# An atrium segmentation network with location guidance and siamese adjustment

Yuhan Xie[1], Zhiyong Zhang[1], Shaolong Chen[1], and Changzhen Qiu[1,*]

School of Electronics and Communication Engineering, Sun Yat-sen University, Guangzhou, China, qiuchzh@mail.sysu.edu.cn

**Abstract.** The segmentation of atrial scan images is of great significance for the three-dimensional reconstruction of the atrium and the surgical positioning. Most of the existing segmentation networks adopt a 2D structure and only take original images as input, ignoring the context information of 3D images and the role of prior information. In this paper, we propose an atrium segmentation network LGSANet with location guidance and siamese adjustment, which takes adjacent three slices of images as input and adopts an end-to-end approach to achieve coarse-to-fine atrial segmentation. The location guidance(LG) block uses the prior information of the localization map to guide the encoding features of the fine segmentation stage, and the siamese adjustment(SA) block uses the context information to adjust the segmentation edges. On the atrium datasets of ACDC and ASC, sufficient experiments prove that our method can adapt to many classic 2D segmentation networks, so that it can obtain significant performance improvements.

**Keywords:** Medical image segmentation,Location guidance,Siamese adjustment,UNet,SwinUNet

## 1 Introduction

Medical image segmentation of atrial region is of great significance for 3D reconstruction, pathological analysis, and surgical positioning based on atrium segmentation results. Before deep learning was widely used, many methods [14,11,6,4,10] based on traditional image processing were derived for segmentation. However, due to the noise in medical image imaging and the shape variability of organs in different cases, it is difficult for traditional methods to segment robustly and produce satisfactory results.

According to the understanding of the difficulty in atrium segmentation shown in **Figure 1**, we can find that the atrium images of MRI imaging have blurred boundaries (myocardium in the ACDC dataset [3]) and large fluctuations in the boundary shape (left atrium in ASC dataset [23]). Therefore, how to use context information to assist in localization is a key point. After deep learning has been widely used, many networks with superior performance have emerged for medical image segmentation, such as: UNet [19], UNet++ [26], TransUNet [7], etc. From the perspective of contextual information utilization, most networks

segment based on 2D slices, ignoring the contextual information in 3D images; while for 3D networks [17,25], larger computing resources and larger datasets are often required, which is very difficult to achieve, so that it is difficult to use in some restricted scenarios.
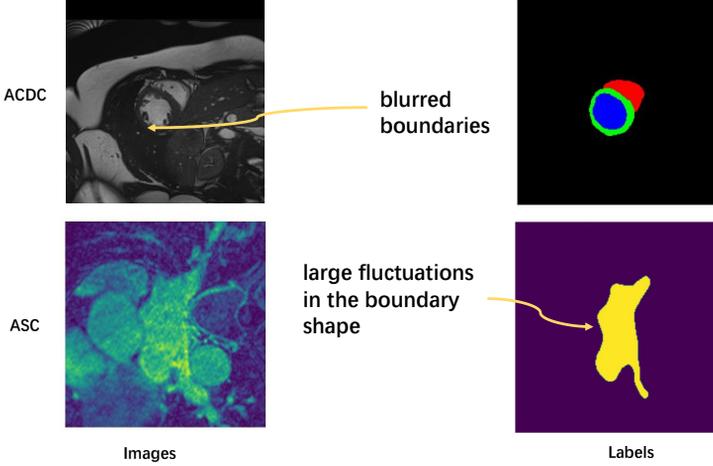


**Fig. 1.** Diffiulty in segmenting atrial images.

From the perspective of physicians manually segmenting medical images, we often use vision for localization first, and then focus on the localization area for detailed segmentation. Some methods use a multi-stage method to obtain human like location and adjustment through tailoring between different stages. However, this method is not end-to-end, which brings difficulty in training and deployment; some methods simply use multiple networks to connect in series to achieve an end-to-end from coarse to fine segmentation, however, it does not take full advantage of coarse localization information.

Inspired by the above problems, we design an atrium segmentation network based on localization guidance and siamese adjustment: location guidance and siamese adjustment Network(LGSANet).

**Contributions:**

1. A two-stage end-to-end method is proposed, using location guidance(LG) block to utilize the coarse localization information and use it in the fine-tuning stage;

2. We adopt a siamese three-layer network structure for the segmentation of three-layer continuous slices, and use siamese adjustment(SA) block between the decoder layers to utilize context information, and fine-tune the segmentation edges through the continuity between slices;

3. The location guidance and siamese adjustment design can be fully applied to most existing excellent 2D networks to improve their performance,such as UNet and SwinUNet. Sufficient experiments have demonstrated the robustness and universality of our method.

## 2    Related Work

### 2.1    Classic segmentation network

The medical image segmentation methods can be mainly divided into the methods based on convolutional neural network [19,26,16,18,8,12] and the methods based on transformer [7,5,25,13,22]. Among the methods based on convolutional neural network, UNet [19] in 2015 established the design direction of medical image segmentation network with the structure of classic encoder and decoder, and proved the effectiveness of skip connection; then UNet++ [26], UNet3+ [16] explored different design of skip connection respectively to achieve a better interaction between encoder and decoder. While AttUNet [18] introduces attention mechanism to the fusion of encoding and decoding features, ResUNet [8] introduces residual design in convolution module, optimizING the design of UNet architecture from different directions. With the in-depth study of transformers [21], the first medical segmentation network TransUNet [7] that introduced transformers appeared, using the transformer architecture to realize the interaction of global information in deep semantic features. Then SwinUNet [5], first medical segmentation network using pure transformers, proved the powerful representation capabilities of transformer.

### 2.2    Coarse-to-fine segmentation

The coarse-to-fine methods can be divided into mainly multi-stage methods [1] and end-to-end series connection methods [9,15]. The former cuts the results of the first stage, and then performs the optimization in the second stage. This method is cumbersome ,complex and cannot provide an end-to-end solution; while the series connection method, like SMCSRNet [9], does not fully consider the positioning information in the coarse positioning stage, the simple connection may not lead to a good result.

### 2.3    Network using context information

MEPDNet [20] uses a multi-encdoer and fusion decoder structure to utilize the context information, but it is directly fused in the decoder part which may lead to loss of context information; in addition,LSTM [2] is introduced to model the sequence relationship between the outputs of different slices, but too long-distance context information may increase the complexity of the model and the difficulty of training and deployment; ConResNet [24] adopts a multi-task method, in which task 1 predicts the segmentation result, task 2 predicts the residual between slices.However,this method also requires great computing resources and large memory.

## 3    Methods

We characterize the 3D medical image as $M \in R^{C \times H \times W}$, take adjacent consecutive three-slice image as input, and describe it as $X = [S_1, S_2, S_3] \in R^{3 \times H \times W}$. Each slice is sent to the LGSA network in parallel, and the output $Y = [M_1, M_2, M_3] \in R^{3 \times H \times W}$ is obtained. The overall expression is shown in **Formula 1**:

$$Y = LGSA(X; \theta) \tag{1}$$

Among them, we will select the output $M_2$ of the center slice as the final output.
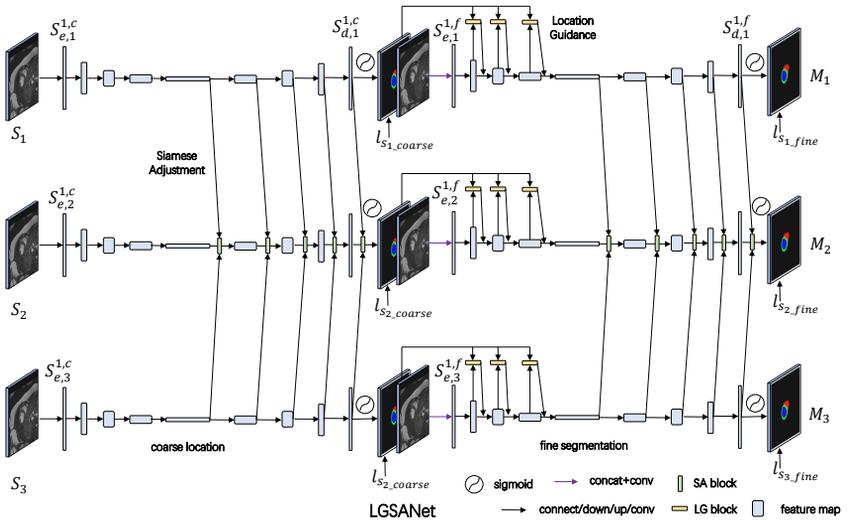
### 3.1    Overall structrue



**Fig. 2.** Overall Structure of LGSANet.

As shown in the **Figure 2**, we input consecutive three-layer 2D slices into three siamese single layer networks, each one is responsible for the segmentation of one-layer slice. The three single layer networks share parameters with each other to constrain the consistency of encoding and decoding. From the perspective of a single layer network, each one is divided into a coarse location and a fine segmentation stage. Between the two stages, location guidance(LG) blocks are used for fusion of multi-scale coarse location information and encoder features of fine segmentation stage, so that the model focuses on the located area in the fine segmentation stage. From the perspective of different slices, context information can be exchanged between different slices in the decoding stage. After

each layer of decoding, an additional cross-slice siamese adjustment(SA) block will be performed so that the decoding features of the middle-layer can fully obtain context information, so as to use the continuity of the upper and lower slices for edge adjustment. The overall process of LGSANet can be expressed as **Formula 2**:

$$S_{e,i}^{h,c} = CL_{\text{encoder}}^{h}\left(S_{e,i}^{h-1,c}\right), i = 1, 2, 3; h = 1, 2 \ldots N$$

$$S_{d,i}^{h-1,c} = CL_{\text{decoder}}^{h}\left(S_{a,i}^{h,c}\right), i = 1, 2, 3; h = 1, 2 \ldots N$$

$$S_{a,i}^{h,c} = CL_{SA}^{h}\left(S_{d,1}^{h,c}, S_{d,2}^{h,c}, S_{d,3}^{h,c}\right), i = 2; h = 1, 2 \ldots N$$

$$S_{e,i}^{h,f} = FS_{\text{encoder}}^{h}\left(S_{l,i}^{h-1,f}\right), i = 1, 2, 3; h = 1, 2 \ldots N \tag{2}$$

$$S_{l,i}^{h,f} = FS_{LG}^{h}\left(S_{e,i}^{h,f}, L_{i}^{h}\right), i = 1, 2, 3; h = 1, 2 \ldots N$$

$$S_{d,i}^{h-1,f} = FS_{\text{decoder}}^{h}\left(S_{a,i}^{h,f}\right), i = 1, 2, 3; h = 1, 2 \ldots N$$

$$S_{a,i}^{h,f} = FG_{SA}^{h}\left(S_{d,i-1}^{h,f}, S_{d,i}^{h,f}, S_{d,i+1}^{h,f}\right), i = 2; h = 1, 2 \ldots N$$

where $CL_{\text{encoder}}^{h}, CL_{\text{decoder}}^{h}, CL_{SA}^{h}, FG_{\text{encoder}}^{h}, FG_{LG}^{h}, FG_{\text{decoder}}^{h}, FG_{SA}^{h}$ represents different model components in LGSANet; $CL$ and $FS$ represent the coarse location and fine segmentation stages respectively; the superscript $h$ represents the h-th layer in the encoders or decoders, with a total of N layers, that is, N-1 times of downsampling are performed while N is 5 in UNet and 4 in SwinUNet. The subscripts *encoder* and *decoder* represent the encoder and decoder in this stage, while $SA$ and $LG$ represent the SA block and LG block in this stage. $S_{e,i}^{h,c}, S_{d,i}^{h,c}, S_{e,i}^{h,f}, S_{d,i}^{h,f}, S_{a,i}^{h,c}, S_{a,i}^{h,f}$ represent the feature maps of different stages appearing in LGSANet; The superscript $c$ or $f$ indicates that the feature belongs to the coarse location or fine segmentation stage; the subscripts $e$, $f$, $a$, and $l$ indicate that they are the feature maps after the encoder, decoder, SA block, and LG block respectively; the subscript $i$ indicates that it belongs to the feature map generated by the i-th slice. $L_{i}^{h}$ represents the multi-scale localization map generated in the coarse localization stage, the superscript $h$ represents the h-th layer, and the subscript $i$ represents the feature map generated by the i-th slice.

### 3.2   Siamese feature encoding and decoding

In the experiment, we mainly use UNet and SwinUNet as backbone. In the encoding process, UNet uses 5 layers of convolution modules as the basic units and performs 16 times downsampling, while SwinUNet uses 4 layers of swin transformer modules as the basic units, and completes 32 times downsampling through patch embedding and patch merging. In the decoding process , UNet uses a 4-layer convolution modules as the basic units and performs 16 times upsampling, while SwinUNet uses a 3-layer swin transformer modules as the basic units, and completes 32 times upsampling through patch expansion. For three-layer inputs, the encoder and decoder share parameters with each other. When

performing siamese adjustment, SwinUNet needs to reshape features from vector form to feature map form. The main reasons why we use UNet and SwinUNet as backbones are: these two networks represent the most basic way of applying CNN and transformer in medical image segmentation, which composed of two different basic components; using them as backbones can effectively verify the performance robustness of our method in both cases.

### 3.3    Location guidance block

The location guidance block mainly uses the location information in stage 1 to guide the coding of the encoder features in stage 2 after multi-scale scaling. Through the introduction of prior information in localization map, the encoder features can be strengthened so that the model can pay more attention to the localization area. Its schematic diagram and formula are shown in **Formula 3,4** and **Figure 3**:
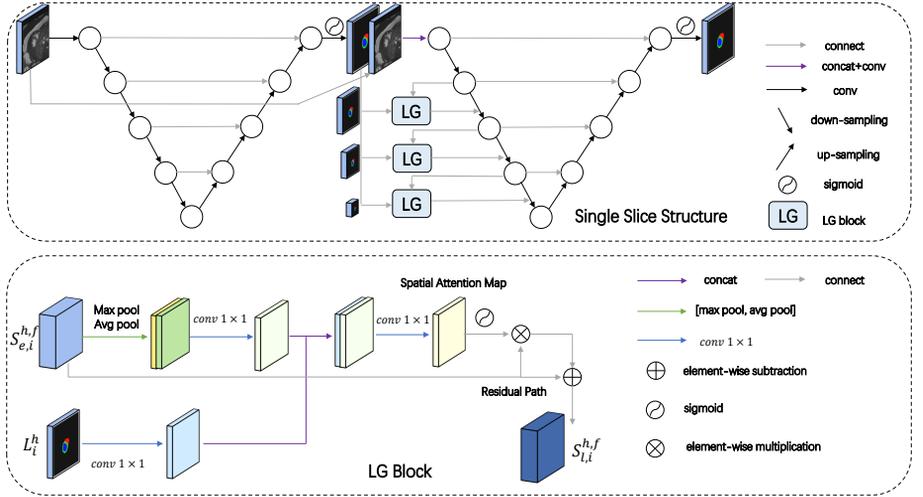


**Fig. 3.** Location guidance block.

$$SA_{map} = conv_{1\times1}(concat(mp\left(S_{e,i}^{h,f}\right), ap\left(S_{e,i}^{h,f}\right), conv_{1\times1}(L_i^h))) \tag{3}$$

$$S_{l,i}^{h,f} = S_{e,i}^{h,f} \times softmax\left(SA_{map}\right) + S_{e,i}^{h,f} \tag{4}$$

where, $mp$ and $ap$ represent maxpooling and avgpooling respectively, and $SA_{map}$ represents the spatial attention map.

### 3.4 Siamese adjustment block

The siamese adjustment block mainly uses the context information of adjacent layers to adjust the output results of the intermediate layers. The input three-layer features are adjacently subtracted to obtain the edge differences, and adjacently multiplied to enlarge the overlapping areas. The edge differences are fused to obtain the edge feature and the overlapping areas are fused to obtain the central feature. Finally, the edge feature and central feature are fused together as output. The center branch uses the coincidence of the context to strengthen the center positioning of the middle layer, and the edge branch uses the edge continuity constraint of the context to fine-tune the edge of the middle layer. Its schematic diagram and formula are shown in **Formula 5,6,7** and **Figure 4**:
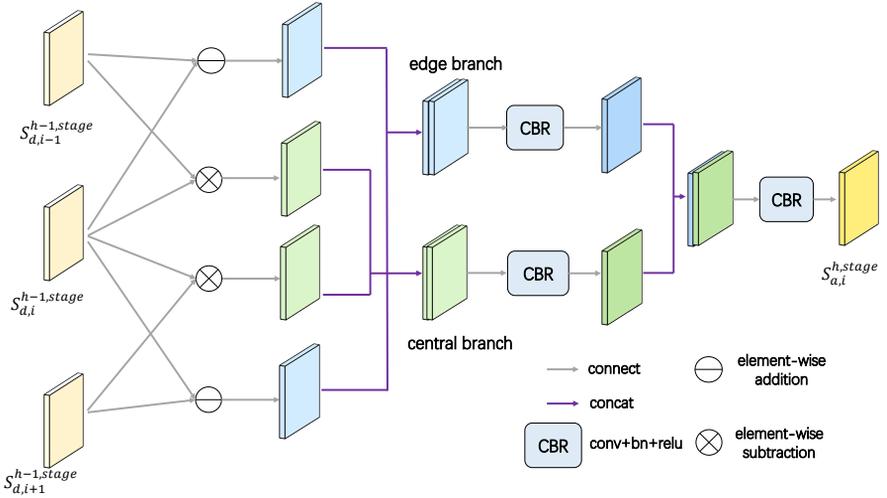


**Fig. 4.** Siamese adjustment block.

$$S_{edge} = CBR(S_{d,1}^{h,stage} - S_{d,2}^{h,stage}, S_{d,2}^{h,stage} - S_{d,3}^{h,stage}) \tag{5}$$

$$S_{central} = CBR(S_{d,1}^{h,stage} \times S_{d,2}^{h,stage}, S_{d,2}^{h,stage} \times S_{d,3}^{h,stage}) \tag{6}$$

$$S_{a,i}^{h,f} = CBR(S_{central}, S_{edge}) \tag{7}$$

where, stage represents the coarse location or fine segmentation stage, and CBR represents the combination of conv3×3, batch normalization, and Relu activation function.

### 3.5    Serial supervision and siamese supervision

In the design of the loss function, we use a combination of serial supervision and siamese supervision. While jointly supervising the coarse location and fine segmentation stages, the outputs of the three slices are all under siamese supervision to achieve the best optimization result. We use the output of the intermediate slices as the final output (except for the first and last layers) because it enjoys the most contextual information. The formula of the overall loss function can be expressed as formula 5,6,7:

$$L_{all} = \beta L_{coarse} + (1 - \beta)L_{fine} \tag{8}$$

$$L_{coarse} = \alpha l_{s_{1_{coarse}}} + (1 - 2\alpha)l_{s_{2_{coarse}}} + \alpha l_{s_{3_{coarse}}} \tag{9}$$

$$L_{fine} = \alpha l_{s_{1_{fine}}} + (1 - 2\alpha)l_{s_{2_{fine}}} + \alpha l_{s_{3_{fine}}} \tag{10}$$

$$l_{s_{i_{stage}}} = \sum_{k \in S_i} \left( -\frac{1}{2}y_k log\hat{y}_k + 1 - \frac{2 \times y_k \times \hat{y}_k}{y_k + \hat{y}_k} \right), stage = coarse, fine \tag{11}$$

Where $k$ represents any point in slice i, and $y_k$, $\hat{y}_k$ represent the groundtruth and prediction result respectively. $\alpha$=0.33, $\beta$=0.5. This is because the quality of location and fine segmentation results, as well as the output results of different layers, affect each other. In order to ensure the final output of the midlle layer as good as possible, the location information needs to be accurate enough, and the adjacent layer outputs that provide fused interaction information also need to be reliable enough, so their weights are equally distributed. This will also be verified in the subsequent ablation experiments. The supervision of a single output is composed of dice loss and bce loss. The weight distribution of dice loss is higher than that of bce loss, because dice loss is more suitable for the segmentation of small targets, which can better overcome the imbalance of foreground and background.

### 3.6    Structure comparison with different baselines

**Figure 5** shows the comparison with baselines.In the design of network architecture, we mainly focus on the utilization of contextual information and coarse localization information. Therefore, from a coarse-to-fine point of view, the baseline we mainly refer to is SMCSRNet, which uses UNet with simple concatenation to complete end-to-end multi-stage segmentation. The difference is that we introduce a location guidance block in the middle of the two stages in order to make the fine-stage encoder pay more attention to the localization area. From the perspective of context information utilization, the baselines we mainly refer to are 3-slice UNet and MEPDNet. The 3-slice UNet uses continuous three-layer slices stacking as input,which is sent to a common encoder and

decoder. MEPDNet uses three independent encoders extracting the features of the three-layer slices, and then performs fusion decoding. The difference is that we use siamese encoder and decoder to ensure the consistency and independence of encoding and decoding, and perform siamese adjustment between decoders to achieve information exchange at the same time. Subsequent experiments demonstrate that our design idea has gains in both aspects.
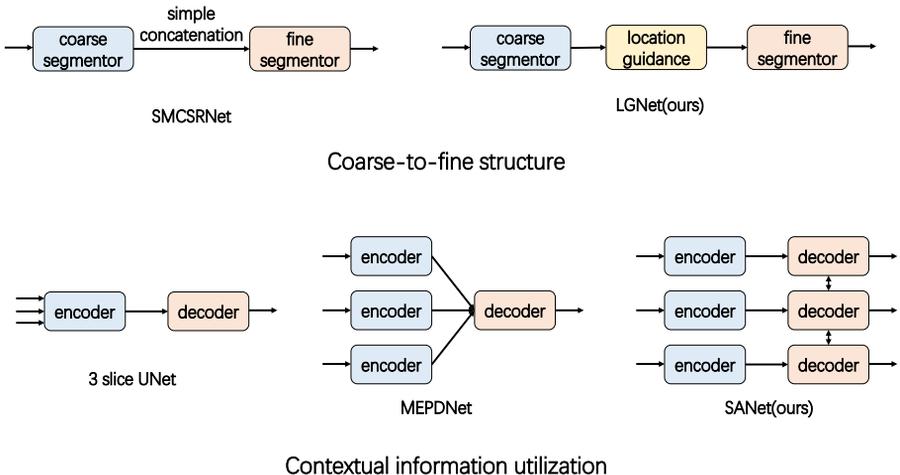


**Fig. 5.** Structure comparison with baselines.

## 4 Experimental results

### 4.1 Datasets

Our framework mainly uses the ACDC dataset [3](2017 Automated cardiac diagnosis challenge) and ASC dataset [23] (the 2018 atrial segmentation challenge) for experiments.

**ACDC**:The ACDC dataset contains 100 three-dimensional cardic MRI images to be segmented, each of which includes three types of manual annotations : left ventricle (LV), right ventricle (RV) and myocardium (MYO). Each case consists of a series of short-axis slices and the slice thickness is of 5 to 8 mm. The short-axis in-plane spatial resolution goes from 0.83 to 1.75 $mm^2$/pixel.There are totally 951 slices included into experiments.

**ASC**:The ASC dataset contains 152 three dimensional MRI images for left atrium (LA) and each of which includes one type of annotation: left atrium(LA).The image resolution is 0.625 × 0.625 × 1.25 $mm^3$. There are totally 13552 slices included into experiments.

## 4.2    Evaluation metrics

We use the 95% Hausdorff Distance (HD95)(mm) , Dice score (DSC)(%), F1
score(%) to characterize the performance of the segmentation. The formula of
DSC and F1 are shown in **Formula 12,13,14,15**:

$$F1 = 2 \times \frac{Precision \times Recall}{Precision + Recall} \tag{12}$$

$$Recall = \frac{TP}{TP + FN} \tag{13}$$

$$Precision = \frac{TP}{TP + FP} \tag{14}$$

$$DSC = 2 \times \frac{|X \cap Y|}{|X| + |Y|} \tag{15}$$

where $X$ denotes the segmentation result of the method, and $Y$ denotes the
ground truth.$TP$ denotes true negative result, $FP$ denotes false positive result
and $FN$ denotes false negative result.

## 4.3    Implementation details

The training, validation and testing processes are all conducted on two RTX3090
cards.The approach is implemented by Python3.7 with Pytorch. the batch size is
set to 24. An Adam optimizer is used in training process, with a learning rate of
5e-4,momentum of o.9 and weight decay of 1e-4. In the experiment,we train on
ACDC for 100 epochs and 50 epochs on ASC, while the early stopping is set to be
20 epochs on ACDC and 10 epochs on ASC. Before conducting the experiments,
we uniformly scale each slice to a size of 224×224. The training set, validation
set, and testing set are divided according to the ratio of 7:1:2. SwinUNet and
its variants are initialized with pre-trained weights, and the rest of the models
are initialized with Gaussian randomization. In order to ensure the reliability of
the experimental results, we repeated the experiments for each category 5 times
and obtained the average of the experimental results. We perform maximum
and minimum normalization on the both datasets in preprocessing, which can
be described by **Formula 16**:

$$I(x, y) = \frac{I(x, y) - I_{min}}{I_{max} - I_{min}} \tag{16}$$

where $I(x, y)$ denotes the grayscale of point (x,y), $Imax$ and $Imin$ denote the
maximum value and minimum value of an image.

In order to compare the difference between the output of the central layer of
LGSANet and other methods, and to maintain the consistency of the comparison
range, we let the head and tail slices of each 3D data not be included in the
testing range; It is worth mentioning that our LGSANet can actually output
the segmentation results of the first and last layers(this will be shown in our
experimental results).

### 4.4   Comparison with the state-of-the-art method

We select CNN-based segmentation networks: UNet, UNet++, DenseUNet, ResUNet and transformer-based segmentation networks: SwinUNet, TransUNet for basic comparison. Besides, we also choose MEPDNet, SMCSRNet and 3-slice UNet as baselines from the aspects of contextual information ultilization and coarse-to-fine segmentation. UNet and SwinUNet are adopted as the two different backbones of our proposed LGSANet respectively. The experimental results are shown in **Table 1,2,3**:

**Table 1.** Experiment results on ACDC dataset.

| Methods | | | RV | | | Myo | | | LV | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | DSC | HD95 | F1 | DSC | HD95 | F1 | DSC | HD95 | F1 |
| One slice input | One stage | ResUNet[8] | 88.45 | 1.39 | 87.18 | 83.41 | 1.37 | 83.63 | 91.37 | 1.00 | 91.66 |
| | | SwinUNet[5] | 90.64 | 1.29 | 90.80 | 83.82 | 1.10 | 84.09 | 94.58 | 0.67 | 94.68 |
| | | UNet[19] | 91.23 | 1.24 | 91.40 | 84.56 | 1.09 | 84.77 | 94.21 | 0.46 | 94.35 |
| | | UNet++[26] | 91.58 | 1.16 | 91.69 | 84.22 | 1.14 | 84.56 | 94.48 | 0.44 | 94.59 |
| | | DenseUNet[12] | 92.53 | 1.06 | 92.94 | 84.89 | 1.08 | 84.96 | 94.50 | 0.39 | 94.84 |
| | | TransUNet[7] | 92.28 | 1.10 | 92.59 | 84.50 | 1.11 | 84.92 | 95.39 | 0.17 | 95.68 |
| | Two stage | SMCSRNet[9] | 91.99 | 1.14 | 92.12 | 84.75 | 1.08 | 84.99 | 95.14 | 0.20 | 95.46 |
| | | LG-SwinUNet | 92.01 | 1.12 | 92.23 | 84.89 | 1.08 | 85.12 | 95.12 | 0.21 | 95.34 |
| | | LG-UNet | **92.48** | **1.09** | **92.53** | **85.47** | **1.05** | **85.59** | **95.54** | **0.32** | **95.60** |
| Three slice input | One stage | 3-slice UNet | 90.35 | 1.40 | 90.46 | 80.89 | 1.66 | 81.21 | 93.77 | 0.82 | 93.95 |
| | | MEPDNet[20] | 91.53 | 1.22 | 91.52 | 82.95 | 1.23 | 83.26 | 94.31 | 0.50 | 93.29 |
| | | SA-SwinUNet | 92.52 | 1.05 | 92.84 | 84.78 | 1.12 | 84.95 | 95.27 | 0.19 | 95.46 |
| | | SA-UNet | **92.94** | **1.00** | **93.11** | **86.65** | **1.00** | **86.84** | **95.53** | **0.17** | **95.70** |
| | Two stage | LGSA-SwinUNet | 92.92 | 1.02 | 93.05 | 85.51 | 1.05 | 85.08 | 95.73 | 0.11 | 95.90 |
| | | LGSA-UNet | **93.21** | **0.83** | **93.36** | **86.88** | **1.00** | **87.01** | **96.58** | **0.08** | **96.62** |

As shown in **Table 1** and **Table 2**, LGSAUNet achieved the best segmentation results in both ACDC and ASC datasets, and performed the best in DSC, HD95, and F1 indicators. It can be seen that the segmentation effect has been significantly improved by using LGSANet structure. On the three targets of RV, Myo, and LV in ACDC, LGSAUNet obtained 1.98%, 2.32%, and 2.37% dice performance improvement compared with UNet respectively. On the ASC dataset, LGSAUNet obtained 1.32% performance improvement compared to UNet. As a variant of LGSANet, LGSA-Swinunet has also improved significantly, with 2.28%, 1.69%, 1.15% on ACDC, and 2.64% on ASC. It can be seen that the design structure of LGSANet has good applicability to the architecture of both CNN and transformer.

For the models that simply using one slice as input and adopt two stage optimization, LGNet that using location guidance achieve better improvements compared with the SMCSRNet that using simply UNet concatenation.For the

**Table 2.** Experiment results on ASC dataset.

| Methods | | | DSC | HD95 | F1 |
|---|---|---|---|---|---|
| One slice input | One stage | Resunet[8] | 82.00 | 2.52 | 83.01 |
| | | Swinunet[5] | 87.53 | 1.70 | 88.17 |
| | | Unet[19] | 90.53 | 1.26 | 90.79 |
| | | Unet++[26] | 90.49 | 1.28 | 90.74 |
| | | DenseUNet[12] | 89.92 | 1.32 | 90.42 |
| | | TransUNet[7] | 90.00 | 1.31 | 90.37 |
| | Two stage | SMCSRNet[9] | 90.60 | 1.22 | 90.75 |
| | | LG-SwinUNet | 89.12 | 1.52 | 89.45 |
| | | LG-UNet | **90.81** | **1.21** | **90.97** |
| Three slice input | One stage | 3-slice UNet | 90.12 | 1.32 | 90.34 |
| | | MEPDNet[20] | 90.18 | 1.21 | 90.66 |
| | | SA-SwinUNet | 89.82 | 1.36 | 89.98 |
| | | SA-UNet | **91.57** | **1.09** | **91.74** |
| | Two stage | LGSA-Swinunet | 90.17 | 1.30 | 90.51 |
| | | LGSA-unet | **91.85** | **1.02** | **92.06** |

models that using three slices as input but simply using one stage optimza-
tion,SANet that using siamese adjustment also performs better than 3-slice UNet
and MEPDNet.

**Table 3.** The dice(%) of different output in LGSANet.

| Datasets | | Coarse location | | | Fine segmentation | | |
|---|---|---|---|---|---|---|---|
| | | 1 | 2 | 3 | 1 | 2 | 3 |
| ACDC | RV | 91.34 | 92.89 | 90.72 | 91.71 | **93.21** | 91.34 |
| | Myo | 86.23 | 86.60 | 84.55 | 86.62 | **86.88** | 84.89 |
| | LV | 95.75 | 96.19 | 95.03 | 96.08 | **96.58** | 95.48 |
| ASC | | 90.89 | 91.47 | 90.77 | 91.31 | **91.85** | 91.26 |

From the results in **Table 3**, it can be seen that the output results of the
middle layer slices are better than the adjacent layers.The output results of the
fine segmentation stage are better than the results of the coarse location stage
as well. It illustrates that the idea of optimizing the central layer from coarse to
fine with the help of context information actually works. Besides,for the output
of the adjacent layers in fine segmentation, although it is a little worse than the
center layer, is still better than the output of its 2D basic network, so it can
also be benefit for the segmentation of the first and last slices.The mean and
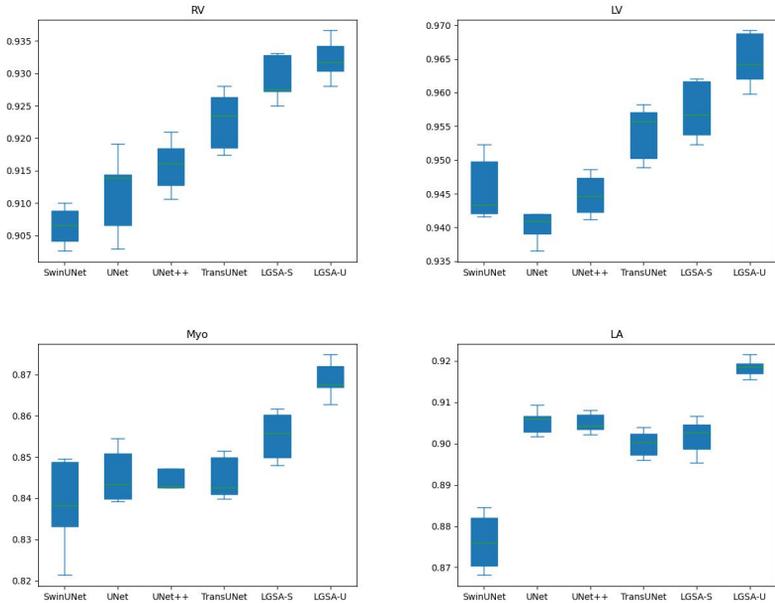fluctuation range of the experimental results are shown in **Figure 6**.

**Fig. 6.** The box and whisker plot on ACDC(LV,Myo,RV) and ASC(LA) dataset.

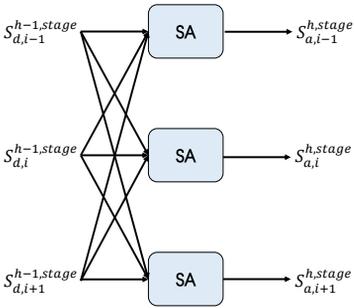**Table 4.** Ablation on different modules used in LGSANet.

| Methods | ACDC | | | ASC | | |
|---|---|---|---|---|---|---|
| | DSC | HD95 | F1 | DSC | HD95 | F1 |
| UNet | 90.00 | 0.93 | 90.17 | 90.53 | 1.26 | 90.79 |
| UNet+LG | 91.16 | 0.82 | 91.24 | 90.81 | 1.21 | 90.97 |
| UNet+SA | 91.70 | 0.72 | 91.88 | 91.57 | 1.09 | 91.74 |
| UNet+LG+SA | **92.22** | **0.64** | **92.33** | **91.85** | **1.02** | **92.06** |

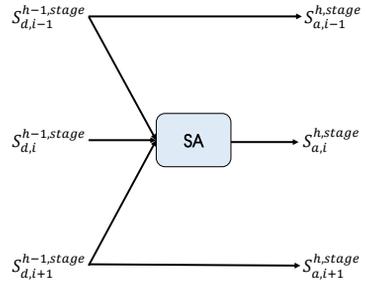### 4.5   Ablation Analysis of Our Method

From the results in **Table 4**, it can be seen that the use of LG block and SA block can gradually improve the performance of the segmentation network. As shown in **Figure 6**, it can be seen that the single-head output results can obtain better benefits than the multi-head output.In multi-head mode, it is probably unreasonable to use the edge information to correct the edge layer features.But the middle layer in the single-head mode can obtain balanced context information, so the effect is relatively better.As shown in **Table 5**, the increase in the number of SI blocks enables both the deep and shallow contextual information to be acquired and adjusted during the decoding process, which also shows that the SA block and skip connection in the 2d network play a similar role. As shown in **Table 7**, the combination of serial supervision and siamese supervision en-

**Table 5.** Ablation on different design of SA block.

| Methods | ACDC | | |
|---|---|---|---|
| | DSC | HD95 | F1 |
| Multi-head | 91.90 | 0.68 | 91.99 |
| Central-head | **92.22** | **0.64** | **92.33** |



Multi-head type SA block                Central type SA Block

**Fig. 7.** Multi-head type and central type SA block.

ables LGSANet to gradually obtain better performance under the structure of LGSANet.

**Table 6.** Ablation on the number of SA block. The number of SA block veries from 1 to 5 in LGSA-UNet.

| Number of | ACDC | | |
|:---:|:---:|:---:|:---:|
| SA block | DSC | HD95 | F1 |
| 1 | 91.27 | 0.78 | 91.54 |
| 3 | 91.92 | 0.68 | 92.05 |
| 5 | **92.22** | **0.64** | **92.33** |

**Table 7.** Ablation on loss function. OS repesents that only ouput of central slice in fine segmentation stage is supervised; SiS means serial supervision and Sis means siamese supervision.

| Supervision | ACDC | | |
|:---:|:---:|:---:|:---:|
| type | DSC | HD95 | F1 |
| OS | 90.57 | 0.88 | 90.84 |
| SeS | 90.89 | 0.79 | 91.21 |
| SiS | 91.84 | 0.70 | 91.97 |
| SeS+SiS | **92.22** | **0.64** | **92.33** |

## 5   Visualization

It can be seen from the visualization results that our approach LGSA-UNet reach the best proformance.In the two datasets, our method can accurately segment the target, making the segmentation result smoother and more accurate. In ACDC dataset, the discontinuity of the segmentation edge is greatly reduced. ASC is a dataset with rich target morphological changes, our method can also better fit the boundaries of the target. Compared with its 2D backbone model, LGSANet can achieve great performance improvement in segmentation task.

## 6   Conclusion and future work

In this paper, we propose an atrium segmentation network based on location guidance and siamese adjustment, which takes consecutive three-layer slices as inputs.It uses location information in stage 1 to guide encoding features in stage 2, and conducts siamese interactions among the three-layer slices to take advantage of contextual information. We use a combination of serial supervision and siamese supervision to obtain the best optimization effect of this network.
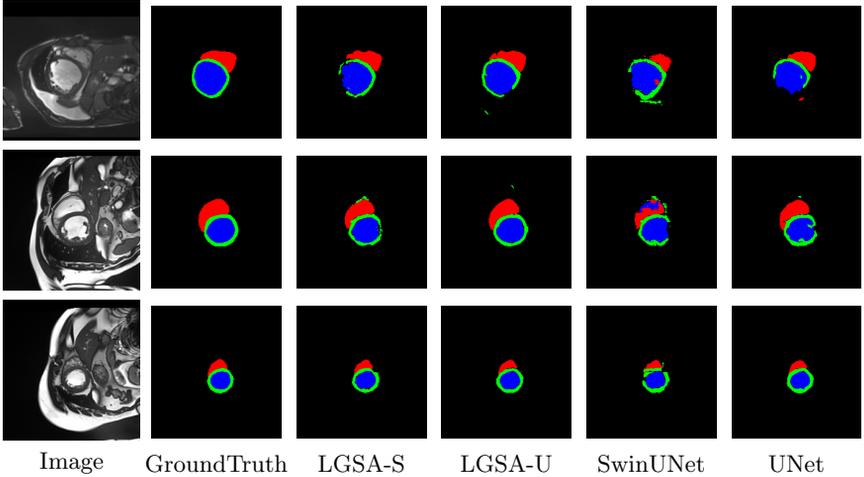
|  |  |  |  |  |  |
|---|---|---|---|---|---|
| Image | GroundTruth | LGSA-S | LGSA-U | SwinUNet | UNet |

**Fig. 8.** Visualization of segmentation results using different methods in ACDC dataset. LGSA-S repesents LGSA-SwinUNet and LGSA-U repesents LGSA-UNet.

Experiments show that our method is suitable for classic 2D networks such as UNet, SwinUNet to achieve a significant performance improvement. In future work, we will further attempt to introduce edge detectors into segmentation tasks to improve the performance of segmentation.

## A    Acknowledgements

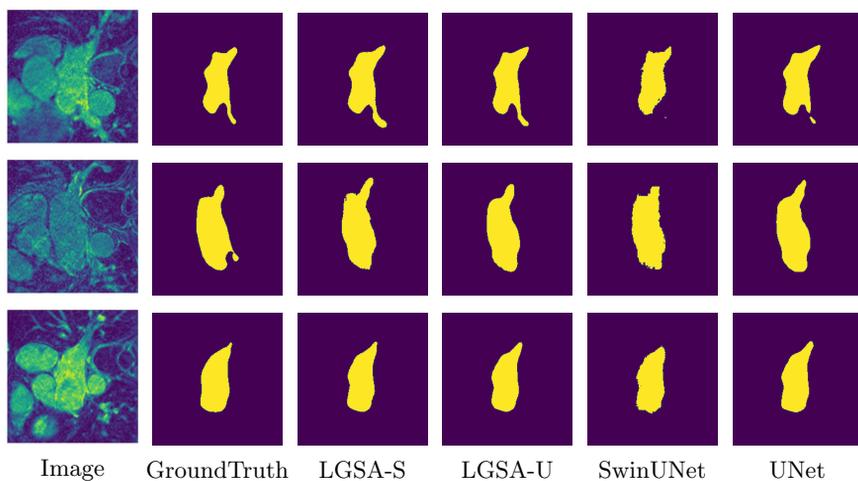| Image | GroundTruth | LGSA-S | LGSA-U | SwinUNet | UNet |

**Fig. 9.** Visualization of segmentation results using different methods in ASC dataset.LGSA-S repesents LGSA-SwinUNet and LGSA-U repesents LGSA-UNet.

# References

1. Astuto, B., Flament, I., K. Namiri, N., Shah, R., Bharadwaj, U., M. Link, T., D. Bucknor, M., Pedoia, V., Majumdar, S.: Automatic deep learning–assisted detection and grading of abnormalities in knee mri studies. Radiology: Artificial Intelligence **3**(3), e200165 (2021)
2. Bai, W., Suzuki, H., Qin, C., Tarroni, G., Oktay, O., Matthews, P.M., Rueckert, D.: Recurrent neural networks for aortic image sequence segmentation with sparse annotations. In: International conference on medical image computing and computer-assisted intervention. pp. 586–594. Springer (2018)
3. Bernard, O., Lalande, A., Zotti, C., Cervenansky, F., Yang, X., Heng, P.A., Cetin, I., Lekadir, K., Camara, O., Ballester, M.A.G., et al.: Deep learning techniques for automatic mri cardiac multi-structures segmentation and diagnosis: is the problem solved? IEEE transactions on medical imaging **37**(11), 2514–2525 (2018)
4. Candemir, S., Jaeger, S., Palaniappan, K., Musco, J.P., Singh, R.K., Xue, Z., Karargyris, A., Antani, S., Thoma, G., McDonald, C.J.: Lung segmentation in chest radiographs using anatomical atlases with nonrigid registration. IEEE transactions on medical imaging **33**(2), 577–590 (2013)
5. Cao, H., Wang, Y., Chen, J., Jiang, D., Zhang, X., Tian, Q., Wang, M.: Swin-unet: Unet-like pure transformer for medical image segmentation. arXiv preprint arXiv:2105.05537 (2021)
6. Castro-Mateos, I., Pozo, J.M., Pereañez, M., Lekadir, K., Lazary, A., Frangi, A.F.: Statistical interspace models (sims): application to robust 3d spine segmentation. IEEE transactions on medical imaging **34**(8), 1663–1675 (2015)
7. 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)
8. Diakogiannis, F.I., Waldner, F., Caccetta, P., Wu, C.: Resunet-a: A deep learning framework for semantic segmentation of remotely sensed data. ISPRS Journal of Photogrammetry and Remote Sensing **162**, 94–114 (2020)
9. Ding, Y., Chen, F., Zhao, Y., Wu, Z., Zhang, C., Wu, D.: A stacked multi-connection simple reducing net for brain tumor segmentation. IEEE Access **7**, 104011–104024 (2019)
10. Dodin, P., Martel-Pelletier, J., Pelletier, J.P., Abram, F.: A fully automated human knee 3d mri bone segmentation using the ray casting technique. Medical & biological engineering & computing **49**(12), 1413–1424 (2011)
11. Engstrom, C.M., Fripp, J., Jurcak, V., Walker, D.G., Salvado, O., Crozier, S.: Segmentation of the quadratus lumborum muscle using statistical shape modeling. Journal of Magnetic Resonance Imaging **33**(6), 1422–1429 (2011)
12. Guan, S., Khan, A.A., Sikdar, S., Chitnis, P.V.: Fully dense unet for 2-d sparse photoacoustic tomography artifact removal. IEEE journal of biomedical and health informatics **24**(2), 568–576 (2019)
13. Hatamizadeh, A., Nath, V., Tang, Y., Yang, D., Roth, H., Xu, D.: Swin unetr: Swin transformers for semantic segmentation of brain tumors in mri images. arXiv preprint arXiv:2201.01266 (2022)
14. Heimann, T., Meinzer, H.P.: Statistical shape models for 3d medical image segmentation: a review. Medical image analysis **13**(4), 543–563 (2009)
15. Hu, J., Wang, H., Gao, S., Bao, M., Liu, T., Wang, Y., Zhang, J.: S-unet: A bridge-style u-net framework with a saliency mechanism for retinal vessel segmentation. IEEE Access **7**, 174167–174177 (2019)

16. Huang, H., Lin, L., Tong, R., Hu, H., Zhang, Q., Iwamoto, Y., Han, X., Chen, Y.W., Wu, J.: Unet 3+: A full-scale connected unet for medical image segmentation. In: ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). pp. 1055–1059. IEEE (2020)
17. Milletari, F., Navab, N., Ahmadi, S.A.: V-net: Fully convolutional neural networks for volumetric medical image segmentation. In: 2016 fourth international conference on 3D vision (3DV). pp. 565–571. IEEE (2016)
18. Oktay, O., Schlemper, J., Folgoc, L.L., Lee, M., Heinrich, M., Misawa, K., Mori, K., McDonagh, S., Hammerla, N.Y., Kainz, B., et al.: Attention u-net: Learning where to look for the pancreas. arXiv preprint arXiv:1804.03999 (2018)
19. 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)
20. Shi, D., Liu, R., Tao, L., He, Z., Huo, L.: Multi-encoder parse-decoder network for sequential medical image segmentation. In: 2021 IEEE International Conference on Image Processing (ICIP). pp. 31–35. IEEE (2021)
21. 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)
22. Xie, Y., Zhang, J., Shen, C., Xia, Y.: Cotr: Efficiently bridging cnn and transformer for 3d medical image segmentation. In: International conference on medical image computing and computer-assisted intervention. pp. 171–180. Springer (2021)
23. Xiong, Z., Xia, Q., Hu, Z., Huang, N., Bian, C., Zheng, Y., Vesal, S., Ravikumar, N., Maier, A., Yang, X., et al.: A global benchmark of algorithms for segmenting the left atrium from late gadolinium-enhanced cardiac magnetic resonance imaging. Medical Image Analysis **67**, 101832 (2021)
24. Zhang, J., Xie, Y., Wang, Y., Xia, Y.: Inter-slice context residual learning for 3d medical image segmentation. IEEE Transactions on Medical Imaging **40**(2), 661–672 (2020)
25. Zhou, H.Y., Guo, J., Zhang, Y., Yu, L., Wang, L., Yu, Y.: nnformer: Interleaved transformer for volumetric segmentation. arXiv preprint arXiv:2109.03201 (2021)
26. Zhou, Z., Siddiquee, M.M.R., Tajbakhsh, N., Liang, J.: Unet++: Redesigning skip connections to exploit multiscale features in image segmentation. IEEE transactions on medical imaging **39**(6), 1856–1867 (2019)