# LViT: Language meets Vision Transformer in Medical Image Segmentation

Zihan Li, Yunxiang Li, Qingde Li, Puyang Wang, Dazhou Guo, Le Lu, Fellow, IEEE, Dakai Jin, Member, IEEE, You Zhang, Qingqi Hong, Member, IEEE

Abstract—Deep learning has been widely used in medical image segmentation and other aspects. However, the performance of existing medical image segmentation models has been limited by the challenge of obtaining sufficient high-quality labeled data due to the prohibitive data annotation cost. To alleviate this limitation, we propose a new text-augmented medical image segmentation model LViT (Language meets Vision Transformer). In our LViT model, medical text annotation is incorporated to compensate for the quality deficiency in image data. In addition, the text information can guide to generate pseudo labels of improved quality in the semi-supervised learning. We also propose an Exponential Pseudo label Iteration mechanism (EPI) to help the Pixel-Level Attention Module (PLAM) preserve local image features in semi-supervised LViT setting. In our model, LV (Language-Vision) loss is designed to supervise the training of unlabeled images using text information directly. For evaluation, we construct three multimodal medical segmentation datasets (image + text) containing X-rays and CT images. Experimental results show that our proposed LViT has superior segmentation performance in both fully-supervised and semisupervised setting. The code and datasets are available at https://github.com/HUANGLIZI/LViT.

Index Terms—Vision-Language, Medical image segmentation, Semi-supervised learning

# I. INTRODUCTION

EDICAL image segmentation is one of the most critical tasks in medical image analysis. In clinical practice, accurate segmentation results are often achieved manually or semi-automatically. It remains a challenging task to extract the desired object accurately, especially when the target organ to be extracted is of high complexity in terms of tissue structures. Recent research shows that deep learning can be a promising approach for automatic medical image segmentation, as the knowledge of experts can be learned and extracted by using a certain deep learning method. A summary of existing solutions is shown in Figure 1(a): (1) one shared encoder followed by two separate decoders [1]; (2) two separate encoders followed

Zihan Li is with Xiamen University and the Department of Computer Science, University of Illinois at Urbana-Champaign, Urbana, IL 61801, USA (e-mail: zl111@illinois.edu).

Yunxiang Li and You Zhang are with the Department of Radiation Oncology, UT Southwestern Medical Center, Dallas, TX 75235, USA.

Qingde Li is with the Department of Computer Science, University of Hull, Hull, HU6 7RX, UK.

Puyang Wang is with DAMO Academy, Alibaba Group, Hangzhou 310024, China.

Dazhou Guo, Le Lu, and Dakai Jin are with DAMO Academy, Alibaba Group, New York, NY 10014, USA.

Qingqi Hong is with Xiamen University, Xiamen 361005, China. (e-mail: hongqq@xmu.edu.cn).

Corresponding author: Qingqi Hong

by one shared decode [2]; (3) two separate encoders followed by a modality interaction model [3]. However, two inherent issues concerning the creation of high quality medical image datasets severely limit the application: one is the difficulty in obtaining high-quality images, and the other one is the high cost of data annotation [4], [5]. These two issues have dramatically limited the performance improvement of medical image segmentation models. Since it is challenging to improve the quantity and quality of medical images themselves, it may be more feasible to use complementary and easy-to-access information to make up for the quality defects of medical images. Thus, we turn our attention to the written medical notes accompanied by medical images. It is well known that text data of medical records are usually generated along with the patients, so no extra cost is needed to access the corresponding text data. The medical text record data and the image data are naturally complementary to each other, so the text information can compensate for the quality deficiency in the medical image data. On the other hand, expert segmentation annotation is often expensive and time-consuming, especially for new diseases like COVID-19, where high-quality annotations are even more difficult to obtain [4], [6], [7]. In order to address the issue of under-annotated data, some approaches have gone beyond traditional supervised learning by training their models using both labeled and more widely available unlabeled data, such as semi-supervised learning [5], [8] and weakly-supervised learning [9]. However, the performance of these approaches is largely determined by the credibility of the pseudo label. This is because the number of pseudo labels is much larger than ground truth labels. Therefore, the critical question to be answered is how to improve the quality of the pseudo label. To effectively address this issue, we develop a model that can be trained using the medical texts written by domain experts. By learning additional expert knowledge from text information, we can improve the quality of pseudo labels.

In summary, the challenges exist in two aspects: 1) How to improve the segmentation performance by using the existing image-text information; 2) How to make full use of text information to guarantee the quality of pseudo labels. To address the first challenge, We propose the LViT model (Figure 1(b)), which is innovative in processing images and text. In LViT, the text feature vector is obtained by using a embedding layer instead of text encoder, which can reduce the number of parameters in the model. In addition, the hybrid CNN-Transformer structure with Pixel-Level Attention Modules (PLAM) is able to better merge text information and encode global features with Transformer while retaining the CNN's

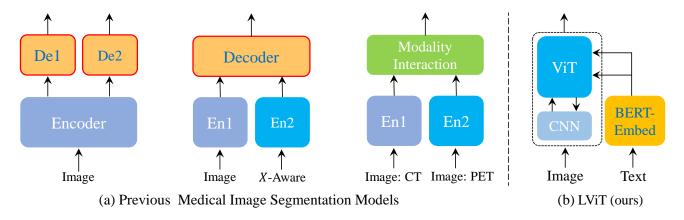


Fig. 1. Comparison of current medical image segmentation models and our proposed LViT model.

ability to extract local features from images. To address the second challenge, we design an Exponential Pseudo label Iteration mechanism (EPI) for the proposed LViT, aiming to cross-utilize the label information of labeled data and the latent information of unlabeled data. The EPI indirectly incorporates text information to refine the pseudo label progressively in the way of Exponential Moving Average (EMA) [10]. In addition, the LV (Language-Vision) loss is designed to utilize text information directly to supervise the training of unlabeled medical images. To validate the performance of LViT, we construct three multimodal medical image segmentation datasets containing CT images (MosMedData+ [11], [12] and ESO-CT) and X-rays (QaTa-COV19 [13]). Results show that LViT has superior segmentation performance, achieving 74.57% Dice score and 61.33% mIoU on the MosMedData+ dataset, 83.66% Dice score and 75.11% mIoU on the QaTa-COV19 dataset, and 71.53% Dice score and 59.94% mIoU on the ESO-CT dataset. And it is worth noting that LViT using 1/4 of the train set labels can still have the same performance as the fully-supervised segmentation method.

## II. RELATED WORK

## A. Semantic segmentation of medical images

Semantic segmentation can be considered as the work for pixel-level image classification, and thus many image classification networks have been extended [14]–[17] to implement semantic segmentation, with fully convolutional network (FCN) [14] being commonly considered as the first end-to-end pixel-to-pixel network for semantic segmentation [18]. Among them, U-Net [19] is considered as a pioneer in medical image segmentation. Based on this, UNet++ [20] improved the skip connection of U-Net. However, most of the above methods are very sensitive to the quantity and quality of the data, resulting in limited generalization performance of the models. Therefore, some approaches [5], [21]–[24] have explored the application of semi-supervised learning in different areas. The problem of lacking data and its annotation can also be further mitigated by introducing multiple modalities into the learning models.

# B. Vision-language model

CLIP [25] is a pioneering work of large-scale vision-language pretraining (VLP) model, which utilized contrast learning to

learn image representations on a dataset of 400 million pairs (image, text) from scratch. By simplifying the processing of visual inputs compared to CLIP, Kim et al. [26] proposed a more parameter-efficient architecture, i.e. ViLT, which allows exploiting the power of interaction layers to process visual features while lacking a separate deep visual embedder. Subsequently, there is a rich line of works on image segmentation [21], [27]–[31] that have begun to use text information to improve the segmentation capabilities of models. For instance, Ding et al. [29] developed a Vision-Language Transformer (VLT) framework for referring segmentation by facilitating deep interactions among multi-modal information and fusing linguistic and visual features. Different from VLT, Language-Aware Vision Transformer (LAVT) [27] framework adopted an early fusion scheme for integrating linguistic features into visual features via a pixel-word attention mechanism, which can effectively exploit the Transformer encoder for modeling multi-modal context. Inspired by VLP in natural image, a few studies have started to utilize text information for assisting with medical image analysis [32]–[34]. However, compared with natural image, medical image has its own characteristics. Unlike the natural image, boundaries between different regions in the medical image are often blurred, and the small gray-scale value differences in the vicinity of the boundaries makes it difficult to extract highly accurate segmentation boundaries. Therefore, it is not applicable to directly apply vision-language model in the natural images to medical image analysis. VTL and LAVT focus on addressing reference segmentation of the natural images, requiring language encoder for the explicit alignment of image and text. Due to the essential differences between medical images and natural images, it is very difficult to achieve strict alignment between text and image in medical images. Similarly, ViLT is designed to address multimodal problems, such as Visual Question Answering (VQA) and retrieval tasks. In addition, ViLT is a pure Transformer model without convolution or region supervision for extracting local features, which is not suitable for medical image segmentation with blurred boundaries. On the other hand, LViT is specifically tailored to medical image segmentation problems. In the LViT model, only the Embedding layer is utilized to transform text features, which requires fewer parameters and lower computational cost. In addition, the hybrid CNN-Transformer

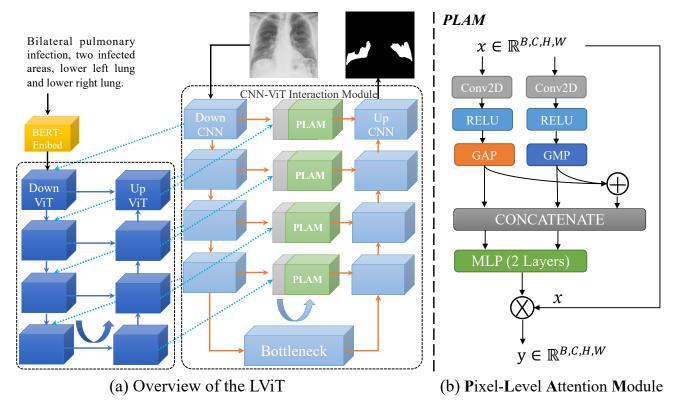


Fig. 2. Illustration of (a) the proposed LViT model, and (b) the Pixel-Level Attention Module (PLAM). The proposed LViT model is a Double-U structure formed by combining a U-shape CNN branch with a U-shaped ViT branch.

structure enables us to retain both local and global features of the image, which is crucial for extracting highly accurate segmentation boundaries on medical images. Furthermore, to address the scarcity of medical image labels, an LV (Language-Vision) loss is designed to supervise the training of unlabeled images using text information directly.

# C. Attention mechanism

Starting from RAN [35], researchers have begun to introduce attention mechanisms into the field of computer vision. Woo et al [36] proposed a well-known Convolutional Block Attention Module (CBAM), where both spatial and channel attentions are used to perform adaptive feature refinement. In addition to the original attention mechanisms [37], self-attention mechanisms [38]-[41] have also begun to enter the field of computer vision. However, since self-attention was originally proposed for solving NLP problems, it faces many challenges, such as high computational cost and neglecting local features of images. Therefore, we propose PLAM to compensate for the lack of attention to local features through self-attention. It also helps the convolutional layer to produce a more effective representation of local features. And to address the high computational problem, we utilize the uniform encoder to encode the vision and language features instead of using separate encoders.

# III. METHOD

As shown in Figure 2, the proposed LViT model is a Double-U structure consisting of a U-shaped CNN branch and a U-shaped Transformer branch. The CNN branch acts

as the source of information input and the segmentation head of prediction output, and the ViT branch is used to merge image and text information, where we exploit the ability of Transformer to process cross-modality information. After a simple vectorization of the text, the text vector is merged with the image vector and send to the U-shaped ViT branch for processing. In the model inference stage, we also need to perform similar processing on text input. And we pass the fusion information of corresponding size back to the U-shape CNN branch at each layer for the final segmentation prediction. In addition, a Pixel-Level Attention Module (PLAM) is set at the skip connection position of the U-shape CNN branch. With PLAM, LViT is able to retain as much local feature information of images as possible. We also conduct ablation experiments to demonstrate the effectiveness of each module.

## A. LViT Model

1) U-shape CNN Branch: As shown in Figure 2(a), the U-shaped CNN branch is used to receive the image information and act as segmentation head to output the prediction mask. The Conv, BatchNorm(BN), and ReLU activation layers are utilized to compose each CNN module. And image feature are downsampled between each DownCNN module using the MaxPool layer. Concatenate operation is performed between each UpCNN module. The specific process of each CNN module is described by Eqn. 1 and 2,

$$D_i = \text{DownCNN}_i = \text{Relu}\left(BN_i\left(\text{Conv}_i(\cdot)\right)\right)$$
 (1)

$$Y_{\text{DownCNN},i+1} = \text{MaxPool}\left(D_i\left(Y_{\text{DownCNN},i}\right)\right)$$
 (2)

where  $Y_{\mathrm{DownCNN},i}$  represents the input of the i-th DownCNN module, which becomes  $Y_{\mathrm{DownCNN},i+1}$  after the downsampling of the i-th DownCNN module and the MaxPool layer. In addition, we design the CNN-ViT interaction module using methods such as upsampling to align the features from ViT, as the details of CNN-ViT interaction module are shown in the appendix. The reconstructed ViT features are also connected with CNN feature by residuals to form CNN-ViT interaction features. In addition, to further improve the segmentation capability for local features, PLAM is designed at the skip connection in the U-shaped CNN branch. So the CNN-ViT interaction features will be fed into PLAM, then the interaction features are transferred to the UpCNN module to give the upward information layer by layer.

2) U-shape ViT Branch: Referring to the U-shaped CNN branch, the U-shaped ViT branch is designed for merging image features and text features. As shown in Figure 2(a), the first layer DownViT module receives the text feature input from BERT-Embed [42] and the image feature input from the first layer DownCNN module. The pretraining model of the BERT-Embed is the BERT\_12\_768\_12 model, which can convert a single word into a 768-dimensional word vector. The specific cross-modal feature merging operation is expressed by the following equations,

$$Y_{\text{DownViT},1} = \text{ViT} (x_{\text{img},1} + \text{CTBN} (x_{\text{text}}))$$
 (3)

$$x_{\text{img, }i} = \text{PatchEmbedding}(Y_{\text{DownCNN},i})$$
 (4)

$$x' = ViT_1(x) = MHSA(LN(x)) + x$$
 (5)

$$Y = \operatorname{ViT}_{2}(x') = \operatorname{MLP}\left(\operatorname{LN}(x')\right) + x' \tag{6}$$

where  $x_{img,i}$  represents the image features from DownCNN,  $x_{text}$  represents the text features, and PatchEmbedding can help  $Y_{DownCNN,i}$  form the embedding features  $x_{img,i}$ . ViT represents the Transformer encoder [39], i.e.,  $Y = ViT(x) = ViT_2(ViT_1(x))$ . ViT consists of the Multi-headed Self-attention (MHSA) module and the MLP layer. And LN represents the normalization layer. The CTBN block also consists of the Conv layer, BatchNorm layer, and ReLU activation layer for aligning the feature dimensions of  $x_{img,1}$  and  $x_{text}$ . The subsequent DownViT modules of layers 2, 3, and 4 receive both feature from the upper DownViT module and the feature from the DownCNN module of the corresponding layer, as shown in Eqn. 7,

$$Y_{\text{DownViT},i+1} = \text{ViT} \left( Y_{\text{DownViT},i} + x_{img,i+1} \right)$$
 (7)

where i=1,2,3. The features of the corresponding size are then transferred back to the CNN-ViT interaction module through the UpViT module. And the feature is merged with the feature from the DownCNN module of the corresponding layer. This will maximize the extraction of image global features and avoid the oscillation of the model performance due to the inaccuracy of text annotation.

3) Pixel-Level Attention Module (PLAM): As shown in Figure 2(b), PLAM is designed to preserve the local features of images and further merge the semantic features in the text. Besides, it can enhance the performance of the convolutional layer in generating a powerful representation of local features. Referring to CBAM [36], our PLAM uses parallel branches for

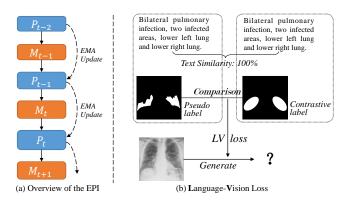


Fig. 3. Illustration of (a) Exponential Pseudo-label Iteration mechanism (EPI), and (b) LV (Language-Vision) Loss.

Global Average Pooling (GAP) and Global Max Pooling (GMP). We also incorporate the concatenate and add operations. The add operation will help merge the corresponding channel features with similar semantics and save computation. In contrast, the concatenate operation can integrate the feature information more intuitively and help preserve the original features of each part. After concatenating the feature information, we use the MLP structure and the multiplication operation to help align the feature size. Generally, our PLAM differs in several aspects from the Pixel-word attention module (PWAM) in LAVT [27]. Firstly, PLAM is designed to enhance local features to mitigate the preference for global features brought by Transformer. In contrast, PWAM is designed to align visual and linguistic representation using cross-attention. Secondly, in terms of implementation, PLAM utilizes a combination of channel attention and spatial attention, while PWAM uses cross self-attention mechanism. Overall, PLAM aims to enhance local features for the purpose of improving the performance of medical image segmentation with text information. On the other hand, PWAM is designed to align the multimodal features for achieving better referring segmentation results.

# B. Exponential Pseudo-label Iteration mechanism

In this section, we propose the Exponential Pseudo label Iteration mechanism (EPI), which is designed to help extend the semi-supervised version of LViT. In EPI, the pseudo label is iteratively updated using the idea of EMA [10], as shown in Figure 3(a) and Eqn. 8,

$$P_t = \beta \cdot P_{t-1} + (1 - \beta) \cdot P_t \tag{8}$$

where  $P_{t-1}$  represents the prediction of model  $M_{t-1}$ , and  $\beta$  is set as the momentum parameter to 0.99. It is worth noting that here  $P_{t-1}$  is an N-dimensional prediction vector, where N represents the number of category classes, and each dimension represents the prediction probability. Therefore, EPI can gradually optimize the segmentation prediction results of the model for each unlabeled pixel and be robustness to noisy labels. This is because we do not simply use the pseudo label predicted by one generation of models as the target for the next-generation model, which can avoid sharp deterioration of the pseudo label quality. The theoretical proof for the effectiveness of EPI algorithm is as follows.

**Proof:** The basic assumption for the EPI algorithm is the model weights will dither around the actual optimum in the last n generations, and therefore the pseudo label predicted by the model will also dither around the mask in the last ngenerations. We expand  $P_t$  around t in Eqn. 8 to Eqn.9,

$$P_{t} = \beta^{n} P_{t-n} + (1 - \beta) \cdot (\beta^{n-1} P_{t-n+1} + \dots + P_{t}). \quad (9)$$

In particular, we let  $n = 1/(1-\beta)$  and  $\beta^n = \beta^{\frac{1}{1-\beta}} \approx \frac{1}{4}$ . So for the first  $1/(1-\beta)$  generations,  $P_t$  decays to a weighted average of 1/e. Further, we introduce an adjustment gradient  $g_{t-1}$  for the predicted label  $P_{t-1}$ , which leads to Eqn. 10,

$$P_t = P_{t-1} - g_{t-1} = \dots = P_1 - \sum_{i=1}^{n-1} g_i.$$
 (10)

Similarly, we extend Eqn. 9 when t = n and  $P_0 \approx P_1$ . Comparing with Eqn. 10, it can be seen from Eqn. (11 - 15) that the EPI algorithm adds the weight coefficient  $1-\beta^{n-i}$  for the gradient descent step of the *i*-th iteration.  $1 - \beta^{n-i}$  will decreases as i increases, so the change of pseudo label is finally stabilized and obtain the pseudo label with high confidence.

$$\stackrel{\sim}{P} = \beta^n P_0 \tag{11}$$

$$P_t = \overset{\sim}{P} + (1 - \beta) \cdot (\beta^{n-1} P_1 + \beta^{n-2} P_2 + \dots + P_n)$$
 (12)

$$P_t = \stackrel{\sim}{P} + (1 - \beta) \cdot \left( \beta^{n-1} P_1 + \dots + \left( P_1 - \sum_{i=1}^{n-1} g_i \right) \right)$$
 (13)

$$P_t = \stackrel{\sim}{P} + (1 - \beta) \cdot \left( \frac{1 - \beta^n}{1 - \beta} P_1 - \sum_{i=1}^{n-1} \frac{(1 - \beta^{n-i}) g_i}{1 - \beta} \right)$$
 (14)

$$P_t \approx P_1 - \sum_{i=1}^{n-1} (1 - \beta^{n-i}) g_i.$$
 (15)

## C. LV (Language-Vision) Loss

To further utilize the text information to guide the pseudolabel generation, we design the LV (Language-Vision) loss function, as shown in Figure 3(b). Generally, the positions of human organs in medical images are relatively fixed. Thus, we can use structured text information to form the corresponding mask (the contrastive label). And we calculate the cosine similarity between the texts, as shown in Eqn. 16,

TextSim = 
$$\frac{x_{\text{text },p} \cdot x_{\text{text},c}}{|x_{\text{text },p}| \times |x_{\text{text },c}|}$$
(16)

where  $x_{text,p}$  represents the text feature vector corresponding to the pseudo label, and  $x_{text,c}$  represents the text feature vector corresponding to the contrastive label. After that, according to TextSim, we select the contrastive text with the highest similarity and find the segmentation mask corresponding to that text; we calculate the cosine similarity between the predicted segmentation pseudo-label and the contrastive label using the label similarity, as shown in Eqn. 17 and 18,

$$ImgSim = \frac{x_{img,p} \cdot x_{img,c}}{|x_{img,p}| \times |x_{img,c}|}$$

$$L_{LV} = 1 - ImgSim$$
(17)

$$L_{TV} = 1 - \text{ImgSim} \tag{18}$$

where  $x_{img,p}$  represents the pseudo-label feature vector, and  $x_{img,c}$  represents the comparison label feature vector. Compared to Euclidean distance, cosine similarity is not sensitive to absolute values and reflects the degree of similarity

more qualitatively, consistent with our task motivation. The contrastive labels mainly provide labeling information of the approximate location instead of refinement for the boundaries. Therefore, the primary purpose of LV loss is to avoid missegmentation or mislabelled cases with significant differences. For this reason, we only use LV loss in the unlabeled case because the contrastive labels are of little help for performance improvement when the data is labeled. And in case of no label, LV loss with consistency supervision can avoid the sharp deterioration of the pseudo-label quality. It is important to note that the Pseudo and Contrastive labels in our LViT aim to address different issues compared to the masked conservative learning in VLT [29]. Firstly, the Pseudo and Contrastive labels are designed for semi-supervised learning, whereas masked conservative learning aims to explore the knowledge of different language expressions pertaining to a single object. Secondly, LViT determines whether a case is similar by calculating text similarity, while VLT achieves this by extracting text features. However, it is difficult to determine the similarity between radiology reports through implicit feature extraction in the medical field, as different radiology reports may have only a few wording changes. Therefore, structured formats are typically used to differentiate between reports. In addition, different from the masked conservative learning, we design an Exponential Pseudo label Iteration mechanism (EPI) to guarantee the quality of pseudo labels with text information, which utilizes the label information of labeled data and the latent information of unlabeled data in a cross-utilized manner.

# D. Proof of CNN-Transformer structure superiority

Unlike the previous Vision-and-Language work, our proposed LViT model is innovative in processing images and text. We do not use text encoder and creatively use the interaction between CNN and ViT to extract features.

**Proof:** For the sake of description, we assume that the patch size in ViT is equal to the kernel size in CNN, which is S. The input matrix is M, and output of convolution is  $Y_{cnn}$ .

$$Y_{cnn,k}(i,j) = \sum_{\xi=0}^{S} \sum_{\eta=0}^{S} f_k(\xi,\eta) * M(i-\xi,j-\eta)$$
 (19)

$$Y_{cnn}(i,j) = [Y_{cnn,1}(i,j); \dots; Y_{cnn,C}(i,j)]$$
 (20)

where  $f_k$  represents the convolution kernel of the k-th channel, and  $Y_{cnn,k}(i,j)$  represents the output of the k-th channel after convolution. The total convolution outputs of C channels form  $Y_{cnn}(i,j)$ . And the convolution operation f(x)is satisfying shift invariance and scale invariance, i.e., if Y(x) = f(x) \* M(x), then we have  $Y(x-\delta) = f(x) * M(x-\delta)$ , and if Y(x) = f(x) \* M(x), then we have  $|\delta|Y(x/\delta) =$  $f(x/\delta) * M(x/\delta)$ . Therefore, CNNs are good at learning shallow features and are affine-invariant. The kernel size is fixed, so each kernel can only learn one aspect of local information, like points, lines, and boundaries. And since the convolutional kernel  $f_k(\xi, \eta)$  of each channel shares the weights on the whole image, convolving the whole image with convolutional kernels that focus on boundary features is equivalent to doing whole-image filtering on the image. Similarly, we set the output after multi-head self-attention as  $Y_{vit}$ , as shown in the Eqn. 21 and 22,

$$Y_{vit,h} = \text{Softmax}\left(\frac{Q_h^T \cdot K_h}{\sqrt{d}}\right) \cdot V_h$$
 (21)

$$Y_{vit} = LN([Y_{vit,1}; Y_{vit,2}; ...; Y_{vit,H}])$$
 (22)

where  $Y_{vit,h}$  denotes the output of the h-th self-attention head, and d prevent the feature gradient from vanishing after Softmax. LN represents the linear layer, which aims to reduce the dimensionality of the output features. And  $Q_h$ ,  $K_h$ , and  $V_h$ in the self-attention mechanism are transformations for their own inputs,  $\operatorname{Softmax}\left(\frac{Q_h^T \cdot K_h}{\sqrt{d}}\right) \cdot V_h$  is computing the similarity between them. So self-attention is essentially focusing on the invariance of input features. The input features M are the whole image for ViT, so ViT is easier to learn the global features with more robustness than CNN.

TABLE I THE SPECIFIC DIVISION OF DIFFERENT DATASETS.

	QaTa-COV19	MosMedData+	ESO-CT
Train set	5716	2183	182
Val set	1429	273	46
Test set	2113	273	58
Total	9258	2729	286

#### IV. EXPERIMENTS

## A. Setup

Three datasets are used in the experiments to evaluate the performance of our method. The first one is the MosMedData+1 dataset [11], [12], which contains 2729 CT scan slices of lung infections. The second one is the OaTa-COV19 dataset [13], which is compiled by researchers from Qatar University and Tampere University. This dataset consists of 9258 COVID-19 chest X-ray radiographs with manual annotations of COVID-19 lesions for the first time. In addition, text annotations for the datasets are extended by us to be used for training the vision-language model. We extend text annotations on the QaTa-COV19 dataset for the first time with the help of professionals. The text annotations focuse on whether both lungs are infected, the number of lesion regions, and the approximate location of the infected areas. For example, "Bilateral pulmonary infection, two infected areas, upper left lung and upper right lung." refers to bilateral lung infection, and there are two infection areas located in the upper left lung and the upper right lung respectively. The text annotations on MosMedData+ dataset mainly contain the same information as QaTa-COV19 dataset, and the text structure is similar, e.g., "Unilateral pulmonary infection, two infected areas, middle lower **left lung."**. The third dataset is the ESO-CT dataset, which consists of 286 cases, and the detail will be presented in the section of generalization study. Those text annotations were provided and verified by two professionals from the Department of Radiation Oncology, UT Southwestern Medical Center. The radiologists independently annotated the same image, and then we compared their annotations to ensure

consistency. Additionally, we conducted a quality check based on the provided mask to ensure that there was no excessive deviation in the text annotations. The loss function we use is shown in Eqn. 26, where  $L_{Dice}$  means dice loss and  $L_{CE}$ means cross-entropy loss. For the unlabeled data, an additional term on the loss  $L_{LV}$  is introduced with  $\alpha = 0.1$ . And for the labeled data,  $\alpha = 0$ . Dice and mIoU are used to evaluate the segmentation performance. And early stop mechanism is used during training phase.

$$L_{Dice} = 1 - \sum_{i=1}^{N} \sum_{j=1}^{C} \frac{1}{NC} \cdot \frac{2|p_{ij} \cap y_{ij}|}{(|p_{ij}| + |y_{ij}|)}$$
(23)

$$L_{CE} = -\sum_{i=1}^{N} \sum_{j=1}^{C} \frac{1}{N} \cdot y_{ij} \log(p_{ij})$$
 (24)

$$L_{sup} = (L_{Dice} + L_{CE})/2 \tag{25}$$

$$L_{unsup} = (L_{Dice} + L_{CE})/2 + \alpha \cdot L_{LV}$$
 (26)

where N represents the number of pixels, C represents the number of categories, which is set to 1 in our experiments.  $p_{ij}$  represents the prediction probability that pixel i belongs to category j,  $y_{ij}$  represents whether pixel i belongs to category j. If pixel i belongs to category j, then  $y_{ij}$  is 1, otherwise 0.

#### B. Evaluation Metrics

For the evaluation metrics, the Dice score and the mIoU metric are used to evaluate the performance of our LViT model and other SOTA methods, as shown in Eqn. 27 and 28,

$$Dice = \sum_{i=1}^{N} \sum_{j=1}^{C} \frac{1}{NC} \cdot \frac{2|p_{ij} \cap y_{ij}|}{(|p_{ij}| + |y_{ij}|)} = 1 - L_{Dice}$$
(27)  
$$mIoU = \sum_{i=1}^{N} \sum_{j=1}^{C} \frac{1}{NC} \cdot \frac{|p_{ij} \cap y_{ij}|}{|p_{ij} \cup y_{ij}|}$$
(28)

$$mIoU = \sum_{i=1}^{N} \sum_{j=1}^{C} \frac{1}{NC} \cdot \frac{|p_{ij} \cap y_{ij}|}{|p_{ij} \cup y_{ij}|}$$
(28)

where N represents the number of pixels, C represents the number of categories,  $p_{ij}$  and  $y_{ij}$  also have the same definition as in the above section.

#### C. Implementation Details

Our proposed approach is implemented using Pytorch. The main parameters of the server are listed below: the operating system is Ubuntu 16.04.12 LTS, the CPU is Intel(R) Xeon(R) Gold 5218, the GPU is a 2-card TESLA V100 32G, and the memory capacity is 128 GB. In terms of dataset division, we split the train set and the validation set from the original train set. Then, the train set is divided into labeled and unlabeled train sets in a specific ratio. The number of samples in each dataset is presented in Table I.

The initial learning rate is set to 3e-4 for the QaTa-COV19 dataset and 1e-3 for the MosMedData+ dataset. We also use an early stop mechanism until the performance of model does not improve for 50 epochs. Different batch sizes are also set for each dataset since they have different data size. The default batch size is 24 for the QaTa-COV19 dataset and the MosMedData+ dataset.

# D. Comparison with State-of-the-Art Methods

We compare the performance of our LViT model with several CNN and Transformer based segmentation models. The number of network parameters and the computational cost of different methods are also reported. Note that the numbers after the methods refer to the ratio of labels used, e.g., LViT-T (1/2)

<sup>&</sup>lt;sup>1</sup>http://medicalsegmentation.com/covid19

TABLE II

PERFORMANCE COMPARISON BETWEEN OUR METHOD (LVIT) AND OTHER STATE-OF-THE-ART METHODS ON THE QATA-COV19 AND MOSMEDDATA+
DATASETS. THE "W" IN LVIT-TW REFERS TO WITHOUT THE TEXT INFORMATION. THE "HYBRID" MEANS CNN-TRANSFORMER STRUCTURE.

						QaTa-	COV19	MosMo	edData+
Method	Backbone	Text	Label ratio	Param (M)	Flops (G)	Dice (%)	mIoU (%)	Dice (%)	mIoU (%)
U-Net [19]	CNN	×	100%	14.8	50.3	79.02	69.46	64.60	50.73
UNet++ [20]	CNN	×	100%	74.5	94.6	79.62	70.25	71.75	58.39
AttUNet [43]	CNN	×	100%	34.9	101.9	79.31	70.04	66.34	52.82
nnUNet [44]	CNN	×	100%	19.1	412.7	80.42	70.81	72.59	60.36
TransUNet [45]	Hybrid	×	100%	105	56.7	78.63	69.13	71.24	58.44
Swin-Unet [46]	Hybrid	×	100%	82.3	67.3	78.07	68.34	63.29	50.19
UCTransNet [47]	Hybrid	×	100%	65.6	63.2	79.15	69.60	65.90	52.69
LViT-TW (1/4)	Hybrid	×	25%	28.0	54.0	79.08	69.42	70.65	58.07
LViT-TW (1/2)	Hybrid	×	50%	28.0	54.0	80.35	70.74	71.89	59.63
LViT-TW	Hybrid	×	100%	28.0	54.0	81.12	71.37	72.58	60.40
ConVIRT [48]	CNN	<b>√</b>	100%	35.2	44.6	79.72	70.58	72.06	59.73
TGANet [34]	CNN	$\checkmark$	100%	19.8	41.9	79.87	70.75	71.81	59.28
CLIP [25]	Hybrid	$\checkmark$	100%	87.0	105.3	79.81	70.66	71.97	59.64
GLoRIA [49]	Hybrid	$\checkmark$	100%	45.6	60.8	79.94	70.68	72.42	60.18
ViLT [26]	Hybrid	$\checkmark$	100%	87.4	55.9	79.63	70.12	72.36	60.15
LAVT [27]	Hybrid	$\checkmark$	100%	118.6	83.8	79.28	69.89	73.29	60.41
LViT-T (1/4)	Hybrid	<b>√</b>	25%	29.7	54.1	80.95	71.31	72.48	60.31
LViT-T (1/2)	Hybrid	$\checkmark$	50%	29.7	54.1	82.73	73.99	73.56	61.05
LViT-T	Hybrid	✓	100%	29.7	54.1	83.66	75.11	74.57	61.33

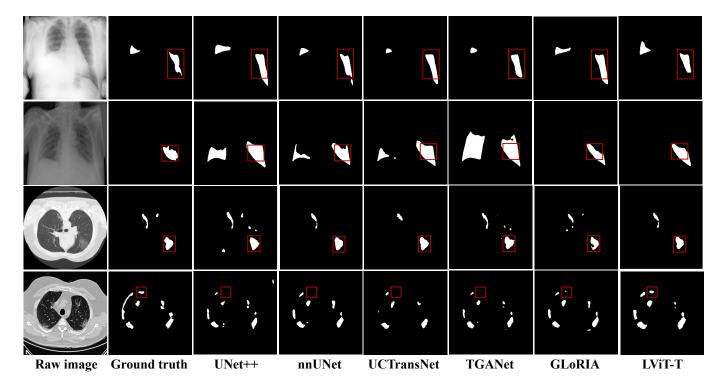


Fig. 4. Qualitative results on the QaTa-COV19 and the MosMedData+ datasets.

refers to the experimental results of the LViT-T model using 1/2 of the train set labels. And LViT-T means the **Tiny** version of LViT. The quantitative experimental results are listed in Table II. Experimental results on the QaTa-COV19 dataset show that LViT-TW/ LViT-T is able to achieve better performance than the previous SOTA method with smaller number of parameters and lower computational cost. In detail, LViT-T improves the Dice score by 3.24% and the mIoU score by 4.3% compared to the suboptimal nnUNet. It is also worth noting that LViT-T

still outperforms other SOTA methods even when only 1/4 of the training labels are used. Similarly, it can be seen that LViTT has a 2.54% higher Dice score and a 4.05% better mIoU score than LViT-TW. This also indicates that introducing text information is able to improve model performance effectively. A similar trend is observed for the MosMedData+ dataset. On the MosMedData+ dataset, compared to GLoRIA, LViT-T improves the Dice value by 2.15% and the mIoU value by 1.15%. Even LViT-TW and LViT-T (1/4) can achieve comparable

TABLE III ABLATION STUDY ON THE EFFECTIVENESS OF SUPERVISED COMPONENTS: DOWNVIT, UPVIT, PLAM, TEXT & SEMI-SUPERVISED COMPONENTS: EPI, TEXT, LOSS  $L_{LV}$  ON THE QATA-COV19 dataset.

Method	CNN	DownViT	UpViT	GAP/PLAM	GMP/PLAM	Text	EPI	${ m L}_{LV}$	Dice (%)	mIoU (%)
nnUNet	✓								80.42	70.81
	<b>√</b>	<b>√</b>							80.73	70.96
	$\checkmark$	$\checkmark$	$\checkmark$						80.85	71.12
	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$					81.03	71.29
LViT-T	$\checkmark$	$\checkmark$	$\checkmark$		$\checkmark$				80.92	71.21
	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$				81.12	71.37
		$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$			80.52	70.43
	$\checkmark$	$\checkmark$	$\checkmark$	✓	✓	$\checkmark$			83.66	75.11
	<b>√</b>	<b>√</b>	✓	✓	<b>√</b>				78.87	69.25
LViT-T (1/4)	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$			80.41	70.79
LVII-I (1/4)	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$		$\checkmark$		79.08	69.42
	$\checkmark$		80.67	71.08						
	$\checkmark$	80.95	71.31							
LViT-T (1/4, sup with $L_{LV}$ )	✓	✓	✓	<b>√</b>	✓	✓	✓	✓	80.98	71.30

TABLE IV
ABLATION STUDY ON DIFFERENT MODEL SIZES: LVIT-T, LVIT-S, LVIT-B. THE DICE AND IOU ARE IN 'MEAN±STD' FORMAT. THE STD STANDS FOR STANDARD DEVIATION IN THREE TIMES RUNS.

Model Size	Param(M)	Flops(G)	QaTa-COV19		MosMe	edData+
		• •	Dice (%)	mIoU (%)	Dice (%)	mIoU (%)
LViT-TW	28.0	54.0	81.12±1.9	71.37±2.4	$72.58 \pm 1.4$	60.40±0.7
LViT-SW	53.1	63.8	$81.54 \pm 1.7$	$71.91 \pm 2.0$	$72.75 \pm 1.2$	$60.52 \pm 0.6$
LViT-BW	69.8	70.4	$81.59 \pm 1.6$	$71.82 \pm 2.1$	$72.84 \pm 1.1$	$60.58 \pm 0.6$
LViT-T	29.7	54.1	83.66±0.8	75.11±1.4	$74.57 \pm 0.8$	61.33±0.5
LViT-S	54.8	63.9	$83.41 \pm 1.0$	$74.84 \pm 1.3$	$74.65 \pm 0.7$	$61.46 \pm 0.5$
LViT-B	71.5	70.5	$83.63 \pm 0.9$	$75.28 \pm 1.1$	$74.76 \pm 0.5$	$61.53 \pm 0.4$

performance to nnUNet and UCTransNet.

The qualitative results of our model and other state-of-the-art methods on the MosMedData+ and QaTa-COV19 datasets are shown in Figure 4, where four baseline methods are selected for comparison. Qualitative results shows that LViT-T has excellent semantic segmentation capabilities, especially when compared to other SOTA methods where the mis-segmentation phenomenon is greatly reduced. As can be seen from the red boxes in Figure 4, UNet++, nnUNet and TransUNet all have more severe mis-segmentation than LViT. It also shows the introduction of text information in our learning mechanism can better guide the training of the model, and consequently lead to more accurate segmentation. In addition, compared with different multimodal segmentation methods, LViT also has obvious advantages as can be shown in Table II and Figure 4. Thanks to the benefits brought by the integration of text and image information in the same encoder, LViT is more delicate in the segmentation boundary.

## E. Ablation Study

A series of ablation experiments are conducted to verify the performance of our LViT model, which is explored in the following four aspects.

1) Effectiveness of Proposed Components: We perform relevant ablation experiments on the effectiveness of both supervised components and semi-supervised components of our LViT model, and the relevant experimental results are presented in Table III. In full supervision, we explore the effectiveness of these four components, i.e., DownViT, UpViT, PLAM, and Text. The Text refers to the text information. And in the absence of text information, we utilize image features as the input of the transformer path. Experimental results illustrate that all of these components are effective, and the performance improvement brought by Text is the most significant. Furthermore, we also conduct ablation experiments on two attention modules (GAP and GMP) in PLAM. Our findings suggest that GAP/PLAM outperforms GMP/PLAM, possibly because GAP is better at integrating diverse information through global average pooling. It is worth noting that combining GAP and GMP yields better results than using them alone. To demonstrate our innovative points in semi-supervision, we explore the effectiveness of these three components, i.e., EPI, Text, and  $L_{LV}$ . From the experimental results, it can be seen that the improvements of EPI, Text, and  $L_{LV}$  are significant. Among them, by incorporating text annotation information, the Dice score and mIoU score are improved by 1.59% and by 1.66%

Method	Text Encoding	Text Format	Param(M)	Flops(G)	QaTa-COV19		MosMedData+	
				•	Dice (%)	mIoU (%)	Dice (%)	mIoU (%)
LViT-T	Text Encoder	Structured	84.9	101.8	83.34	74.76	74.29	61.18
LViT-T	Embedding Layer	Structured	29.7	54.1	83.66	75.11	74.57	61.33
LViT-T	Text Encoder	Unstructured	84.9	101.8	82.53	73.09	73.71	61.09
LViT-T	Embedding Layer	Unstructured	29.7	54.1	82.41	72.92	73.68	61.08

respectively when using 1/4 of the train set labels. And by introducing EPI mechanism, the semi-supervised performance of LViT is guaranteed to be comparable to the fully-supervised performance of U-Net. LViT (1/4) with EPI yields a 0.26% increase in the Dice score with text annotations and a 0.21% increase in the Dice score without text annotations. Finally, continuous improvements of model performance are also ensured by introducing  $L_{LV}$  on unlabeled data. On the other hand, for labeled data, we already have the accurate mask for supervised learning. Therefore, using  $L_{LV}$  on labeled data does not have significant benefits, as presented in the last row of Table III.

TABLE VI
ABLATION STUDY WITH LVIT-T ON DIFFERENT HYPER-PARAMETERS:
BATCH SIZE AND LEARNING RATE.

Hyper-Parameters		QaTa-	COV19	MosMedData+		
		Dice (%)	mIoU (%)	Dice (%)	mIoU (%)	
	16	82.72	73.96	73.98	61.10	
Batch Size	20	82.83	74.02	74.34	61.21	
	24	83.66	75.11	74.57	61.33	
	3e-4	83.66	75.11	74.52	61.31	
Learning Rate	1e-3	82.25	73.69	74.57	61.33	
	3e-3	82.20	73.53	74.46	61.27	

2) Ablation Study on Model Size: We conduct ablation experiments three times for model sizes to investigate the specific performance of LViT with different model sizes. Experiments are conducted on two datasets, QaTa-COV19 and MosMedData+, with six different model sizes, namely, LViT-TW/LViT-T, LViT-SW/LViT-S, LViT-BW/LViT-B, where "W" refers to without the text annotation, "T" refers to the tiny model, "S" refers to the small model, and "B" refers to the base model. In the original vision transformer (ViT) [39], each ViT module has L=12 transformer layers. And the differences between different versions of LViT are in the number of Transformer layers in the DownViT module and UpViT module, where LViT-TW/LViT-T has only 1 Transformer layer per ViT module, LViT-SW/LViT-S has 4 Transformer layers per ViT module, and LViT-BW/LViT-B has 6 Transformer layers per ViT module. Experimental results are reported in Table V. As observed in the table, it is worth noting that LViT with the text annotation has only 1.7M more parameters and 0.1G more computation than LViT-W, while the improvement of segmentation performance brought by the text information is significant. Besides, there is also an interesting observation that larger models do not consistently

improve accuracy over small models. We believe that it is related to the data distribution of the dataset. When the dataset distribution is relatively uniform, and the image is easy to segment, increasing the size of the model may not bring about a consistent improvement in performance. Conversely, if the dataset distribution presents notable differences and image segmentation is challenging, increasing the model size can lead to performance improvements. Nevertheless, it is noteworthy that as the model size increases, the model's performance jitter reduces, indicating the model becomes more robust.

3) Ablation Study on Hyper-Parameters: Ablation experiments are conducted on two aspects of hyper-parameters: Batch Size and Learning Rate. For Batch Size, we set 16, 20, and 24 on the QaTa-COV19 dataset and the MosMedData+ dataset. According to Table VI, LViT is optimal for batch size of 24 and learning rate of 3e-4 on the QaTa-COV19 dataset. And LViT is optimal for batch size of 24 and learning rate of 1e-3 on the MosMedData+ dataset. It is worth noting that the impact of hyper-parameters on the model performance is more considerable than the model size.

4) Ablation Study on Text Encoder and Embedding Layer: To further explore the encoding capabilities of text encoders and text embedding layers for text, as well as the differences in their applications to multimodal features, we conduct two sets of experiments for analysis. One set focuses on existing well-structured texts, while the other set focuses on poorly structured texts. We construct unstructured text annotations by randomly swapping the positions of phrases within the text. For example, we modify the description "Bilateral pulmonary infection, two infected areas, all left lung and middle lower right lung" to "All left lung and middle lower right lung, two infected areas, bilateral pulmonary infection". The experimental results, presented in Table V, reveal that utilizing a text encoder requires nearly three times as many parameters and nearly twice as much computation compared to using a text embedding layer. However, despite the increased complexity, the model performance does not improve and even decreases in wellstructured reports. This finding supports our decision to employ a text embedding layer in our LViT model. It is worth noting that the model performance with text embedding layer is slightly lower than that of a text encoder for poorly structured reports. We believe this discrepancy can be attributed to the better encoding ability and robustness of text encoders when dealing with the more diverse radiology reports. However, it is important to acknowledge that the resulting parameter and computational costs are not cost-effective.

5) Ablation Study on Semi-Supervision: Multiple semisupervised experiments are conducted to verify the model performance in semi-supervised learning. These experiments cover two different label ratios, i.e., 25% and 50%, to explore the performance changes under different label ratios. Additionally, experiments are conducted with and without text information. We compare our method with both the traditional SOTA semi-supervised medical image segmentation methods, such as DTC [50], PLCT [51], and MC-Net+ [52], and the multimodal methods, such as LAVT [27] and GLoRIA [49]. The experimental results are presented in Table VI, which demonstrate that our proposed LViT model achieves superior segmentation performance to other methods. This is attributed to the Exponential Pseudo-label Iteration mechanism and LV loss, regardless of whether text information is included in the pipeline or not.

TABLE VII
ABLATION STUDY ON SEMI-SUPERVISION WITH LV1T-T AND OTHER
METHODS ON THE QATA-COV19 DATASET

Method	Text	Label ratio	Dice (%)	mIoU (%)
DTC [50]		25%	76.07	66.04
DTC [50]	×	/-		
PLCT [51]	×	25%	76.65	66.71
MC-Net+ [52]	×	25%	76.93	67.02
LViT-TW (1/4)	×	25%	79.08	69.42
LAVT [27]	<b>√</b>	25%	77.08	67.21
GLoRIA [49]	$\checkmark$	25%	77.32	67.48
LViT-T (1/4)	$\checkmark$	25%	80.95	71.31
DTC [50]	×	50%	77.23	67.42
PLCT [51]	×	50%	77.66	68.04
MC-Net+ [52]	×	50%	77.91	68.47
LViT-TW (1/2)	×	50%	80.35	70.74
LAVT [27]	<b>√</b>	50%	77.96	68.53
GLoRIA [49]	$\checkmark$	50%	78.49	68.97
LViT-T (1/2)	$\checkmark$	50%	82.73	73.99

TABLE VIII
PERFORMANCE COMPARISON BETWEEN OUR METHOD (LViT) AND OTHER
METHODS ON THE ESO-CT DATASET.

Method	Param (M)	Flops (G)	Dice (%)	mIoU (%)
U-Net	14.8	50.3	66.75	56.31
U-Net++	74.5	94.6	66.70	56.59
nnUNet	19.1	412.7	68.38	56.10
TransUNet	105	56.7	65.94	55.78
LViT-TW	28.0	54.0	68.27	57.02
LViT-T	29.7	54.1	71.53	59.94

# F. Practical Application and Generalization Study on Esophageal CT Dataset

To explore the generalization performance of LViT and demonstrate how text annotations can help in practical scenarios, we conduct experiments on the ESO-CT dataset. This dataset is an esophageal cancer segmentation dataset collected by us, which consists of 286 cases. Each case contains mask and clinical information manually annotated by radiologist. In clinical information, the radiologist will divide the esophagus

into four sections (Cervical, Upper Thorax, Middle Thorax, Lower Thorax) from top to bottom, and then give the rough location of tumor. The rough location of tumor is utilized as the text input in LViT. We compare our LViT model with several SOTA methods, including UNet, UNet++, nnUNet, and TransUNet. As presented in Table. VIII, the performance of LViT-TW and nnUNet is comparable and superior to other methods. It should be noted that nnUNet has an extremely complex preprocessing for images, while LViT does not need it. In addition, with text innoations, LViT-T outperforms other methods by a large margin, thereby verifying the effectiveness and extensibility of LViT across different datasets. Therefore, we believe that with the approximate location (text annotations), the LViT model is able to segment the lesion area to form a more precise radiotherapy target area to better kill cancer cells.

## G. Interpretability Study

The interpretability study is performed on the QaTa-COV19 dataset to explore whether the LViT network can notice lesion regions and whether the introduction of text information can enhance the attention to lesion regions. In order to provide a more intuitive display of the changes in the region of interest of the model, we utilize GradCAM [53] to compare the activation for regions of attention. Figure 5 shows that UNet++, nnUNet, UCTransNet, and GLoRIA all have different degrees of misactivation regions. For example, in the case of total lung infection, these methods can only activate half of the lung region. In contrast, by introducing text information, our model can activate more regions, and the edge profile of the activated regions of interest is more consistent with the ground truth. Therefore, the localization of lesions can be learned by the text input.

Besides, to better explore the benefits of incorporating text information into the pipeline, we conduct more experiments on another two cases with more segmentation areas, as shown in Figure 6. We perform activation mapping in DownCNN1, DownCNN2, DownCNN3, DownCNN4, and DownViT1, respectively, where DownCNN1 and DownViT1 represent the first layer DownCNN and the first layer DownViT, respectively. The text information is input to the model through DownViT1, thus the difference in activation regions between DownViT1 and DownCNN1 can be approximated as the difference brought by the text information. It can be seen that the activation effect of the region of interest of DownViT1 is similar to that of DownCNN4. It is also worth noting that image feature of DownViT1 comes from DownCNN1, which failed to activate the lesion region but only activated the lung boundary. However, DownViT1 can directly activate the relevant lesion region by introducing the text information. And it indicates that the text information can effectively help locate lesion region in the lung, thus prompting the network to pay more attentions on the region indicated by the text information. The CAM output of LViT-TW and LViT-T shows the final activation difference caused by the text information. By comparing the regions of interest for LViT-TW/LViT-T, we believe that the text information can help reduce the probability of mis-segmentation.

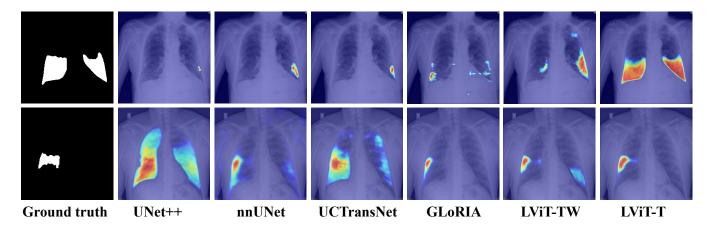


Fig. 5. Saliency map for interpretability study of different approaches on the QaTa-COV19 dataset. The language input of the first row is "Bilateral pulmonary infection, two infected areas, lower left lung and lower right lung". The language input of the second row is "Unilateral pulmonary infection, one infected area, middle left lung".

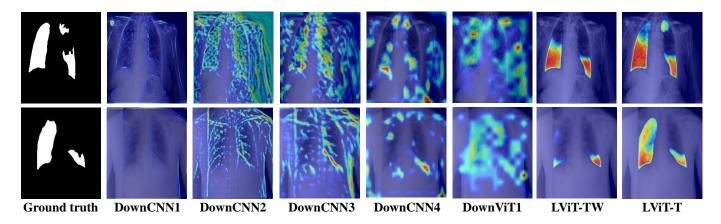


Fig. 6. Saliency map for interpretability study of different layers of LViT on the QaTa-COV19 dataset. The language input of the first row is "Bilateral pulmonary infection, three infected areas, all left lung and upper lower right lung". The language input of the second row is "Bilateral pulmonary infection, two infected areas, all left lung and lower right lung".

# V. CONCLUSION

In this paper, we propose a new vision-language medical image segmentation model LViT, which leverages medical text annotation to compensate for the quality deficiency in image data and guide to generate pseudo labels of improved quality in the semi-supervised learning. Multimodal medical segmentation datasets (image + text) are constructed to evaluate the performance of LViT, and experimental results show that our model has superior segmentation performance in both fullysupervised and semi-supervised settings. In addition, we present an example application on the diagnosis and treatment of earlystage esophageal cancer to demonstrate how text annotations can help in practical scenarios. Currently, the proposed model is a 2D Segmentation model. In our future work, we will extend our model to 3D and conduct experiments on more medical data to further verify its generality. Besides, in the current version of our LViT model, it is necessary to supply text inputs during the inference stage. Therefore, our another future work is to generate text annotation automatically according to the provided image information. As the text annotations are structured, we can transform the problem of text annotation generation into a classification problem in the future version

of LViT. This will enable us support inference either with or without text input, thereby enhancing the usability of our model.

## REFERENCES

- S. Chen, K. Ma, and Y. Zheng, "Med3d: Transfer learning for 3d medical image analysis," arXiv preprint arXiv:1904.00625, 2019.
- [2] J. Wang, L. Wei, L. Wang, Q. Zhou, L. Zhu, and J. Qin, "Boundary-aware transformers for skin lesion segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 206–216, Springer, 2021. I
- [3] Z. Zhu et al., "Lymph node gross tumor volume detection and segmentation via distance-based gating using 3d ct/pet imaging in radiotherapy," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 753–762, Springer, 2020. I
- [4] X. Zhang et al., "Self-supervised tumor segmentation through layer decomposition," arXiv preprint arXiv:2109.03230, 2021. I
- [5] Y. Li, L. Luo, H. Lin, H. Chen, and P. A. Heng, "Dual-consistency semisupervised learning with uncertainty quantification for covid-19 lesion segmentation from ct images," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 199–209, Springer, 2021. I, II-A
- [6] Z. Li et al., "Tfcns: A cnn-transformer hybrid network for medical image segmentation," in *International Conference on Artificial Neural Networks*, pp. 781–792, Springer, 2022. I
- [7] X. Yu, J. Wang, Q. Q. Hong, R. Teku, S. H. Wang, and Y. D. Zhang, "Transfer learning for medical images analyses: A survey," *Neurocomputing*, vol. 489, pp. 230–254, 2022. I

- [8] L. Yu, S. Wang, X. Li, C. W. Fu, and P. A. Heng, "Uncertainty-aware self-ensembling model for semi-supervised 3d left atrium segmentation," in *International Conference on Medical Image Computing and Computer-*Assisted Intervention, pp. 605–613, Springer, 2019. I
- [9] X. Feng, J. Yang, A. F. Laine, and E. D. Angelini, "Discriminative localization in cnns for weakly-supervised segmentation of pulmonary nodules," in *International conference on medical image computing and* computer-assisted intervention, pp. 568–576, Springer, 2017. I
- [10] J. B. Grill et al., "Bootstrap your own latent-a new approach to self-supervised learning," Advances in Neural Information Processing Systems, vol. 33, pp. 21271–21284, 2020. I, III-B
- [11] S. P. Morozov et al., "Mosmeddata: Chest ct scans with covid-19 related findings dataset," arXiv preprint arXiv:2005.06465, 2020. I, IV-A
- [12] J. Hofmanninger, F. Prayer, J. Pan, S. Röhrich, H. Prosch, and G. Langs, "Automatic lung segmentation in routine imaging is primarily a data diversity problem, not a methodology problem," *European Radiology Experimental*, vol. 4, no. 1, pp. 1–13, 2020. I, IV-A
- [13] A. Degerli, S. Kiranyaz, M. E. Chowdhury, and M. Gabbouj, "Osegnet: Operational segmentation network for covid-19 detection using chest x-ray images," *arXiv preprint arXiv:2202.10185*, 2022. I, IV-A
- [14] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE conference on* computer vision and pattern recognition, pp. 3431–3440, 2015. II-A, II-A
- [15] L. C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," *IEEE transactions on* pattern analysis and machine intelligence, vol. 40, no. 4, pp. 834–848, 2017. II-A
- [16] F. Xu et al., "Mrdff: A deep forest based framework for ct whole heart segmentation," Methods, 2022. II-A
- [17] F. Xu, L. Lin, D. Li, Q. Hong, K. Liu, et al., "A multi-resolution deep forest framework with hybrid feature fusion for ct whole heart segmentation," in 2021 IEEE International Conference on Bioinformatics and Biomedicine (BIBM), pp. 1119–1124, IEEE, 2021. II-A
- [18] W. Yuan, Y. Peng, Y. Guo, Y. Ren, and Q. Xue, "Dcau-net: dense convolutional attention u-net for segmentation of intracranial aneurysm images," Vis. Comput. Ind. Biomed. Art, vol. 5, no. 9, 2022. II-A
- [19] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*, pp. 234–241, Springer, 2015. II-A, II
- [20] Z. Zhou, M. M. Rahman Siddiquee, N. Tajbakhsh, and J. Liang, "Unet++: A nested u-net architecture for medical image segmentation," in *Deep learning in medical image analysis and multimodal learning for clinical decision support*, pp. 3–11, Springer, 2018. II-A, II
- [21] Y. Li et al., "Gt u-net: A u-net like group transformer network for tooth root segmentation," in *International Workshop on Machine Learning in Medical Imaging*, pp. 386–395, Springer, 2021. II-A, II-B
- [22] Q. Liu, L. Yu, L. Luo, Q. Dou, and P. A. Heng, "Semi-supervised medical image classification with relation-driven self-ensembling model," *IEEE* transactions on medical imaging, vol. 39, no. 11, pp. 3429–3440, 2020. II-A
- [23] Z. Li, W. Chen, Z. Wei, X. Luo, and B. Su, "Semi-wtc: A practical semi-supervised framework for attack categorization through weight-task consistency," arXiv preprint arXiv:2205.09669, 2022. II-A
- [24] Y. Xia et al., "Uncertainty-aware multi-view co-training for semi-supervised medical image segmentation and domain adaptation," Medical Image Analysis, vol. 65, p. 101766, 2020. II-A
- [25] A. Radford et al., "Learning transferable visual models from natural language supervision," in *International Conference on Machine Learning*, pp. 8748–8763, PMLR, 2021. II-B, II
- [26] W. Kim, B. Son, and I. Kim, "Vilt: Vision-and-language transformer without convolution or region supervision," in *International Conference* on Machine Learning, pp. 5583–5594, PMLR, 2021. II-B, II
- [27] Z. Yang, J. Wang, Y. Tang, K. Chen, H. Zhao, and P. H. Torr, "Lavt: Language-aware vision transformer for referring image segmentation," in 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 18134–18144, 2022. II-B, II-B, III-A3, II, IV-E5, VII
- [28] H. Ding, C. Liu, S. Wang, and X. Jiang, "Vision-language transformer and query generation for referring segmentation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 16321–16330, 2021. II-B
- [29] H. Ding, C. Liu, S. Wang, and X. Jiang, "Vlt: Vision-language transformer and query generation for referring segmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–16, 2022. II-B, II-B,

- [30] W. Yin, Y. Liu, C. Shen, A. V. D. Hengel, and B. Sun, "The devil is in the labels: Semantic segmentation from sentences," arXiv preprint arXiv:2202.02002, 2022. II-B
- [31] J. Xu et al., "Groupvit: Semantic segmentation emerges from text supervision," arXiv preprint arXiv:2202.11094, 2022. II-B
- [32] R. Bhalodia et al., "Improving pneumonia localization via cross-attention on medical images and reports," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 571–581, Springer, 2021. II-B
- [33] P. Müller, G. Kaissis, C. Zou, and D. Rückert, "Joint learning of localized representations from medical images and reports," arXiv preprint arXiv:2112.02889, 2021. II-B
- [34] N. K. Tomar, D. Jha, U. Bagci, and S. Ali, "Tganet: Text-guided attention for improved polyp segmentation," arXiv preprint arXiv:2205.04280, 2022. II-B, II
- [35] F. Wang et al., "Residual attention network for image classification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3156–3164, 2017. II-C
- [36] S. Woo, J. Park, J. Y. Lee, and I. S. Kweon, "Cbam: Convolutional block attention module," in *Proceedings of the European conference on computer vision (ECCV)*, pp. 3–19, 2018. II-C, III-A3
- [37] C. Huang et al., "A deep segmentation network of multi-scale feature fusion based on attention mechanism for ivoct lumen contour," *IEEE/ACM Transactions on computational biology and bioinformatics*, vol. 18, no. 1, pp. 62–69, 2020. II-C
- [38] A. Vaswani et al., "Attention is all you need," Advances in neural information processing systems, vol. 30, 2017. II-C
- [39] A. Dosovitskiy et al., "An image is worth 16x16 words: Transformers for image recognition at scale," arXiv preprint arXiv:2010.11929, 2020. II-C, III-A2, IV-E2
- [40] Y. Li et al., "Agmb-transformer: Anatomy-guided multi-branch transformer network for automated evaluation of root canal therapy," IEEE Journal of Biomedical and Health Informatics, 2021. II-C
- [41] D. Shan, Z. Li, W. Chen, Q. Li, J. Tian, and Q. Hong, "Coarse-to-fine covid-19 segmentation via vision-language alignment," arXiv preprint arXiv:2303.00279, 2023. II-C
- [42] J. Devlin, M. W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," arXiv preprint arXiv:1810.04805, 2018. III-A2
- [43] O. Oktay et al., "Attention u-net: Learning where to look for the pancreas," arXiv preprint arXiv:1804.03999, 2018. II
- [44] F. Isensee et al., "nnu-net: a self-configuring method for deep learning-based biomedical image segmentation," Nature methods, vol. 18, no. 2, pp. 203–211, 2021. II
- [45] J. Chen et al., "Transunet: Transformers make strong encoders for medical image segmentation," arXiv preprint arXiv:2102.04306, 2021. II
- [46] H. Cao et al., "Swin-unet: Unet-like pure transformer for medical image segmentation," arXiv preprint arXiv:2105.05537, 2021. II
- [47] H. Wang, P. Cao, J. Wang, and O. R. Zaiane, "Uctransnet: Rethinking the skip connections in u-net from a channel-wise perspective with transformer," arXiv preprint arXiv:2109.04335, 2021. II
- [48] Y. Zhang, H. Jiang, Y. Miura, C. D. Manning, and C. P. Langlotz, "Contrastive learning of medical visual representations from paired images and text," arXiv preprint arXiv:2010.00747, 2020. II
- [49] S. C. Huang, L. Shen, M. P. Lungren, and S. Yeung, "Gloria: A multimodal global-local representation learning framework for labelefficient medical image recognition," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3942–3951, 2021. II, IV-E5, VII
- [50] X. Luo, J. Chen, T. Song, and G. Wang, "Semi-supervised medical image segmentation through dual-task consistency," in *Proceedings of* the AAAI Conference on Artificial Intelligence, vol. 35, pp. 8801–8809, 2021. IV-E5, VII
- [51] K. Chaitanya, E. Erdil, N. Karani, and E. Konukoglu, "Local contrastive loss with pseudo-label based self-training for semi-supervised medical image segmentation," *Medical Image Analysis*, p. 102792, 2023. IV-E5, VII
- [52] Y. Wu, Z. Ge, D. Zhang, M. Xu, L. Zhang, Y. Xia, and J. Cai, "Mutual consistency learning for semi-supervised medical image segmentation," *Medical Image Analysis*, vol. 81, p. 102530, 2022. IV-E5, VII
- [53] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," in *Proceedings of the IEEE international* conference on computer vision, pp. 618–626, 2017. IV-G