [12] by presenting the color-coded ribbon results on gastrectomy dataset. As shown in Fig. 4, our STAR-Net

outperforms both PhaseNet [12] and Trans-SVNet [21], and is the closest to ground truth. In this way, these qualitative results confirm the superiority of our STAR-Net in surgical video analysis.

### C. Comparison on Cholec80 Dataset

To further evaluate the performance of phase recognition, we perform the comparison with more state-of-the-arts [13], [24] on the public Cholec80 benchmark in terms of performance and efficiency. In Table II, our STAR-Net achieves the overwhelming performance with the best AC of 91.2%, PR of 91.6% and JA of 79.5%. Furthermore, our STAR-Net demonstrates superior efficiency in comparison to existing algorithms [13], [17], [18], [20], [21], [24] with the minimal parameters and computation except for the frame-wise 2D CNN [12]. These competitive experimental results confirm the superiority of our STAR-Net on surgical phase recognition.

### D. Qualitative Analysis of Surgical Temporal Action

To analyze the surgical temporal action, we further visualize the multi-scale action features  $a_{\rm ms}$  of MS-STA, as shown in Fig. 5. Compared with the current frame, the MS-STA can accurately capture the surgical actions from several previous frames, where multi-scale action features  $a_{\rm ms}$  highlight the instrument motions on gastrectomy and Cholec80 datasets. For example, the motion of the ultrasound knife, grasper and hook is perceived by the multi-scale action features of MS-STA in Fig. 5. In this way, the MS-STA provides visual features with the spatial and temporal information of surgical actions for the STAR-Net, thereby facilitating the phase recognition tasks.

### IV. CONCLUSION

In this work, we propose the STAR-Net to promote online surgical phase recognition efficiently. Specifically, we first devise the MS-STA module to integrate the visual features with the multi-scale temporal knowledge of surgical actions, which enables the STAR-Net to process the surgical video sequence with more abundant surgical information. Moreover, we introduce the DSR to regularize the training of STAR-Net over the frame prediction of video sequences using an auxiliary classifier. Extensive experiments on gastrectomy and cholecystectomy surgical datasets confirm the remarkable advantages of our STAR-Net over state-of-the-art works in terms of performance and efficiency, as well as the perception of surgical temporal actions.

### REFERENCES

- L. Maier-Hein, M. Eisenmann, D. Sarikaya, K. März, T. Collins, A. Malpani, J. Fallert, H. Feussner, S. Giannarou, P. Mascagni *et al.*, "Surgical data science–from concepts toward clinical translation," *Med. Image Anal.*, vol. 76, p. 102306, 2022.
- [2] Z. Chen, Q. Guo, L. K. Yeung, D. T. Chan, Z. Lei, H. Liu, and J. Wang, "Surgical video captioning with mutual-modal concept alignment," in *MICCAI*. Springer, 2023, pp. 24–34.
- [3] Y. Zhai, Z. Chen, Z. Zheng, X. Wang, X. Yan, X. Liu, J. Yin, J. Wang, and J. Zhang, "Artificial intelligence for automatic surgical phase recognition of laparoscopic gastrectomy in gastric cancer," *Int. J. Comput. Assist. Radiol. Surg.*, 2023.

- [4] C. R. Garrow, K.-F. Kowalewski, L. Li, M. Wagner, M. W. Schmidt, S. Engelhardt, D. A. Hashimoto, H. G. Kenngott, S. Bodenstedt, S. Speidel *et al.*, "Machine learning for surgical phase recognition: a systematic review," *Annals of surgery*, vol. 273, no. 4, pp. 684–693, 2021.
- [5] S. S. Panesar, M. Kliot, R. Parrish, J. Fernandez-Miranda, Y. Cagle, and G. W. Britz, "Promises and perils of artificial intelligence in neurosurgery," *Neurosurgery*, vol. 87, no. 1, pp. 33–44, 2020.
- [6] Z. A. Abdalkareem, A. Amir, M. A. Al-Betar, P. Ekhan, and A. I. Hammouri, "Healthcare scheduling in optimization context: a review," *Health and Technology*, vol. 11, pp. 445–469, 2021.
- [7] L. R. Kennedy-Metz, P. Mascagni, A. Torralba, R. D. Dias, P. Perona, J. A. Shah, N. Padoy, and M. A. Zenati, "Computer vision in the operating room: Opportunities and caveats," *IEEE Trans. Med. Robot. Bionics*, vol. 3, no. 1, pp. 2–10, 2020.
- [8] A. Kirubarajan, D. Young, S. Khan, N. Crasto, M. Sobel, and D. Sussman, "Artificial intelligence and surgical education: A systematic scoping review of interventions," *Journal of Surgical Education*, vol. 79, no. 2, pp. 500–515, 2022.
- [9] F. Yi and T. Jiang, "Hard frame detection and online mapping for surgical phase recognition," in *MICCAI*. Springer, 2019, pp. 449–457.
- [10] Y. Zhang, S. Bano, A.-S. Page, J. Deprest, D. Stoyanov, and F. Vasconcelos, "Retrieval of surgical phase transitions using reinforcement learning," in *MICCAI*. Springer, 2022, pp. 497–506.
- [11] X. Ding, Z. Liu, and X. Li, "Free lunch for surgical video understanding by distilling self-supervisions," in *MICCAI*. Cham: Springer Nature Switzerland, 2022, pp. 365–375.
- [12] A. P. Twinanda, S. Shehata, D. Mutter, J. Marescaux, M. De Mathelin, and N. Padoy, "Endonet: a deep architecture for recognition tasks on laparoscopic videos," *IEEE Trans. Med. Imaging*, vol. 36, no. 1, pp. 86–97, 2016.
- [13] Y. Jin, H. Li, Q. Dou, H. Chen, J. Qin, C.-W. Fu, and P.-A. Heng, "Multitask recurrent convolutional network with correlation loss for surgical video analysis," *Med. Image Anal.*, vol. 59, p. 101572, 2020.
- [14] I. Funke, S. Bodenstedt, F. Oehme, F. von Bechtolsheim, J. Weitz, and S. Speidel, "Using 3d convolutional neural networks to learn spatiotemporal features for automatic surgical gesture recognition in video," in *MICCAI*. Springer, 2019, pp. 467–475.
- [15] B. Zhang, A. Ghanem, A. Simes, H. Choi, A. Yoo, and A. Min, "Swnet: Surgical workflow recognition with deep convolutional network," in *Medical imaging with deep learning*. PMLR, 2021, pp. 855–869.
- [16] B. Zhang, A. Ghanem, A. Simes, H. Choi, and A. Yoo, "Surgical workflow recognition with 3dcnn for sleeve gastrectomy," *Int. J. Comput. Assist. Radiol. Surg.*, vol. 16, no. 11, pp. 2029–2036, 2021.
- [17] Y. Jin, Q. Dou, H. Chen, L. Yu, J. Qin, C.-W. Fu, and P.-A. Heng, "Sv-rcnet: workflow recognition from surgical videos using recurrent convolutional network," *IEEE Trans. Med. Imaging*, vol. 37, no. 5, pp. 1114–1126, 2017.
- [18] T. Czempiel, M. Paschali, M. Keicher, W. Simson, H. Feussner, S. T. Kim, and N. Navab, "Tecno: Surgical phase recognition with multistage temporal convolutional networks," in *MICCAI*. Springer, 2020, pp. 343–352.
- [19] T. Czempiel, M. Paschali, D. Ostler, S. T. Kim, B. Busam, and N. Navab, "Opera: Attention-regularized transformers for surgical phase recognition," in *MICCAI*. Springer, 2021, pp. 604–614.
- [20] Y. Jin, Y. Long, C. Chen, Z. Zhao, Q. Dou, and P.-A. Heng, "Temporal memory relation network for workflow recognition from surgical video," *IEEE Trans. Med. Imaging*, vol. 40, no. 7, pp. 1911–1923, 2021.
- [21] X. Gao, Y. Jin, Y. Long, Q. Dou, and P.-A. Heng, "Trans-svnet: accurate phase recognition from surgical videos via hybrid embedding aggregation transformer," in *MICCAI*. Springer, 2021, pp. 593–603.
- [22] Y. A. Farha and J. Gall, "Ms-tcn: Multi-stage temporal convolutional network for action segmentation," in CVPR, 2019, pp. 3575–3584.
- [23] J. Lin, C. Gan, and S. Han, "Tsm: Temporal shift module for efficient video understanding," in *ICCV*, 2019, pp. 7083–7093.
- [24] X. Ding, X. Yan, Z. Wang, W. Zhao, J. Zhuang, X. Xu, and X. Li, "Less is more: Surgical phase recognition from timestamp supervision," *IEEE Trans. Med. Imaging*, 2022.
- [25] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, "Pytorch: An imperative style, high-performance deep learning library," *arXiv preprint arXiv*:1912.01703, 2019.
- [26] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in CVPR, 2016, pp. 770–778.