# Rethinking Surgical Captioning: End-to-End Window-Based MLP Transformer Using Patches

Mengya Xu[1,2,3*], Mobarakol Islam[4*], and Hongliang Ren[1,2,3**]

[1] Dept. of Electronic Engineering, The Chinese University of Hong Kong, Hong Kong SAR, China
[2] Dept. of Biomedical Engineering, National University of Singapore, Singapore
[3] National University of Singapore (Suzhou) Research Institute (NUSRI), China
[4] BioMedIA Group, Dept. of Computing, Imperial College London, London, UK
`mengya@u.nus.edu, m.islam20@imperial.ac.uk, hlren@ee.cuhk.edu.hk`

**Abstract.** Surgical captioning plays an important role in surgical instruction prediction and report generation. However, the majority of captioning models still rely on the heavy computational object detector or feature extractor to extract regional features. In addition, the detection model requires additional bounding box annotation which is costly and needs skilled annotators. These lead to inference delay and limit the captioning model to deploy in real-time robotic surgery. For this purpose, we design an end-to-end detector and feature extractor-free captioning model by utilizing the patch-based shifted window technique. We propose **S**hifted **Win**dow-Based **M**ulti-**L**ayer **P**erceptrons **Trans**former **Cap**tioning model (SwinMLP-TranCAP) with faster inference speed and less computation. SwinMLP-TranCAP replaces the multi-head attention module with window-based multi-head MLP. Such deployments primarily focus on image understanding tasks, but very few works investigate the caption generation task. SwinMLP-TranCAP is also extended into a video version for video captioning tasks using 3D patches and windows. Compared with previous detector-based or feature extractor-based models, our models greatly simplify the architecture design while maintaining performance on two surgical datasets. The code is publicly available at `https://github.com/XuMengyaAmy/SwinMLP_TranCAP`.

## 1 Introduction

Automatic surgical captioning is a prerequisite for intra-operative context-aware surgical instruction prediction and post-operative surgical report generation. Despite the impressive performance, most approaches require heavy computational resources for the surgical captioning task, which limits the real-time deployment. The main-stream captioning models contain an expensive pipeline of detection and feature extraction modules before captioning. For example, Meshed-Memory Transformer [4], self-sequence captioning [13] entail bounding box from detection

---

[*] Equal technical contribution.
[**] Corresponding author.

model (Faster R-CNN [12]) to extract object feature using a feature extractor of ResNet-101 [7]. On the other hand, bounding box annotations are used to extract regional features for surgical report generation [17, 18]. However, these kinds of approaches arise following issues: (i) require bounding box annotation which is challenging in the medical scene as it requires the use of professional annotators, (ii) region detection operation leads to high computational demand, and (iii) cropping regions may ignore some crucial background information and destroy the spatial correlation among the objects. Thus recent vision-language studies [20] are moving toward the detector-free trend by only utilizing image representations from the feature extractor. Nonetheless, the feature extractor still exists as an intermediate module which unavoidably leads to inadequate training and long inference delay at the prediction stage. These issues restrict the application for real-time deployment, especially in robotic surgery.

To achieve end-to-end captioning framework, ViTCAP model [6] uses the Vision Transformer (ViT) [5] which encodes image patches as grid representations. However, it is very computing intensive even for a reasonably large-sized image because the model has to compute the self-attention for a given patch with all the other patches in the input image. Most recently, the window-based model Swin Transformer [10] outperforms ViT [5] and reduces the computation cost by performing local self-attention between patches within the window. The window is also shifted to achieve information sharing across various spatial locations. However, this kind of approach is yet to explore for captioning tasks.

So far, many Transformer-variants are still based on the common belief that the attention-based token mixer module contributes most to their performance. However, recent studies [15] show that the attention-based module in Transformer can be replaced by spatial multi-layer perceptron (MLPs) and the resulting models still obtain quite good performance. This observation suggests that the general architecture of the transformer, instead of the specific attention-based token mixer module, is more crucial to the success of the model. Based on this hypothesis, [19] replaces the attention-based module with the extremely simple spatial pooling operator and surprisingly finds that the model achieves competitive performance. However, such exploration is mainly concentrated on the image understanding tasks and remains less investigated for the caption generation task.

**Our contributions can be summed up as the following points:** 1) We design SwinMLP-TranCAP, a detector and feature extractor-free captioning models by utilizing the shifted window-based MLP as the backbone of vision encoder and a Transformer-like decoder; 2) We also further develop a Video SwinMLP-TranCAP by using 3D patch embedding and windowing to achieve the video captioning task; 3) Our method is validated on publicly available two captioning datasets and obtained superior performance in both computational speed and caption generation over conventional approaches; 4) Our findings suggest that the captioning performance mostly relies on transformer type architecture instead of self-attention mechanism.
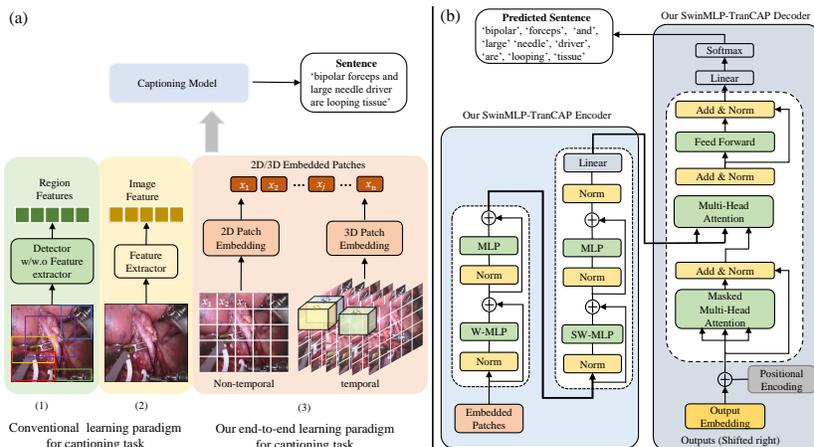
**Fig. 1.** (a) Comparisons of conventional learning paradigm and our learning paradigm for captioning task and (b) our proposed end-to-end SwinMLP-TranCAP model. (1) Captioning models based on an object detector w/w.o feature extractor to extract region features. (2) To eliminate the detector, the feature extractor can be applied as a compromise to the output image feature. (c) To eliminate the detector and feature extractor, the captioning models can be designed to take the patches as the input representation directly.

## 2 Methodology

### 2.1 Preliminaries

**Vision Transformer** Vision Transformer (ViT) [5] divides the input image into several non-overlapping patches with a patch size of $16 \times 16$. The feature dimension of each patch is $16 \times 16 \times 3 = 768$. We compute the self-attention for a given patch with all the other patches in the input image. This becomes very compute-intensive even for a reasonably large-sized image: $\Omega(MSA) = 4hwC^2 + 2(hw)^2$, where hw indicates the number of patches, C stands for the embedding dimension.

**Swin Transformer** The drawback of Multi-Head Self Attention (MSA) in ViT is attention calculation is very compute-heavy. To overcome it, Swin Transformer [10] introduces "window" to perform local self-attention computation. Specifically, a single layer of transformer is replaced by two layers which design Window-based MSA (W-MSA) and Shifted Window-based MSA (SW-MSA) respectively. In the first layer with W-MSA, The input image is divided into several windows and we compute self-attention between patches within that window. The intuition is local neighborhood patches are more important than attending patches that are far away. In the second layer with SW-MSA, the shifted window is designed to achieve information sharing across various spatial locations. The computation cost can be formulated as $\Omega(\text{W-MSA}) = 4hwC^2 + 2M^2hwC$, where
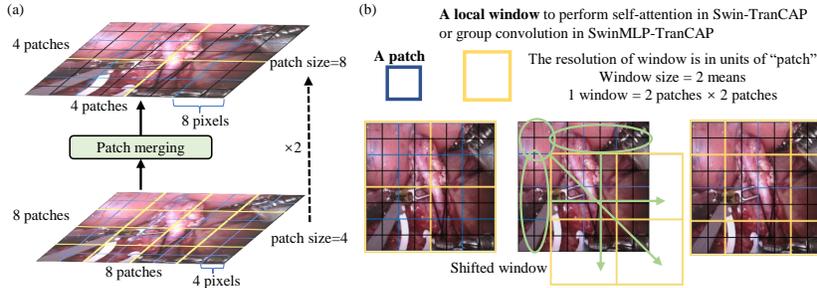
**Fig. 2.** (a) Patch merging layer in our Swin-TranCAP and SwinMLP-TranCAP. (b) shifted window to perform self-attention in Swin-TranCAP and group convolution in SwinMLP-TranCAP. For space without any pixels, cycle shifting is utilized to copy over the patches on top to the bottom, from left to right, and also diagonally across to make up for the missing patches.

M stands for window size and $M^2$ indicates the number of patches inside the window.

## 2.2   Our Model

We propose the end-to-end SwinMLP-TranCAP model for captioning tasks, as shown in Fig. 1 (b), which takes the embedded patches as the input representation directly. SwinMLP-TranCAP consists of a hierarchical fully multi-layer perceptron (MLP) architecture-based visual encoder and a Transformer-like decoder. Our model creates a new end-to-end learning paradigm for captioning tasks by using patches to get rid of the intermediate modules such as detector and feature extractor (see Fig. 1 (a)), reduces the training parameters, and improves the inference speed. In addition, we designed the temporal version of SwinMLP-TranCAP, named Video SwinMLP-TranCAP to implement video captioning.

**2D/3D Patch Embedding Layer** The raw image $[3, H, W]$ is partitioned and embedded into N discrete non-overlapping patches $[N, C]$ with the embedding dimension $C$ by using a 2D patch embedding layer which consists of a Conv2d projection layer (See Equation 1) followed by flattening and transpose operation and LayerNorm. The size of each patch is $[p, p, 3]$ and the number of patches $N$ is $\frac{H}{p} \times \frac{W}{p}$. The raw feature dimension $p \times p \times 3$ of each patch will be projected into an arbitrary embedding dimension $C$. Thus, the output size from the Conv2d projection layer is $[C, \frac{H}{p}, \frac{W}{p}]$. Next, it is flattened and transposed into embedded patches of $[N, C]$. Eventually, LayerNorm is applied to embedded patches. Similarly, in the 3D patch embedding layer, the video clip $[3, T, H, W]$ is partitioned and embedded into $N$ 3D patches of $[C, \frac{T}{t}, \frac{H}{p}, \frac{W}{p}]$ (t stands for the sequence length of a cubic patch, T stands for the sequence length of a video clip) via a Conv3d projection layer.

$$Conv2d(in = 3, out = C, kernel\_size = p, stride = p) \tag{1}$$

**Swin-TranCAP** We utilize Swin Transformer [10] as the backbone of the vision encoder and Transformer-like decoder to implement the window-based self-attention computation and reduce the computation cost.

**SwinMLP-TranCAP** SwinMLP-TranCAP consist of a SwinMLP encoder and a Transformer-like decoder, as shown in Fig. 2 (b). Embedded patches $[\frac{H}{p} \times \frac{W}{p}, C]$ are used as the input representation and each patch is treated as a "token". In SwinMLP encoder, we replace the Window/Shifted-Window MSA with Window/Shifted-Window Multi-Head Spatial MLP (W-MLP and SW-MLP) with the number of heads $n$ and window size $M$ via group convolution for less-expensive computation cost, which can be formulated as

$$Conv1d(in = nM^2, out = nM^2, kernel\_size = 1, groups = n) \qquad (2)$$

In two successive Swin MLP blocks, the first block with MLP is applied to image patches within the window independently and the second block with MLP is applied across patches by shifting window. The vision encoder consists of 4 stages and a linear layer. The Swin MLP block in the first stage maintain the feature representation as $[\frac{H}{p} \times \frac{W}{p}, C]$. Each stage in the remaining 3 stages contains a patch merging layer and multiple ($2\times$) Swin MLP Block. The patch merging layer (see Fig. 2 (a)) concatenates the features of each group of $2 \times 2$ neighboring patches and applies a linear layer to the $4C$ dimensional concatenated features. This decreases the amount of "tokens" by $2 \times 2 = 4$, and the output dimension is set to $2C$. Thus the feature representation after patch merging layer is $[\frac{H}{2\times p} \times \frac{W}{2\times p}, 2C]$. Afterward, Swin MLP blocks maintain such a shape of feature representation. Swin MLP blocks are created by replacing the standard MSA in a Transformer block with the window-based multi-head MLP module while keeping the other layer the same. The block in the vision encoder consists of a window-based MLP module, a 2-layer MLP with GELU nonlinearity. Before each windowed-based MLP module and each MLP, a LayerNorm layer is applied, and a residual connection is applied after each module. Overall, the number of tokens is reduced and the feature dimension is increased by patch merging layers as the network gets deeper. The output of 4 stages is $[\frac{H}{2^3\times p} \times \frac{W}{2^3\times p}, 2^3 \times C]$ (3 means the remaining 3 stages). Eventually, a linear layer is applied to produce the $[\frac{H}{2^3\times p} \times \frac{W}{2^3\times p}, 512]$ feature representation to adapt the Transformer-like decoder. The decoder is a stack of 6 identical blocks. Each block consists of two multi-head attention layers (an encoder-decoder cross attention layer and a decoder masked self-attention layers) and a feed-forward network.

**Video SwinMLP-TranCAP** We form the video clip $x = \{x_{t-T+1}, ..., x_t\}$ consisting the current frame and the preceding $T - 1$ frames. The video clip is partitioned and embedded into $\frac{T}{t} \times \frac{H}{p} \times \frac{W}{p}$ 3D patches and each 3D patch has a $C$-dimensional feature via a 3D Patch Embedding layer. Each 3D patch of size $[t, p, p]$ is treated as a "token". The major part is the video SwinMLP block which is built by the 3D shifted window-based multi-head MLP module while keeping the other components same. Given a 3D window of $[P, M, M]$, in two consecutive layers, the multi-head MLP module in the first layer employs the regular window partition approach to create non-overlapping 3D windows of $[\frac{T}{t\times P} \times \frac{H}{p\times M} \times \frac{W}{p\times M}]$.

In the second layer with multi-head MLP module, the window is shifted along the temporal, height, and width axes by $[\frac{P}{2}, \frac{M}{2}, \frac{M}{2}]$ tokens from the first layer's output.

**Our models depart from ViT [5] and Swin Transformer [10] with several points:** 1) no extra "class" token; 2) no self-attention blocks in the visual encoder: it is replaced by MLP via group convolutions; 3) Designed for captioning tasks. The window-based visual encoder is paired with a Transformer-like decoder via a linear layer which reduces the dimensional to 512 to adapt to the language decoder; 4) use window size of 14, instead of 7 used in Swin Transformer [10].

## 3    Experiments

### 3.1    Dataset Description

**DAISI dataset** The Database for AI Surgical Instruction (DAISI) [14] [5] contains 17339 color images of 290 different medical procedures, such as laparoscopic sleeve gastrectomy, laparoscopic ventral hernia repair, tracheostomy, open cricothyroidotomy, inguinal hernia repair, external fetal monitoring, IVC ultrasound, etc. We use the filtered DAISI dataset cleaned by [20] which removes noisy and irrelevant images and text descriptions. We split the DAISI dataset into 13094 images for training, and 1646 images for validation by following [20].

**EndoVis 2018 Dataset** is from the MICCAI robotic instrument segmentation dataset[6] of endoscopic vision challenge 2018 [1]. The training set consists of 15 robotic nephrectomy procedures acquired by the da Vinci X or Xi system. We use the annotated caption generated by [17] which employs 14 sequences out of the 15 sequences while disregarding the 13th sequence due to the less interaction. The validation set consists of the 1st, 5th, and 16th sequences, and the train set includes the 11 remaining sequences following the work [17].

### 3.2    Implementation details

Our models are trained using cross-entropy loss with a batch size of 9 for 100 epochs. We employ the Adam optimizer with an initial learning rate of $3e - 4$ and follow the learning rate scheduling strategy with 20000 warm-up steps. Our models are implemented on top of state-of-the-art captioning architecture, Transformer [5]. The vanilla Transformer architecture is realized with the implementation [7] from [5]. All models are developed with Pytorch and trained with NVIDIA RTX3090 GPU. In our models, the embedding dimension $C = 128$, patch size $p = 4$, and window size $M = 14$. In Video SwinMLP-TranCAP, we use the sequence length $T$ of 4 and 3D patch of $[2, 4, 4]$ and window size of $[2, 14, 14]$.

---

[5] https://engineering.purdue.edu/starproj/_daisi/
[6] https://endovissub2018-roboticscenesegmentation.grand-challenge.org/
[7] https://github.com/ruotianluo/ImageCaptioning.pytorch

## 4   Results and Evaluation

**Table 1.** Comparison of hybrid models and our models. The DAISI dataset does not have object annotation and video captioning annotation. Thus the hybrid of YLv5 w. RN18 and Transformer, and Video SwinMLP-Tran experiments are left blank.

| Model | | | DAISI Dataset | | | | EndoVis18 Dataset | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Det | FE | Captioning Model | B4 | MET | SPI | CID | B4 | MET | SPI | CID |
| FasterRCNN [12] | RN18 [7] | Tran [5] | | | ✗ | | 0.363 | 0.323 | 0.512 | 2.017 |
| | | Self-Seq [13] | | | | | 0.295 | 0.283 | 0.496 | 1.801 |
| | | AOA [8] | | | | | 0.377 | 0.371 | 0.58 | 1.811 |
| YLv5x [9] | RN18 [7] | Tran [5] | | | ✗ | | 0.427 | 0.328 | 0.577 | **3.022** |
| ✗ | RN18 [7] | Self-Seq [13] | 0.296 | 0.207 | 0.330 | 2.827 | 0.446 | 0.353 | 0.531 | 2.674 |
| | | AOA [8] | 0.349 | 0.246 | 0.403 | 3.373 | 0.427 | 0.322 | 0.533 | 2.903 |
| | | Tran [20] | 0.454 | 0.308 | **0.479** | **4.283** | 0.426 | 0.335 | 0.524 | 2.826 |
| ✗ | 3DRN18 | Tran [5] | | | ✗ | | 0.406 | 0.345 | 0.586 | 2.757 |
| Ours | | Swin-TranCAP | 0.346 | 0.237 | 0.378 | 3.280 | **0.459** | 0.336 | 0.571 | 3.002 |
| | | SwinMLP-TranCAP | **0.459** | **0.308** | 0.478 | 4.272 | 0.403 | 0.313 | 0.547 | 2.504 |
| | | V-SwinMLP-TranCAP | | | ✗ | | 0.423 | **0.378** | **0.619** | 2.663 |

We evaluate the captioning models using the BLEU-4 (B4) [11], METEOR (MET) [3], SPICE (SPI) [2], and CIDEr (CID) [16]. We employ the ResNet18 [7] pre-trained on the ImageNet dataset as the feature extractor (FE), YOLOv5 (YLv5) [9], and FasterRCNN [12] pre-trained on the COCO dataset as the detector (Det). The captioning models includes self-sequence captioning model (Self-Seq) [13], Attention on Attention (AOA) model [8], Transformer model [5]. Self-sequence and AOA originally take the region features extracted from the object detector with feature extractor as input. In our work, we design the hybrid style for them by sending image features extracted by the feature extractor only. 3D ResNet18 is employed to implement the video captioning task. We define the end-to-end captioning models take the image patches directly as input, trained from scratch as **Ours**, including Swin-TranCAP and SwinMLP-TranCAP for image captioning and Video SwinMLP-TranCAP (V-SwinMLP-TranCAP) for video captioning.

**Table 2.** Proof of "less computation cost" of our approaches with FPS, Num of Parameters, and GFLOPs.

| Model | | | Proof of less-computation cost | | |
|---|---|---|---|---|---|
| Det | FE | Captioning Model | FPS | N_Parameters(M) | GFLOPs |
| FasterRCNN [12] | RN18 [7] | Tran [5] | 8.418 | 28.32+46.67 | 251.84+25.88 |
| YLv5x [9] | RN18 [7] | Tran [5] | 9.368 | 97.88+46.67 | 1412.8+25.88 |
| ✗ | RN18 [7] | Tran [20] | 11.083 | 11.69+46.67 | 1.82+25.88 |
| Ours | | Swin-TranCAP | 10.604 | 165.51 | 19.59 |
| | | SwinMLP-TranCAP | **12.107** | **99.11** | **14.15** |

We now briefly describe how to use the detector and feature extractor extract features for captioning model. The spatial features from YOLOv5 w. ResNet are
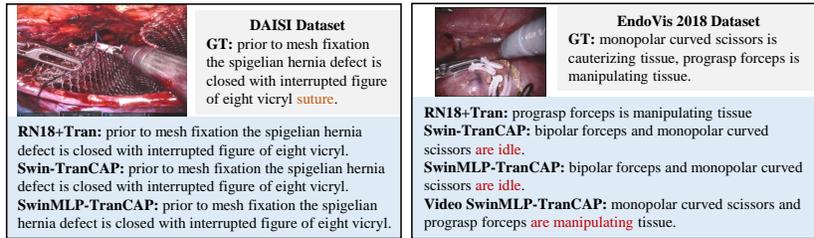
**Fig. 3.** Qualitative results with our model and hybrid models.

$(X, 512)$. $X$ indicates the number of predicted regions, and $X$ varies from image to image. The first $N$ predicted regions are sent to the captioning model. Zero appending is used to deal with those images where $N > X$ ($N = 6$ in our work). When only using the feature extractor, we take the 2D adaptive average pooling of the final convolutional layer output, which results in $14 \times 14 \times 512$ output for each image. It is reshaped into $196 \times 512$ before sending into the captioning model. We find that our models can achieve decent performance even without the need for a detector and feature extractor (see Table 1 and Fig. 3). On the DAISI dataset, our SwinMLP-TranCAP preserves the performance compared with the hybrid Transformer model and outperforms the FC and AOA model including 11% gain in BlEU-4. On EndoVis 2018 Dataset, our Swin-TranCAP model shows 4.7% improvement in SPICE compared with the hybrid Transformer model. Although the performance of SwinMLP-TranCAP drops a bit, it has less computation. In Table 2, less computation cost of our approaches is proven by evaluating FPS, Num of Parameters, and GFLOPs. Our approach achieves better captioning performance/efficiency trade-offs. It is worth highlighting that our purpose is not to provide better results but to remove the object detector and feature extractor from the conventional captioning system training pipeline for more flexible training, less computation cost, and faster inference speed without sacrificing performance. Surprisingly, our approach also obtained a slightly better quantitative performance.

**Table 3.** Tune patch size $p$ for vanilla Transformer.

| | DAISI | | EndoVis18 | |
|---|---|---|---|---|
| p | B4 | CID | B4 | CID |
| 4 | 0.192 | 1.703 | 0.339 | 2.105 |
| 8 | 0.247 | 2.195 | 0.364 | 1.886 |
| 16 | **0.302** | **2.776** | **0.416** | **3.037** |

**Table 4.** Tune embedding dimension $C$ with patch size $p$ of 4 and window size $M$ of 14, for our models.

| | Para. | | DAISI | | EndoVis18 | |
|---|---|---|---|---|---|---|
| Model | C | M | B4 | CID | B4 | CID |
| Swin-Tran-S | 96 | 14 | 0.329 | 3.109 | **0.471** | **3.059** |
| Swin-Tran-L | 128 | 14 | 0.346 | 3.280 | 0.459 | 3.002 |
| SwinMLP-Tran-S | 96 | 14 | 0.455 | 4.188 | 0.398 | 2.322 |
| SwinMLP-Tran-L | 128 | 14 | **0.459** | **4.272** | 0.403 | 2.504 |
| Swin-Tran-L | 128 | 7 | 0.433 | 4.049 | 0.389 | 2.932 |
| SwinMLP-Tran-L | 128 | 7 | 0.434 | 4.046 | 0.403 | 2.707 |

## 5    Ablation Study

We simply incorporate the patch embedding layer into the vanilla Transformer captioning model and report the performance in Table 3. The embedding dimension $C$ is set to 512. And we also investigate the performance of different configurations for our models (see Table 4). $C = 96$ paired with layer number $(2, 2, 6, 2)$ is **S**mall version and $C = 128$ paired with layer number $(2, 2, 18, 2)$ is **L**arger version. We find that the large version with a window size of 14 performs slightly better.

## 6    Discussion and Conclusion

We present the end-to-end SwinMLP-TranCAP captioning model, and take the image patches directly, to eliminate the object detector and feature extractor for real-time application. The shifted window and multi-head MLP architecture design make our model less computation. Video SwinMLP-TranCAP is also developed for video captioning tasks. Extensive evaluation of two surgical captioning datasets demonstrates that our models can maintain decent performance without needing these intermediate modules. Replacing the multi-head attention module with multi-head MLP also reveals that the generic transformer architecture is the core design, instead of the attention-based module. Future works will explore whether the attention-based module in a Transformer-like decoder is necessary or not. Different ways to implement patch embedding in the vision encoder are also worth studying.

## References

1. Allan, M., Kondo, S., Bodenstedt, S., Leger, S., Kadkhodamohammadi, R., Luengo, I., Fuentes, F., Flouty, E., Mohammed, A., Pedersen, M., et al.: 2018 robotic scene segmentation challenge. arXiv preprint arXiv:2001.11190 (2020)
2. Anderson, P., Fernando, B., Johnson, M., Gould, S.: Spice: Semantic propositional image caption evaluation. In: European conference on computer vision. pp. 382–398. Springer (2016)
3. Banerjee, S., Lavie, A.: Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In: Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization. pp. 65–72 (2005)
4. Cornia, M., Stefanini, M., Baraldi, L., Cucchiara, R.: Meshed-memory transformer for image captioning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10578–10587 (2020)

5. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929 (2020)

6. Fang, Z., Wang, J., Hu, X., Liang, L., Gan, Z., Wang, L., Yang, Y., Liu, Z.: Injecting semantic concepts into end-to-end image captioning. arXiv preprint arXiv:2112.05230 (2021)

7. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)

8. Huang, L., Wang, W., Chen, J., Wei, X.Y.: Attention on attention for image captioning. In: International Conference on Computer Vision (2019)

9. Jocher, G., Chaurasia, A., Stoken, A., Borovec, J., Kwon, Y., et al.: ultralytics/yolov5: v6.1 - TensorRT, TensorFlow Edge TPU and OpenVINO Export and Inference (Feb 2022). https://doi.org/10.5281/zenodo.6222936, https://doi.org/10.5281/zenodo.6222936

10. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. International Conference on Computer Vision (ICCV) (2021)

11. Papineni, K., Roukos, S., Ward, T., Zhu, W.J.: Bleu: a method for automatic evaluation of machine translation. In: Proceedings of the 40th annual meeting of the Association for Computational Linguistics. pp. 311–318 (2002)

12. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. Advances in neural information processing systems **28** (2015)

13. Rennie, S.J., Marcheret, E., Mroueh, Y., Ross, J., Goel, V.: Self-critical sequence training for image captioning. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 7008–7024 (2017)

14. Rojas-Muñoz, E., Couperus, K., Wachs, J.: Daisi: Database for ai surgical instruction. arXiv preprint arXiv:2004.02809 (2020)

15. Tolstikhin, I.O., Houlsby, N., Kolesnikov, A., Beyer, L., Zhai, X., Unterthiner, T., Yung, J., Steiner, A., Keysers, D., Uszkoreit, J., et al.: Mlp-mixer: An all-mlp architecture for vision. Advances in Neural Information Processing Systems **34** (2021)

16. Vedantam, R., Lawrence Zitnick, C., Parikh, D.: Cider: Consensus-based image description evaluation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4566–4575 (2015)

17. Xu, M., Islam, M., Lim, C.M., Ren, H.: Class-incremental domain adaptation with smoothing and calibration for surgical report generation. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 269–278. Springer (2021)

18. Xu, M., Islam, M., Lim, C.M., Ren, H.: Learning domain adaptation with model calibration for surgical report generation in robotic surgery. In: 2021 IEEE International Conference on Robotics and Automation (ICRA). pp. 12350–12356. IEEE (2021)

19. Yu, W., Luo, M., Zhou, P., Si, C., Zhou, Y., Wang, X., Feng, J., Yan, S.: Metaformer is actually what you need for vision. arXiv preprint arXiv:2111.11418 (2021)

20. Zhang, J., Nie, Y., Chang, J., Zhang, J.J.: Surgical instruction generation with transformers. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 290–299. Springer (2021)