

TFCNs: A CNN-Transformer Hybrid Network for Medical Image Segmentation

Zihan Li^{#1}, Dihan Li^{#1}, Cangbai Xu¹, Weice Wang¹, Qingqi Hong^{*1,4}[0000-0002-9996-6870], Qingde Li²[0000-0001-5998-7565], and Jie Tian³[0000-0003-0498-0432]

¹ Xiamen University, Xiamen, 361005, China

{zihanli,dihanli,cangbaixu,wangweice}@stu.xmu.edu.cn hongqq@xmu.edu.cn

² University of Hull, Hull, HU6 7RX, UK Q.li@hull.ac.uk

³ Chinese Academy of Sciences, Beijing, 100190, China tian@ieee.org

⁴ State Key Laboratory of Virtual Reality Technology and Systems, Beihang University, China

Abstract. Medical image segmentation is one of the most fundamental tasks concerning medical information analysis. Various solutions have been proposed so far, including many deep learning-based techniques, such as U-Net, FC-DenseNet, etc. However, high-precision medical image segmentation remains a highly challenging task due to the existence of inherent magnification and distortion in medical images as well as the presence of lesions with similar density to normal tissues. In this paper, we propose TFCNs (Transformers for Fully Convolutional denseNets) to tackle the problem by introducing ResLinear-Transformer (RL-Transformer) and Convolutional Linear Attention Block (CLAB) to FC-DenseNet. TFCNs is not only able to utilize more latent information from the CT images for feature extraction, but also can capture and disseminate semantic features and filter non-semantic features more effectively through the CLAB module. Our experimental results show that TFCNs can achieve state-of-the-art performance with dice scores of 83.72% on the Synapse dataset. In addition, we evaluate the robustness of TFCNs for lesion area effects on the COVID-19 public datasets. The Python code will be made publicly available on <https://github.com/HUANGLIZI/TFCNs>.

Keywords: Medical image segmentation · CNN-Transformers · Attention mechanism

1 Introduction

Medical image segmentation plays a critical role in clinical diagnosis and assisting doctors to evaluate patient’s reactions to treatment. Various algorithms based on convolutional neural networks (CNNs) [9] have been applied to image segmentation. And with a U-shaped network design, U-Net [15] has achieved

means equal contribution, * means corresponding author

great success in various medical imaging applications. Following this technical route, many algorithms have been developed for medical image and volume segmentation, such as U-Net++ [27]. In order to solve the degradation problem, He et al. proposed ResNets [5], which aims to simplify very deep networks by introducing a residual block that sums two input signals. Then a new CNN architecture called DenseNets [7] has been developed by the composition of dense blocks and pooling operations. In FC-DenseNet [23], the up-sampling path was introduced to restore the input resolution. Recently, inspired by the great success of Transformers in the field of natural language processing (NLP) [3], researchers have tried to introduce Transformers into the field of computer vision [10,26]. Vision transformer (ViT) [4] has been proposed to achieve object detection tasks.

Currently, there are two problems: 1). As shown in Fig.1, since the convolution operation collects information by layer, which leads to too much focus on local feature information. In the field of medical image segmentation, the lack of global information often leads to false category of segmentation. Therefore, a visual transformer was introduced, which can reflect complex spatial transformations and long-distance feature dependencies, which are regarded as global representations. Currently, although Chen et al. proposed Transunet [2] to solve problems such as lack of high-level details. However, we found that the direct feeding of CNN-style features into the transformer for recoding tends to bring limited improvement. 2). In U-shape networks, the skip connection between encoder and decoder is crucial. However, semantic-independent features tend to be fed to the decoder with direct transmission, which will interfere with image segmentation. Our main motivation is how to preserve image features more completely.

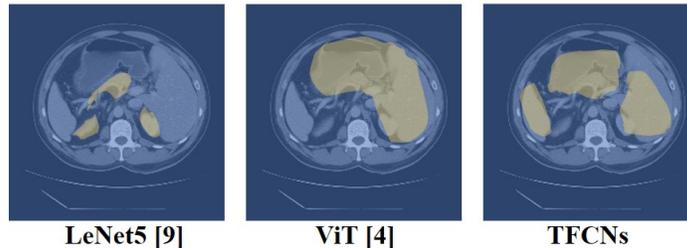


Fig. 1. Class activation maps in LeNet5 [9], ViT [4] and TFCNs by using the CAM method [25]. In which we set a fixed value as the activation intensity threshold.

To tackle the problem 1), we utilize DenseBlocks to facilitate the propagation of feature information to the transformer part while adding a residual structure to the MLP to further effectively preserve the global representation information. To tackle the problem 2), we decide to fuse spatial and channel attention on the original skip connection to transmit information more efficiently. Our main contributions are as follows:

1. A new deep neural network framework (TFCNs) is proposed, to the best of our knowledge, which is the first network to introduce Transformers into FC-DenseNet and improve the internal structure of Transformers.

2. A general attention module CLAB (Convolutional Linear Attention Block) is proposed to improve segmentation performance, which includes two types of attentions: (a) attention over the spatial extent of image and (b) attention over the CNN-style feature channels.

2 Related Work

In the field of semantic segmentation, FCNs [13] innovatively proposed a model structure in the form of encoder-decoder. And to solve the problem of information loss in the encoding process, it utilized the form of residual connection to combine the encoding process. In addition to the semantic segmentation of real objects, more and more attention has also been paid to medical image segmentation. Based on FCNs, U-Net [15] was proposed and applied to medical image segmentation. This model makes use of a mutually symmetrical encoding-decoding design. Another example was FC-DenseNet [8], where they extended the work in DenseNets [7] by introducing the DenseBlock in the process of upsampling, which not only alleviated the problem of dimensional explosion in the deep encoder but also retained contextual information better. Some Transformer-based methods have also been proposed for semantic segmentation, object detection, and instance segmentation, such as SETR, DETR [24,1]. Inspired by the previous breakthroughs, TransUNet [2] embedded the Transformers in the down-sampling process to extract the information in the original image. More recently, a Gated Axial-Attention model was proposed in MedT [17] to extend some existing attention-based schemes. There are also other variants to the Transformers such as the Swin Transformer [12], which utilize a sliding window to limit self-attention calculations to non-overlapping partial windows.

3 Method

3.1 Overall structure of TFCNs

As described in the first section, the conventional U-shaped structured network lacks global contextual information to perform high-precision medical image segmentation. Given this, we propose TFCNs (Fig.2), which takes FC-DenseNet [23] as the backbone network, with an RL-Transformer Layer being added to the encoder to enhance the segmentation capability of the network. In addition, CLAB (Convolutional Linear Attention Block) in the skip connection part is introduced to enhance the spatial recovery of the focused segmentation region. The CNN-Transformer hybrid model acts as an encoder and CLAB as the upper and lower connecting part between DenseBlock and Transition-Down, which helps to filter non-semantic features. Compared with TransUNet [2], TFCNs not only replaces the traditional convolutional layers with Dense Blocks, but also changes the feature encoding method. The details of each part of the structure will be described in the next two subsections. More specifically, the RL-Transformer (ResLinear-Transformer) is described in section 3.2, the CLAB (Convolutional Linear Attention Block) is described in section 3.3.

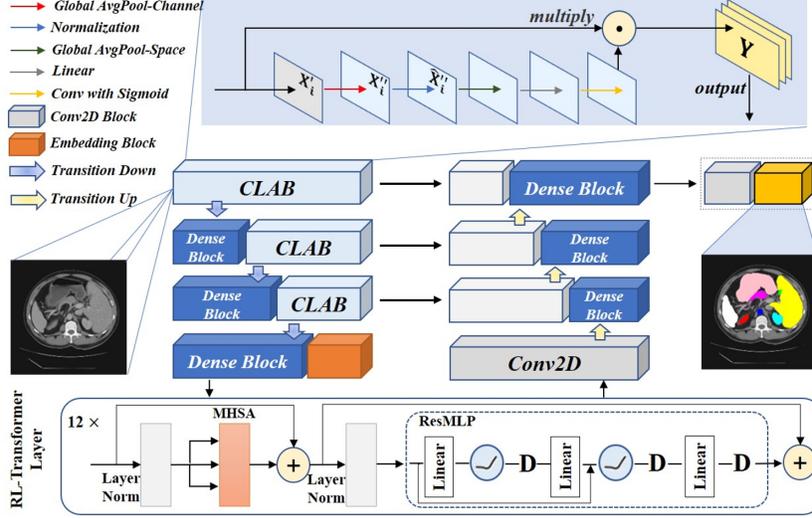


Fig. 2. Overview of the proposed TFCNs. RL-Transformer module at the end of encoder gives access to a receptive field containing images and CLAB modules are dedicated to filtering non-semantic features by including spatial and channel attention.

3.2 RL-Transformer

Referring to the ViT [4] implementation, we propose the ResLinear-Transformer (RL-Transformer) and apply it to the encoder of TFCNs. Of particular note is that the 2D patch x_P^i is expanded linearly and its projection is mapped to the D dimension using a trainable linear layer, as shown in Eqn.1. The output of this projection is called Patch Embedding.

$$z_0 = [x_{class}; x_P^1 E; x_P^2 E; \dots; x_P^N E] + E_{pos} \quad (1)$$

where E is the patch embedding projection which is located before entering the RL-transform layer and E_{pos} is the position embedding. RL-Transformer Layer consists of alternating ResMLP blocks and L layers of Multi-Head Self-Attention (MHSA). The expressions are shown in Eqn.2 and Eqn.3, where LN denotes Layer Norm.

$$z'_\ell = MHSA(LN(z_{\ell-1})) + z_{\ell-1} \quad (2)$$

$$z_\ell = ResMLP(LN(z'_\ell)) + z'_\ell \quad (3)$$

The RL-Transformer encoder consists of alternating multi-headed self-attentive layers and ResMLP blocks. As shown in Figure 2, our proposed ResMLP consists of two GELU [6] nonlinear layers, three Linear layers, and three dropout layers alternating with a residual connection to the source input before the second

GELU [6] layer. ResMLP can be expressed in Eqn.4 and Eqn.5 shown below:

$$z_\ell'' = LN(z_\ell') \quad (4)$$

$$y = z_\ell'' + L(\alpha GELU(L(z_\ell''))) \quad (5)$$

where GELU represents the GELU[6] nonlinear layer, L represents the Linear layer, and α represents the associated weight parameters which is a learned parameter. Finally, the state of the sequence at the output of the RL-Transformer encoder is utilized as the image features.

3.3 CLAB (Convolutional Linear Attention Block)

Inspired by TTD (Test-Time Dropout) [20] and TTA (Test-Time Augmentation) [11], we propose **Convolutional Linear Attention Block (CLAB)**. Nowadays TTD [20] can utilize dropout layers in the reasoning process and generate multiple predictions for each data instance. Compared with CUAB [19], we change the order of extracting channel and spatial attention, as indicated by the experimental results of CBAM [21]. In addition, the normalization operation is utilized between the two attention operations to accelerate the training convergence process. Finally, inspired by the local field of view, the weak influence outside the local area is directly reduced to zero influence by using the linear layer.

As shown in Fig.2, the global average pooling layer with the linear layer is added between the two convolutional layers. In CLAB, we first use 1×1 convolutional layers, where each convolutional layer has K kernels to generate the sequence $X_i' = \mathbb{R}^{H \times W \times K}$, $i \in \{1, \dots, N\}$. Then a global average pooling operation is performed on X_i' in channel dimensions to obtain X_i'' , and a normalization operation is performed on X_i'' to homogenize the data to obtain $\hat{X}_i'' = \frac{X_i' - \mu_B}{\sqrt{\sigma_B^2 + \varepsilon}}$ to accelerate the convergence and accuracy improvement of the training process, where μ_B denotes the mean of a batch, σ_B denotes the standard deviation of the batch, and ε is a minimal positive value to ensure that the denominator is not zero. All $\hat{X}_i'' = \mathbb{R}^{H \times W}$, $i \in \{1, \dots, N\}$, are concatenated to form $X_m = \mathbb{R}^{H \times W \times N}$, which is then input sequentially into a submodule consisting of a global average pooling layer (with respect to dimensionality), a linear layer and a convolutional layer with sigmoid. The result is then multiplied with the source input to obtain the final output feature Y . We later perform ablation analysis for CLAB as well, and the experimental results show that CLAB has a significant effect on the model segmentation performance improvement.

4 Experimental Results and Discussion

4.1 Implementation Details

For all experiments, we perform a simple data augmentation, i.e., by performing random rotation or flipping operations. The optimizer chosen is the SGD optimizer, with a momentum of 0.9 and a weight decay of 1e-4. The learning rate

selected for the experiment for our method is 0.005 and is set to 0.001 for other models. The learning rate is made to decay after 30,000 iterations. The batch size is set to 3 for Segtran [24], and 12 for all other models. The epoch is set to 150 on all datasets.

For Patch Size, it is worth noticing from Table 1 that the performance of the model is optimal when the patch size is set to 16. When the patch size is set to 8, the included area is too small, so some relatively large organs such as liver or kidney will be divided into many different patches for encoding, which splits some important information. This makes it difficult for the decoder to perform well, thereby reducing the performance of the model. When the patch size is set to 32, because the coding area is relatively large, it contains many interfering information that makes the model to misjudge. Although the CLAB module functions as a filter, the remaining redundant information is still greater than when the patch size is 16, so it also affects the judgments of the model.

Table 1. Ablation study on different patch sizes in transformer(Dice score% and Hausdorff distance in mm and Jaccard score%).

Patch_Size	Dice \uparrow	Hd95 \downarrow	Jaccard \uparrow
8	78.79	38.11	65.95
16	83.72	17.26	72.78
32	80.00	24.41	67.84

Moreover, in our experiment, the combination of the Dice coefficient and cross-entropy is taken as the loss function for all methods. And the weight of these two components is 0.5.

4.2 Comparison with other SOTA methods

We conduct our experiments on Synapse[‡] dataset and COVID-19[§] dataset. The experimental results are analyzed by taking the Dice coefficient, Hausdorff distance, and Jaccard coefficient as evaluation factors. All the State-of-the-art methods are implemented using original paper without any deviations.

Synapse Dataset

As shown in Table 2, these methods including Transformers, i.e. TransUNet [2], Segtran [24], and our TFCNs achieves better performance when compare with other methods. Especially, our method has a slight lead in terms of Dice coefficient and Jaccard coefficient when compared to TransUNet [2] and Segtran [24]. We believe this improvement is brought by the Dense Block we add, which is able to enhance the transmission of semantic information in the main pipeline.

From Table 2, it can be observed that our method achieves more accurate results than other approaches on some organs that are difficult to be segment,

[‡] <https://www.synapse.org/#!/Synapse:syn3193805/wiki/217789>

[§] <https://aistudio.baidu.com/aistudio/datasetdetail/34221>

Table 2. Performance comparison between our method and other state-of-the-art methods on Synapse dataset (Dice score% and Hausdorff distance in mm and Jaccard score%, and Dice score % for each organ). Avg means average result of all testing cases and the Dice coefficient% on each organ.

Method	Dice(avg)↑	Hd95(avg)↓	Jaccard(avg)↑	Aorta↑	Gallbladder↑	Kidney(L)↑	Kidney(R)↑	Liver↑	Pancreas↑	Spleen↑	Stomach↑
U-Net [15]	81.81	26.81	71.21	86.99	67.31	88.64	82.81	94.15	60.99	90.50	83.10
Fc-DenseNet [23]	81.62	21.83	70.62	86.68	66.31	88.14	82.07	95.26	61.73	92.11	80.68
AttU-Net [14]	81.05	29.09	70.71	89.63	67.05	88.46	77.08	94.52	56.89	92.13	82.69
Res-UNet [22]	78.33	58.66	66.61	86.16	59.63	86.55	83.93	94.49	48.94	86.96	79.96
U-Net++ [27]	81.60	28.31	70.48	88.15	67.29	86.31	83.62	94.00	61.71	89.51	82.18
DDANet [16]	79.60	21.29	67.87	83.74	64.93	88.93	83.83	94.99	51.99	90.68	77.69
TransUNet [2]	82.33	19.88	71.18	88.50	62.96	91.23	90.03	94.90	59.90	90.53	80.59
Segtran [24]	83.02	14.73	72.68	87.10	62.88	92.66	89.94	95.47	61.76	91.47	83.40
TFCNs	83.72	17.26	72.78	89.69	67.75	90.11	88.30	94.74	64.08	92.22	82.84

such as the Gallbladder and Pancreas. Since these tiny organs occupy a relatively small proportion in the original image, other approaches are easy to wrap other interfering information when extracting features from the areas containing these tiny organs. Conversely, through the utilization of the designed CLAB module, our model is able to pay more attention to these tiny organs themselves instead of those irrelevant areas. Meanwhile, for other organs, the results achieved by our method are also in the top 5.

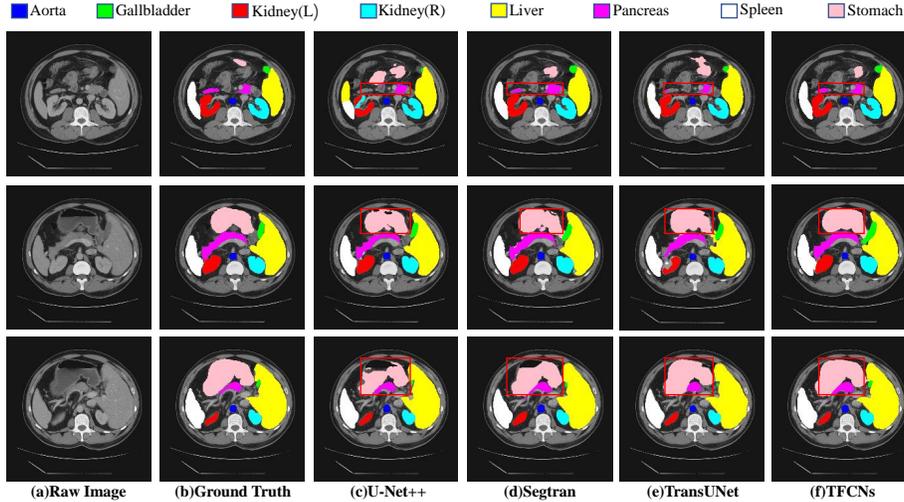


Fig. 3. Qualitative results on Synapse Dataset. All columns respectively represent: (a) Raw Image; (b) Ground Truth; (c) U-Net++ [27]; (d) Segtran [24]; (e) TransUNet [2]; (f) TFCNs. Top row describes corresponding color of each organ in raw image.

Analysis of segmentation for COVID-19 infected areas

The performance of our model at the fine-grained feature is explored using the second dataset, i.e. the COVID-19 dataset. Since lesion features normally

are more refined and scattered than organs which cover a large proportion of the image. Table 3 indicates that the capability of our method on the fine-grained target also reaches the SOTA level.

Table 3. Performance comparison between our method and other state-of-the-art methods on Covid-19 dataset (Dice score% and Hausdorff distance in mm and Jaccard score%). Avg means result average of all testing cases.

Method	Dice(Avg) \uparrow	Hd95(Avg) \downarrow	Jaccard(Avg) \uparrow
Segtran [24]	75.35	41.18	60.35
Res-UNet [22]	73.53	47.54	59.86
DDANet [16]	75.52	39.36	60.52
U-Net [15]	73.96	45.19	59.99
FC-DenseNet [23]	71.13	54.72	55.52
AttU-net [14]	74.70	48.36	60.26
TransUNet [2]	72.19	52.51	57.22
TFCNs	75.55	37.32	60.74

Combining with the results visualized in Fig.4, it can be seen that DDANet [16], which is very close to our method in Dice coefficient and Jaccard coefficient, is very prone to under-segment, indicating that its segmented strategy may be choosing to ignore the pixels that are hard to distinguish. This strategy avoids the decrease inaccuracy caused by erroneous prediction, but it can easily cause false-negative that should be as least as possible in the diagnosis. TransUNet [2] shows the over-segmentation in all rows. This may be the result of excessive utilization of contextual information in the local area which is suitable for continuous and regular objects such as organs, but these pieces of information are still too coarse for the lesion features. However, we solve this problem by adding the CLAB module.

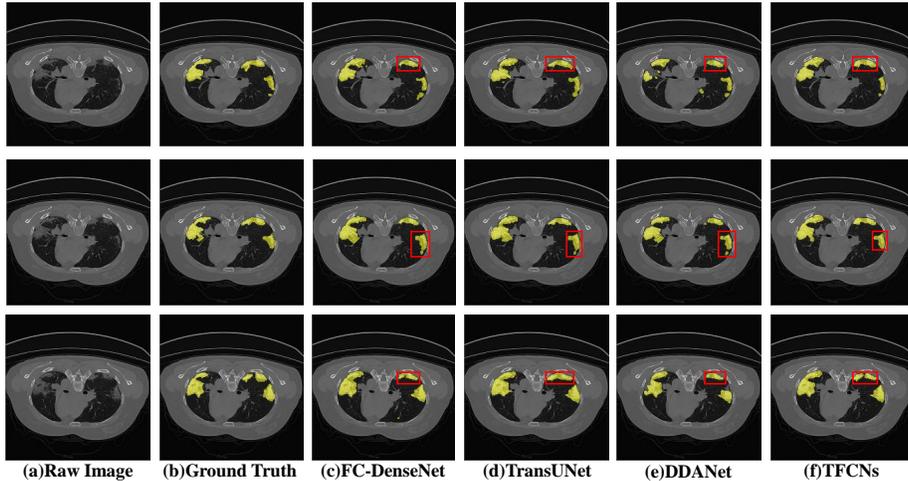


Fig. 4. Qualitative results on COVID-19 dataset. All columns respectively represent: (a) Raw Image; (b) Ground Truth; (c) FC-DenseNet [23]; (d) TransUNet [2]; (e) DDANet [16]; (f) TFCNs. Infected area in lung are labeled in yellow.

4.3 Interpretability Analysis

Interpretability analysis is conducted on the COVID-19 dataset using Class Activation Mapping (CAM) [25] to explore what causes the model to make the decision on every pixel and which area of feature map the model pays the most attention to when it segments the infected area. As shown in Fig.5, TFCNs focuses its attention on the lung area and pay least attention in the area surrounding the lung which means our model can act as an expert that focuses on the lung area at the beginning, and then pays more and more attention to the infected area gradually (in Fig.5, the lung area is light red, and the infected area is dark red) when it predicts.

As for TransUNet [2], although it can also focus its attention on the infected area when predicting, it disperses other attention to the entire image instead of gathering it to the lungs. This reflects the importance of our CLAB module because it relocates the attention of the model to the key area.

Moreover, even though DDANet [16] achieves good results, it ignores most areas in the image except the infected area, indicating that it does not have the concept of lungs and only makes some mechanical predictions based on ground truth which will make it very difficult to predict at the edge of the infected area, resulting in a large number of false-negatives.

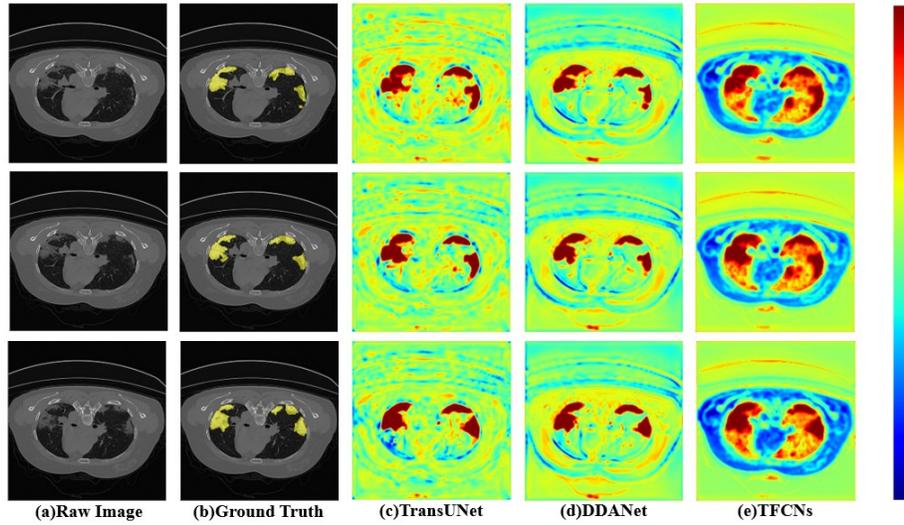


Fig. 5. Heat map for interpretability analysis of different approaches on COVID-19 dataset. All columns respectively represents: (a) Raw Image; (b) Ground Truth; (c) TransUNet [2]; (d) DDANet [16]; (e) TFCNs. The colormap in right presents the degree of attention which increases from bottom to top. Infected area in lung are labeled using yellow in ground truth.

4.4 Ablation Studies

Effectness of ResMLP

It can be seen from Table 4 that the Dice coefficient increases by 1.87% after using ResMLP. In addition, the Hausdorff distance drops by 8.38mm, indicating that ResMLP makes the semantic features in the Transformers propagate more complete, thereby increasing the contextual information extracted by the encoder, which promotes the performance of the overall structure.

Table 4. Ablation study on verifying the effectness of ResMLP (Dice score% and Hausdorff distance in mm and Jaccard score%).

	Dice↑	Hd95↓	Jaccard↑
ResMLP	83.72	17.26	72.78
MLP	81.85	25.64	70.57

Type of Attention Block

As shown in Table 5, the performance of the model is greatly improved after using the attention blocks (no matter what type it is), which means these attention blocks play a critical role in removing irrelevant and redundant information at skip connections. Moreover, it can be seen that after using the CLAB module we designed, the performance of the model is continuously improved, which demonstrates that our CLAB module can accurately locate the area containing more effective information in the feature map.

Table 5. Ablation study on different types of attention block in skip connection (Dice score% and Hausdorff distance in mm and Jaccard score%).

Type of Attention Block	Dice↑	Hd95↓	Jaccard↑
None	80.16	28.40	68.53
CUAB[19]	81.83	25.68	70.72
CLAB	83.72	17.26	72.78

5 Conclusion

To improve the performance of medical image segmentation, in this paper, we propose TFCNs based on Transformer [18] and FC-DenseNet [23]. And the internal structure of Transformers is modified by introducing residual connections to form RL-Transformer. This change can help enhance the receptive field and improve the coding ability of the model. In addition, a common module—CLAB, which is mainly composed of global average pooling and linear mapping, is designed in the network to filter out non-semantic features. Experimental results show that TFCNs which is the best among the baselines achieves a score of

83.72% on the Dice coefficient and a score of 72.78% on the Jaccard coefficient in terms of the Synapse dataset. The experiments are also conducted on the COVID-19 public dataset, and results show that TFCNs also has the state-of-the-art performance.

6 Acknowledge

This work was supported in part by the Natural Science Foundation of Fujian Province of China (No. 2020J01006), the National Natural Science Foundation of China (No. 61502402), and the Open Project Program of State Key Laboratory of Virtual Reality Technology and Systems, Beihang University (No.VRLAB2022 AC04).

References

1. Carion, N., et al.: End-to-end object detection with transformers. In: European Conference on Computer Vision. pp. 213–229. Springer (2020)
2. Chen, J., Lu, Y., Yu, Q., Luo, X., Adeli, E., Wang, Y., Lu, L., Yuille, A.L., Zhou, Y.: Transunet: Transformers make strong encoders for medical image segmentation. arXiv preprint arXiv:2102.04306 (2021)
3. Chowdhury, G.G.: Natural language processing. *Annual review of information science and technology* **37**(1), 51–89 (2003)
4. 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)
5. He, K., Zhang, X., Ren, S., Sun, J.: Identity mappings in deep residual networks. In: European conference on computer vision. pp. 630–645. Springer (2016)
6. Hendrycks, D., Gimpel, K.: Gaussian error linear units (gelus). arXiv preprint arXiv:1606.08415 (2016)
7. Iandola, F., Moskewicz, M., Karayev, S., Girshick, R., Darrell, T., Keutzer, K.: Densenet: Implementing efficient convnet descriptor pyramids. arXiv preprint arXiv:1404.1869 (2014)
8. Jégou, S., Drozdal, M., Vazquez, D., Romero, A., Bengio, Y.: The one hundred layers tiramisu: Fully convolutional densenets for semantic segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition workshops. pp. 11–19 (2017)
9. LeCun, Y., Bottou, L., Bengio, Y., Haffner, P.: Gradient-based learning applied to document recognition. *Proceedings of the IEEE* **86**(11), 2278–2324 (1998)
10. Li, Z., Li, Y., Li, Q., Zhang, Y., Wang, P., Guo, D., Lu, L., Jin, D., Hong, Q.: Lvit: Language meets vision transformer in medical image segmentation. arXiv preprint arXiv:2206.14718 (2022)
11. Litjens, G., Toth, R., van de Ven, W., Hoeks, C., Kerkstra, S., van Ginneken, B., Vincent, G., Guillard, G., Birbeck, N., Zhang, J., et al.: Evaluation of prostate segmentation algorithms for mri: the promise12 challenge. *Medical image analysis* **18**(2), 359–373 (2014)

12. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 10012–10022 (2021)
13. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3431–3440 (2015)
14. Oktay, O., et al.: Attention u-net: Learning where to look for the pancreas. arXiv preprint arXiv:1804.03999 (2018)
15. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: International Conference on Medical image computing and computer-assisted intervention. pp. 234–241. Springer (2015)
16. Tomar, N.K., Jha, D., Ali, S., Johansen, H.D., Johansen, D., Riegler, M.A., Halvorsen, P.: Ddanet: Dual decoder attention network for automatic polyp segmentation. In: International Conference on Pattern Recognition. pp. 307–314. Springer (2021)
17. Valanarasu, J.M.J., et al.: Medical transformer: Gated axial-attention for medical image segmentation. arXiv preprint arXiv:2102.10662 (2021)
18. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. *Advances in neural information processing systems* **30** (2017)
19. Wang, C.S., Su, F.Y., Lee, T.L.M., Tsai, Y.S., Chiang, J.H.: Cuab: Convolutional uncertainty attention block enhanced the chest x-ray image analysis. arXiv preprint arXiv:2105.01840 (2021)
20. Wang, G., Li, W., Aertsen, M., Deprest, J., Ourselin, S., Vercauteren, T.: Aleatoric uncertainty estimation with test-time augmentation for medical image segmentation with convolutional neural networks. *Neurocomputing* **338**, 34–45 (2019)
21. Woo, S., Park, J., Lee, J.Y., Kweon, I.S.: Cbam: Convolutional block attention module. In: Proceedings of the European conference on computer vision (ECCV). pp. 3–19 (2018)
22. Xiao, X., Lian, S., Luo, Z., Li, S.: Weighted res-unet for high-quality retina vessel segmentation. In: 2018 9th international conference on information technology in medicine and education (ITME). pp. 327–331. IEEE (2018)
23. Zhang, R., Zhao, L., Lou, W., Abrigo, J.M., Mok, V.C., Chu, W.C., Wang, D., Shi, L.: Automatic segmentation of acute ischemic stroke from dwi using 3-d fully convolutional densenets. *IEEE transactions on medical imaging* **37**(9), 2149–2160 (2018)
24. Zheng, S., et al.: Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6881–6890 (2021)
25. Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., Torralba, A.: Learning deep features for discriminative localization. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2921–2929 (2016)
26. Zhou, D., Kang, B., Jin, X., Yang, L., Lian, X., Jiang, Z., Hou, Q., Feng, J.: Deepvit: Towards deeper vision transformer. arXiv preprint arXiv:2103.11886 (2021)
27. Zhou, Z., Rahman Siddiquee, M.M., Tajbakhsh, N., Liang, J.: 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)