# Selective Information Passing for MR/CT Image Segmentation

**Qikui Zhu · Liang Li# · Jiangnan Hao# · Yunfei Zha# · Yan Zhang# · Yanxiang Cheng# · Fei Liao# · Pingxiang Li#**

**Abstract** Automated medical image segmentation plays an important role in many clinical applications, which however is a very challenging task, due to complex background texture, lack of clear boundary and significant shape and texture variation between images. Many researchers proposed an encoder-decoder architecture with skip connections to combine low-level feature maps from the encoder path with high-level feature maps from the decoder path for automatically segmenting medical images. The skip connections have been shown to be effective in recovering fine-grained details of the target objects

Qikui Zhu
School of Computer Science, Wuhan University, Wuhan, China. E-mail: QikuiZhu@whu.edu.cn

Liang Li
Department of Radiology, Renmin Hospital of Wuhan University, Wuhan, China. E-mail: liliang_082@163.com

Jiangnan Hao
Xi'an Aeronautical University. E-mail: haojiangnan@mail.nwpu.edu.cn

Yunfei Zha
Renmin Hospital of Wuhan University, Wuhan, Hubei Province, China. E-mail: zhayunfei999@126.com

Yan Zhang
Department of Clinical Laboratory, Renmin Hospital of Wuhan University, Wuhan, China. E-mail: peneyyan@mail.ustc.edu.cn

Yanxiang Cheng
Department of Obstetrics and Gynecology, Renmin Hospital of Wuhan University, Wuhan, China. E-mail: yanxiangcheng@whu.edu.cn

Fei Liao
Department of Gastroenterology, Renmin Hospital of Wuhan University, Wuhan, China. E-mail: feiliao@whu.edu.cn

Pingxiang Li
Renmin Hospital of Wuhan University, Wuhan, China. E-mail: pxli@whu.edu.cn

'#' indicates co-corresponding authors.

arXiv:2010.04920v1 [eess.IV] 10 Oct 2020

and may facilitate the gradient back-propagation. However, not all the feature maps transmitted by those connections contribute positively to the network performance. In this paper, to adaptively select useful information to pass through those skip connections, we propose a novel 3D network with self-supervised function, named selective information passing network (SIP-Net). We evaluate our proposed model on the MICCAI Prostate MR Image Segmentation 2012 Grant Challenge dataset, TCIA Pancreas CT-82 and MICCAI 2017 Liver Tumor Segmentation (LiTS) Challenge dataset. The experimental results across these data sets show that our model achieved improved segmentation results and outperformed other state-of-the-art methods. The source code of this work is available at https://github.com/ahukui/SIPNet.

**Keywords** Medical image segmentation · convolutional neural network · attention-focused module

## 1 Introduction

Medical image segmentation is an essential part of medical image analysis. Accurate segmentation of medical image provides very useful information for computer aided diagnosis and treatment of cancers as well as other diseases[1]. For instance, segmentation of the liver and tumors plays an important role in hepatocellular carcinoma diagnosis [2]. Accurate prostate segmentation is useful for treatment planning and therapeutic procedures for prostate cancer[3, 4,5]. However, automated medical image segmentation is very challenging for several reasons. Taking prostate segmentation as an example: First, due to many slices only have small part of segmented tissues specifically at the apex and base, which always led to those slices lack of clear boundary and make the automated segmentation fail. Second, imaging artifacts always distribute in the whole image randomly, which negatively influence the process of segmentation. Third, tissues can have a wide variation in size and shape among different slices, which adds to the complexity of segmentation. Fourth, the complex background and fuzzy boundary also make the segmentation process challenging. Furthermore, different from natural images dataset, the size of available medical image dataset is limited. Fig.1 shows examples of prostate MR images. Fig.1(a) shows the phenomenon that imaging artifacts locate in prostate region. Fig.1(b) shows prostate region lacks clear boundary. Fig.1(c) shows the prostate and surrounding tissues have similar intensity distribution. All of above phenomena bring challenges for automated medical image segmentation.

To overcome the above challenges, over the past few decades, various methods have been developed for medical image segmentation, including machine learning based methods [6,7,8,9,10,11,12,13], level sets[14], atlas-based methods[15,16,17], super-pixel segmentation[18] and active shape model[19, 20]. Recently, deep convolutional neural networks (CNNs) have become the dominant machine learning approach due to their superior performance. CNNs have achieved state-of-the-art performances in many fields including computer
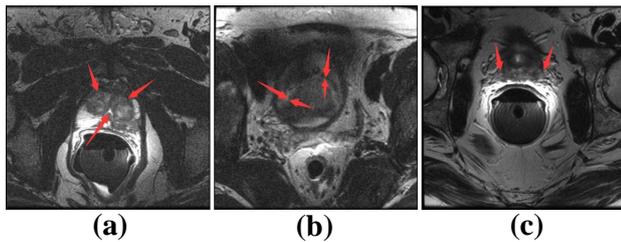
**Fig. 1** Challenges in segmenting the prostate from MR images. (a) Noise inside prostate. (b) Weak boundary. (c) Surrounding tissues having similar intensity distribution with prostate.

vision [21, 13, 22, 23, 24, 25], natural language processing (NLP)[26, 27, 28], medical image analysis[29]. The superiority of CNNs[30] can be partially attributed to the ability of learning hierarchical representation of the data.

However, medical image segmentation has a higher-level requirement of accuracy than natural image segmentation, where many excellent networks, such as VGG [31] and FCN [32], cannot be directly utilized. To obtain accurate segmentation results and overcome the problems specific to medical imaging, specific models have been proposed for medical image analysis. For instance, Milletari et al. [33] proposed a network architecture based on the volumetric CNNs, which can segment prostate volumes in a fast and accurate manner. Yu et al. [34] proposed a novel volumetric CNN with mixed long and short residual connections for automated prostate segmentation. Gibson et al. [35] proposed a network called DenseVNet, which can segment the pancreas, esophagus, stomach, liver, spleen, gallbladder, left kidney and duodenum accurately. Li et al. [36] proposed a novel hybrid densely connected U-Net for liver and tumor segmentation. One thing that these medical image segmentation networks have in common is an encode and decode architecture with skip connections for combining low-level feature maps from the encoder path with high-level feature maps from the decoder path. There is no doubt that the skip connections are effective in recovering fine-grained details of the target objects and help the gradient back-propagation. However, as a lot of information can be passed through those skip connections, do all the feature maps transmitted by those connections always contribute positively to the network performance?

To answer this question, we analyzed the behavior of the classical U-Net [37] with and without the long skip connections on the task of prostate segmentation. The segmentation results are shown in Fig. 2. Compared with ground truth segmentation, U-Net can obtain finer details and higher accuracy in general. However, the segmentation result of fully convolutional network (FCN) [32] is smoother and that of U-Net picks up non-prostate regions when those areas are highly inhomogeneous. To make the long skip connections inside the network select the useful information and further improve medical image segmentation performance, in this paper, we propose a novel 3D convolutional network, named SIP-Net. Our proposed SIP-Net adopts Densely-connected
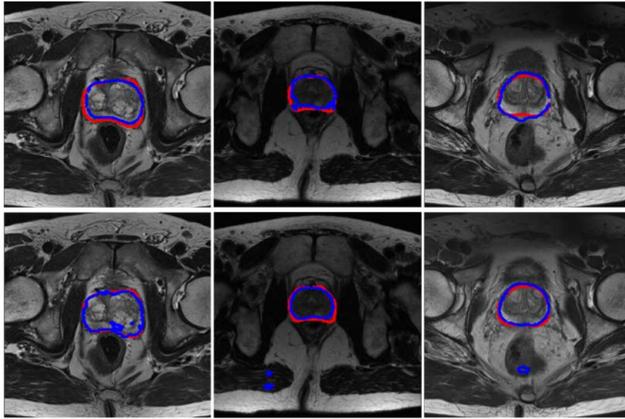
**Fig. 2** Top row: Segmentation results of U-Net without long skip connections, which is in fact equivalent to the original FCN; Bottom row: Segmentation results of U-Net. The red and blue contours indicate the ground truth and segmentation results, respectively.

Residual Blocks (DRBs) and Attention-focused Modules (AMs). The contributions of this work are summarized as follows.

– Inspired by the attention mechanism, we propose to integrate attention-focused modules into our model to make the long connections transmit mainly useful features and reduce the negative impact of noise from feature maps. That makes the long connections focus more on the regions to be segmented and the irrelevant noise features from the background and surrounding tissues may be suppressed during feature transmission.
– In the same time, to overcome the problem of small size of medical image data, we integrate three different types of connections seamlessly into our proposed model. Together with the above attention-focused modules, these connections improve training efficiency and feature extraction capability of the network by enhancing information propagation and encouraging feature reuse.
– To reduce the computational load and more importantly the number of network parameters for alleviating the potential overfitting problem, we design a modified dense block to construct deeper network, which possesses more than 90 convolutional layers but fewer parameters. Our experimental results show that the proposed model is effective in addressing the problems of complex background, fuzzy boundary and large shape variations.

The remainder of the paper is organized as follows. Section 2 provides a brief survey of related works. Section 3 describes the details of the proposed 3D segmentation network model. In Section 4, various experiments of segmenting prostate MR images, pancreas CT images, liver CT images are performed to validate the proposed model. The performance of the proposed method is further discussed through ablation studies in Section 5. Finally, several concluding remarks are drawn in Section 6.

## 2 Related Works

In this section, we give a brief review of deep learning techniques for semantic image segmentation. We first review the methods for natural image segmentation and then discuss the ones specialized for medical image segmentation.

### 2.1 Deep Learning for Semantic Segmentation

Semantic segmentation is a critical component in image understanding. The task of semantic segmentation is to assign a categorical label to every pixel in an image. Over the past few years, deep learning based methods and in particular convolutional networks (CNNs) have improved segmentation results remarkably in pixel-wise semantic segmentation tasks. This success can be attributed to the ability of hierarchical representation of CNNs. Fully convolutional networks (FCNs) mark a major milestone in CNN based semantic segmentation [32], which is trained end-to-end to perform pixels-to-pixels segmentation. Since then, FCNs have dominated the field of semantic image segmentation with a number of extensions. For instance, Li et al. [38] extended the FCN model for instance-aware semantic segmentation. The model significantly improves the segmentation performance in both accuracy and efficiency.

In the same time, researchers develop deeper and more powerful CNN models to extract more discriminating and complex representation features. For example, Simonyan et al. [39] proposed a 19-layer network, the famous VGG-19 model, to investigate the effect of the depth of CNNs on their accuracy in large-scale image recognition. He et al. [31] presented a residual learning framework to ease the training of very deep networks. Based on this framework, the author proposed a 101-layer model (ResNet-101) and a 152-layer model (ResNet-152) and won the first places in several tracks in ILSVRC & COCO 2015 competitions[1]. Soon after that, Wu et al. [40] proposed a method for high-performance semantic image segmentation based on the deep residual networks, which achieves the state-of-the-art performance.

### 2.2 Deep Learning for Medical Image Segmentation

Recently, deep CNNs have also become the dominant approach for medical image segmentation. Many researchers have employed various CNN models to segment images from different medical imaging modalities. In our previous work, we proposed a deeply supervised CNN model [41], which employs additional supervised layers and utilizes the residual information to segment the prostate from MR image. To exploit the information from different views of volumetric images but without using 3D convolutions, Mortazi et al. [42] proposed a multi-view CNN to segment structures from cardiac MR images by using an adaptive fusion strategy. Han [43] proposed a 2.5D model to segment
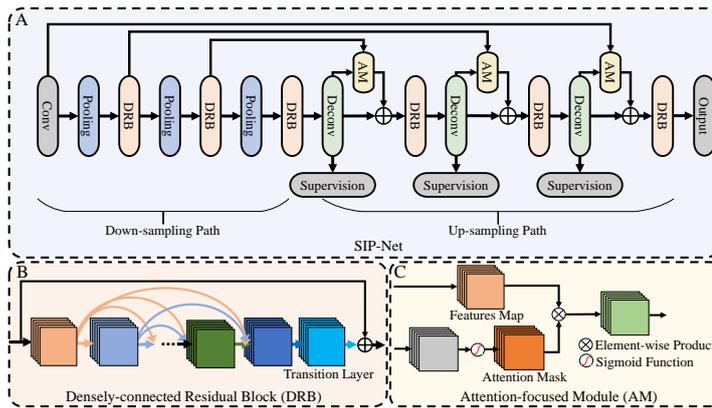
---

[1] http://image-net.org/challenges/ilsvrc+mscoco2015

**Fig. 3** The illustration of the pipeline for medical image segmentation. A. The proposed SIP-Net. B. The structure of Densely-connected Residual Block (DRB). C. The structure of attention-focused module (AM).

liver tumors, which takes a stack of adjacent slices as input and produces the segmentation map corresponding to the center slice.

To fully exploit 3D spatial information in volumetric MR images, a few studies employed 3D convolutions to build CNNs. For example, Li et al. [36] proposed a novel hybrid densely connected U-Net for liver and tumor segmentation. The proposed model consists of a 2D Dense-U-Net and a 3D counterpart, which can extract intra-slice features and hierarchically aggregate 3D contexts under the spirit of the auto-context algorithm [44]. Chen et al. [45] extended deep residual learning into a 3D for 3D brain segmentation. This model also seamlessly integrates the low-level image appearance features, implicit shape information and high-level context together for further improving the 3D segmentation performance. Recently, Yu et al. [46] proposed a novel densely-connected volumetric CNN, which adopts the 3D fully convolutional architecture to automatically segment cardiac and vascular structures from 3D cardiac MR images.

Compared with 2D networks, these 3D networks were able to achieve better segmentation performance. However, 3D CNNs have a much larger number of parameters and computational complexity than 2D networks. Due to the limited size of typical medical image dataset, it makes the network difficult to train. Furthermore, the trained network easily suffers from overfitting. Therefore, there is still much need in pushing the potential of CNNs by effectively extracting the information from limited training data to improve the segmentation performance and also reduce the complexity of the networks to avoid overfitting.

## 3 Methods

In this section, we first give an overview of the proposed SIP-Net and then discuss each module of the model in detail.

### 3.1 SIP-Net

In order to fully use the 3D spatial contextual information of volumetric data to accurately segment medical images, in this paper, we propose a 3D CNN with densely-connected residual blocks (DRBs) and attention-focused modules (AMs), named SIP-Net. The overall structure is shown in Fig. 3. The proposed SIP-Net contains two paths: down-sampling path and up-sampling path. The down-sampling path consists of one convolutional block, three DRBs and three average pooling layers. The pooling layers use stride of 2, which gradually reduces the resolution of feature map and increases the receptive field of the convolutional layers. To obtain accurate segmentation result in the original image resolution, an up-sampling path is implemented, which contains three deconvolutional layers and three DRBs. The deconvolutional layers gradually up-sample the feature maps until reaching the original input size. The overall illustration and detailed structure of proposed network are shown in Fig. 3(A) and Table 1, respectively.

In our proposed SIP-Net, we could have used the long connections between the down-sampling path and up-sampling path to connect the blocks in the same resolution level in the down-sampling and up-sampling paths. However, our study shows that simply adding the long connections may cause noisy segmentation by considering part of noise as shown in Fig. 2. To make the network focus more on the segmented region and reduce the negative influence from background and surrounding tissues, in this paper, we employ the attention mechanism in our proposed model. Inspired by the attention mechanism in residual attention network [47], three attention-focused modules are used in up-sampling path, which reduces irrelevant noise in background and surrounding tissues and holds segmenting features from down-sampling path and make the network focus more on the areas to be segmented in the up-sampling path.

In addition, to enforce the attention-focused modules to act effectively as information pass filters, we also integrate a deep supervision mechanism [48] for the attention-focused modules. An additional supervision layer is added after each deconvolutional layer. Each of the three additional supervision layers consists of one up-sampling layer for enlarging the feature map to its original size and one convolutional layer for obtaining the segmentation output as shown in Fig. 3(A). Those additional supervision layers bring two advantages. First, it helps to supervise the attention-focused modules to produce accurate attention masks to guide information passing. Second, it can accelerate the network convergence speed during training due to the shorter backpropagation paths from the additional supervision outputs.

In total, our proposed SIP-Net has more than 100 layers in depth including convolutional layers, pooling layers, layers in dense blocks, transitional layers, dropout layers and deconvolutional layers. The dense layers contain different numbers of BN-ReLU-Conv($1\times1\times1$)-BN-ReLU-Conv($3\times3\times3$) with growth rate of $k = 32$. The transition layer is implemented using a BN-ReLU-Conv($1\times1\times1$) layer. After each Conv($3\times3\times3$) layer, a dropout layer with 0.3 dropout rate is added to help deal with the potential overfitting problem. The designs of DRB and AM are shown in Fig. 3(B) and Fig. 3(C) and the details are given in the rest of this section.

### 3.2 Densely-connected Residual Block (DRB)

Let $x_l$ be the output of the $l^{th}$ convolutional layer, which can be considered as the result of applying a non-liner transformation $H_l$ defined as a convolution followed by a batch-normalization and a rectified linear unit (ReLU) in the $l^{th}$ layer. And $x_0$ denotes the input data sample passed to the CNN. For a classical CNN layer with straightforward connection, $x_l$ can be modelled as

$$x_l = H_l(x_{l-1}),\tag{1}$$

where $x_{l-1}$ is the output of the $(l-1)^{th}$ layer. However, when a network goes deeper, the network suffers from the degradation problem - the gradient may vanish or explode. This phenomenon leads to large training errors and the network training may not converge.

To alleviate the problem by promoting information propagation within the network, in this paper, we propose a new block by combining dense block[49] with residual connection as show in Fig. 3(B). The dense connected layers provide a directly connects with all subsequent layers. The feature maps produced by all the preceding layers are concatenated as input for the subsequent layers. Consequently, the $l^{th}$ layer receives all feature maps produced by $[0, 1, ..., l-1]$ layers as inputs. The output of the $l^{th}$ layer is then defined as

$$x_l = H_l([x_0, x_1, ..., x_{l-1}]),\tag{2}$$

where $[x_0, x_1, \ldots, x_{l-1}]$ represent the concatenation of the feature maps.

To reduce the number of features and efficiently fuse the features from dense layers, a transition layer is added at the end of each dense block. The transition layer consists of a $1\times1$ convolution layer, a batch-normalization and a ReLU. The out of the transition layer is

$$x_t = H_t(H_l([x_0, x_1, ..., x_{l-1}])),\tag{3}$$

where $H_t$ is a non-liner transformation of transition layer. To further promote information propagation and make the network easier to optimize, we also employ residual connection into our block.

**Table 1** Detailed structure of the SIP-Net.

| | Feature Size | SIP-Net (k=32) |
|---|---|---|
| Input | 96×96×16×1 | |
| Convolution1 | 96×96×16×64 | 3×3×3 conv |
| Pooling | 48×48×8×64 | 2×2×2 avg. pool stride=2 |
| DRB1 | 48×48×8×192 | 1×1×1 conv 3×3×3 conv num=4 |
| TransLayer1 | 48×48×8×128 | 1×1×1 conv |
| Pooling | 24×24×4×128 | 2×2×2 avg. pool stride=2 |
| DRB2 | 24×24×4×384 | 1×1×1 conv 3×3×3 conv num=8 |
| TransLayer2 | 24×24×4×256 | 1×1×1 conv |
| Pooling | 12×12×2×256 | 2×2×2 avg. pool stride=2 |
| DRB3 | 12×12×2×768 | 1×1×1 conv 3×3×3 conv num=16 |
| TransLayer3 | 12×12×2×512 | 1×1×1 conv |
| Deconvolution1 | 24×24×4×256 | 3×3×3 conv stride=2 |
| DRB4 | 24×24×4×512 | 1×1×1 conv 3×3×3 conv num=8 |
| TransLayer4 | 24×24×4×256 | 1×1×1 conv |
| Deconvolution2 | 48×48×8×128 | 3×3×3 conv stride=2 |
| DRB5 | 48×48×8×256 | 1×1×1 conv 3×3×3 conv num=4 |
| TransLayer5 | 48×48×8×128 | 1×1×1 conv |
| Deconvolution3 | 96×96×16×64 | 3×3×3 conv stride=2 |
| DRB6 | 96×96×16×128 | 1×1×1 conv 3×3×3 conv num=2 |
| TransLayer6 | 96×96×16×64 | 1×1×1 conv |
| Convolution2 | 96×96×16×1 | 1×1×1 conv |

3.3 Attention-focused Module (AM)

To make the network focus more on the region to be segmented and to reduce
noise features from the surrounding region, we introduce an attention-focused
module in our model. The structure of attention-focused module is shown in
Fig. 3(C), which consists of a sigmoid layer and an element-wise multiplication
layer. The output of AM is the element-wise multiplication of input feature-
maps and attention masks. The attention masks are produced by sigmoid
layer:

$$M_t(x) = f(H_t(x)) \tag{4}$$

$$f(x) = \frac{1}{1 + e^{-x}} \tag{5}$$

where $M_t(x)$ denotes the attention mask, whose values range in $[0, 1]$, $H_t(x)$
denotes the feature map from long connection.

## 4 Experiments

To evaluate the performance of our proposed model, we applied the devel-
oped method on the MICCAI Prostate MR Image Segmentation 2012 Grant
Challenge dataset[2], TCIA Pancreas CT-82[3] and MICCAI 2017 Liver Tumor
Segmentation (LiTS) Challenge dataset[4] for image segmentation.

4.1 Implementation Details

The proposed method is implemented using the open source deep learning
library Keras. Our network is trained end-to-end by using the Stochastic Gra-
dient Descent (SGD) optimization method. In the training phase, the learning
rate is initially set to 0.0001 and decreased with a weight decay of 10e-6. The
momentum is set to 0.9. Experiments are carried out on a NVIDIA GTX
1080ti GPU with 11GB memory.

4.2 Prostate Segmentation From MR Image

We first evaluated our proposed method on MICCAI 2012 Prostate MR Image
Segmentation (PROMISE12) challenge dataset. There are in total 50 transver-
sal T2-weighted MR images of the prostate and the corresponding ground truth
segmentation, which were checked and corrected by a radiological resident with
more than 6 years of experience in prostate MRI. These images are a represen-
tative set of the types of MR images acquired in different hospitals. And these

---

[2]  https://promise12.grand-challenge.org/
[3]  https://wiki.cancerimagingarchive.net/display/Public/Pancreas-CT
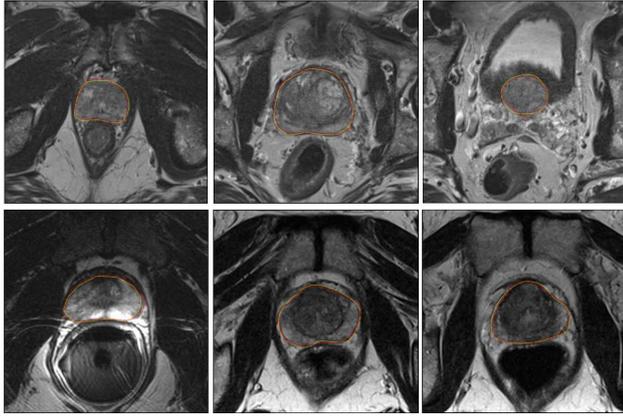[4]  https://competitions.codalab.org/competitions/17094

**Fig. 4** Segmentation results of the prostate from MR images. The yellow and red contours indicate the ground truth and our segmentation results, respectively. Note that these results are directly obtained from challenge website.

**Table 2** Quantitative evaluation results of the proposed method and other methods on prostate MR segmentation.

| User | | QuIIL | aslm | GeertLitjens | lanqier xl | tbrosch | **Ours** |
|---|---|---|---|---|---|---|---|
| ABD [mm] | Whole | 1.71 | 1.53 | 1.71 | 1.59 | 1.49 | **1.31** |
| | Base | 1.96 | 1.64 | 1.96 | 1.88 | 1.73 | **1.60** |
| | Apex | 1.62 | 1.93 | 1.56 | 1.67 | 1.73 | **1.39** |
| HD [mm] | Whole | 4.92 | 4.62 | 5.13 | 4.63 | 4.68 | **3.97** |
| | Base | 5.07 | **4.34** | 5.22 | 5.22 | 4.90 | 4.75 |
| | Apex | 3.97 | 5.16 | 4.17 | 4.26 | 4.49 | **3.70** |
| DSC [%] | Whole | 89.02 | 90.24 | 89.43 | 89.69 | 90.46 | **91.42** |
| | Base | 86.04 | 88.98 | 86.42 | 86.79 | 88.51 | **89.41** |
| | Apex | 86.39 | 83.31 | 86.81 | 86.79 | 85.29 | **88.51** |
| aRVD [%] | Whole | 7.26 | 7.98 | 6.95 | 7.58 | **6.59** | 6.97 |
| | Base | 13.57 | 12.68 | 11.04 | 11.63 | 9.64 | **8.53** |
| | Apex | 16.70 | 18.92 | 15.18 | 14.92 | 18.51 | **13.05** |
| **Overall score** | | 86.71 | 86.89 | 87.15 | 87.21 | 87.67 | **89.18** |

images are from multiple vendors and have different acquisition protocols and variations in voxel size, dynamic range, position, field of view and anatomic appearance. To evaluate the proposed algorithms, the organizers provide 30 testing MR images and the corresponding ground truth is held out.

Before training the network, we resampled all MR volumes into a fixed resolution of $0.625 \times 0.625 \times 1.5$mm and then normalized them as zero mean and unit variance. To facilitate network training, we applied data augmentation operations including rotation, scaling and flipping. During training, we adopted a random cropping strategy, where sub-volumes in the size of $16 \times 96 \times 96$ $(d \times w \times h)$ voxels are randomly cropped from the training data during every iteration. In the testing phase, similar to the works in [34, 46], we used overlapping sliding windows to crop sub-volumes and used the average of the probability maps of these sub-volumes to get the whole volume prediction.

**Table 3** Performance of CNN based CT pancreas segmentation methods, which are trained and evaluated using the same number of training and testing images.

| Methods | DSC [%] |
|---|---|
| Holistically Nested 2D FCN Stage-1 [50] | 76.8 ± 11.1 |
| Holistically Nested 2D FCN Stage-2 [50] | 81.2 ± 7.3 |
| 2D FCN [51] | 80.3 ± 9.0 |
| 2D FCN + Recurrent Network [51] | 82.3 ± 6.7 |
| Single Model 2D FCN [52] | 75.7 ± 10.5 |
| Multi-Model 2D FCN [52] | 82.2 ± 5.7 |
| Attention U-Net [53] | 81.5 ± 6.2 |
| SIP-Net (Our) | **83.9 ± 4.5** |

The sub-volume size was also 16×96×96 and the stride was 8×48×48. Due to the limitation of the memory, we used the mini-batch size of 4. The number of parameters of SIP-Net was 3.16M, and the prediction time was approximately 1 minutes for one MR volume.

Several sample results of our proposed method are shown in Fig. 4. It can be seen that our model can accurately segment the prostate and obtain smooth and continuous prostate boundaries. Quantitative evaluation was also performed. The evaluation metrics used in PROMISE12 challenge include Dice Similarity Coefficient (DSC), percentage of the absolute difference between the volumes (aRVD), average over the shortest distance between the boundary points of the volumes (ABD) and Hausdorff Distance (HD). All the evaluation metrics are calculated in 3D. In addition to evaluating these metrics over the entire prostate segmentation, the challenge organizers also calculated the boundary measures specifically for the apex and base parts of the prostate, because those parts are difficult to segment but in the same time very important for many clinical applications. The apex and base the prostate are determined by dividing the prostate into three approximately equal sized parts along the axial direction (the first 1/3 as apex and the last 1/3 as base). Then an overall score will be computed by taking all the criteria into consideration rank the algorithms.

The results of our proposed method and the competitors are shown in Table 2. Only the top 10 teams are listed. Note that all the results reported in this section were obtained directly from the challenge website. As it can be seen from the table, our overall performance was the best and therefore ranked the first place among all the teams (by May 22, 2018)[5] with the score of 89.18. From Table 2, it can be seen that our proposed model achieved the best performance in several measures. The segmentation results of our model were the best not only for whole prostate segmentation, but also in the base and apex areas, which demonstrates the effectiveness of the proposed 3D model with DRBs and AM modulated long connections.
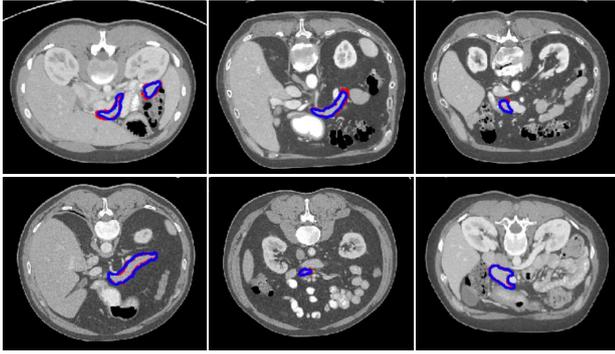
---

[5] https://promise12.grand-challenge.org/evaluation/results/

**Fig. 5** Sample segmentation results of the pancreas CT images. The red and blue contours are the ground truth and our segmentation results, respectively.

### 4.3 Pancreas Segmentation

The proposed model is also evaluated on another publicly available dataset – TCIA Pancreas CT-82. This dataset contains 82 contrast enhanced 3D CT scans, which have resolutions of $512\times512$ pixels with varying pixel sizes and slice thickness between 1.5-2.5 mm, acquired on Philips and Siemens MDCT scanners [54]. The dataset is publicly available and commonly used to benchmark CT pancreas segmentation frameworks. In our experiments, the 82 scans are randomly split with 62 images for training and 20 images for testing. Before training the model, we resampled all volumes into a fixed resolution of 1.0mm$\times$1.0mm$\times$1.0mm. Then all the scans are normalized to have zero mean and unit variance. We again applied data augmentation operations including rotation, scaling and flipping. We also employed the random cropping strategy, where sub-volumes in the size of $64\times96\times96$ ($d \times w \times h$) voxels are randomly cropped from the training data during every iteration. In the testing phase, we used overlapping sliding windows to crop sub-volumes and used the average probability maps of these sub-volumes to get the whole volume prediction. The sub-volume size was also $64\times96\times96$ and the stride was $32\times48\times48$. The architecture of network was same as that utilized on prostate segmentation. The prediction time was approximately 1 minutes for one CT volume.

To evaluate the proposed architecture, we compare the performance of the model against other state-of-the-art CT pancreas segmentation methods. The results are summarized in Table 3. It can be seen that our proposed model achieved 83.9 $\pm$ 4.51 in DSC for pancreas labels, which outperform other state-of-the-art methods. Several example segmentation results of our proposed method are shown in Fig. 5. Our proposed model can accurately segment the pancreas from CT images. It is worth noting that we only employ a single model to segment pancreas and our model does not require multiple CNN models as in [50].

**Table 4** Quantitative evaluation results of the proposed method and other methods on MICCAI 2017 LiTS Challenge Dataset.

| Methods | Per Case DSC [%] | Global DSC [%] |
| --- | --- | --- |
| H-DenseUNet [36] | 96.1 | 96.5 |
| CascadedResNet[55] | **96.3** | **96.7** |
| SIP-Net(2D) (Ours) | 95.9 | 96.3 |
| SIP-Net(3D) (Ours) | 94.2 | 94.6 |

4.4 Liver Segmentation

We also tested our proposed model on the competitive dataset of MICCAI 2017 LiTS Challenge, which contains 131 contrast enhanced 3D abdominal CT scans with radiologist hand-drawn ground truths for training and the rest 70 used for testing with unreleased ground truth. Since the data were acquired from different clinical sites, which have different scanners and protocols, the scans have largely varying in-plane resolution (0.55mm-1.0mm) and slice spacing (0.45mm-6.0mm). Before training the model, we truncated the image intensity values of all scans to the range of [-200,200] to remove the irrelevant details and then normalized each volume. In addition to 3D model, we also evaluate the 2D model with same network structs for evaluating the influence of parameters.
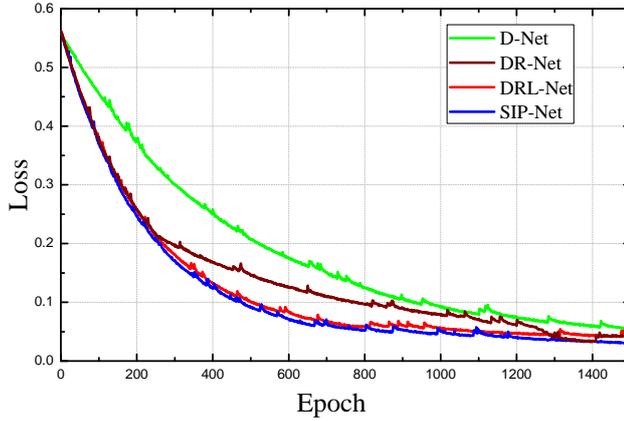
During the network training, we randomly cropped patches in the size of 224×224×16 pixels for 3D model (224×224 pixels for 2D model) from the training data during every iteration. In the testing phase, we used overlapping sliding windows to crop sub-volumes and used the average probability maps of these sub-volumes to get the whole volume prediction. The cropped size was also 224×224×16 pixels for 3D model (224×224 pixels for 2D model) and the stride was 112×112×8 for 3D model (112×112 for 2D model). The number of parameters of SIP-Net (2D) was 1.43M, and the prediction time was approximately 2 minutes for one CT volume.

There were more than 60 submissions for the MICCAI LiTS Challenge. The segmentation performances of the teams are listed on the leaderboard[6] and we were among the top seven teams (by November 15, 2018, team of Qikui_sigma-RPI). We compared the performance of our model with two published top-performance models: H-DenseUNet [36] and CascadedResNet [55]. H-DenseUNet employed a simple ResNet to process the original data, which makes the network subject to the performance of pro-processing. In addition, H-DenseUNet employed 3D convolutional layers inside the model with much more parameters and thus increased training difficulty. CascadedResNet, on the other hand, achieved good results but took approximately 7 days on two Titan X GPUs for training. Our proposed method performs similarly to the above two approaches with negligible differences, however, can be trained much more efficiently than those methods. Comparing the performance of 2D and 3D models reveals that the 2D model can even obtain better performance.

---

[6] https://competitions.codalab.org/competitions/17094#results

**Table 5** Performance of the proposed model in different configurations.

| Configurations | Global DSC [%] |
|---|---|
| D-Net | 86.0 |
| DR-Net | 86.9 |
| DRL-Net | 88.8 |
| SIP-Net (Ours) | **89.8** |



**Fig. 6** Training loss of networks with different structures.

This indicates that the network architecture is the key for the performance gain and a larger number of network parameters may lead to performance decrease.

## 5 Discussions

In this section, we provide in-depth discussions of the effects of some of our proposed components.

5.1 Ablation Study of Network Structure

In order to evaluate the effectiveness of the residual connections in dense blocks, the long connections and attention-focused modules used in our model, we performed a set of ablation study experiments. The prostate MR image dataset was used. We randomly selected 10 patients for validation and the rest 40 patients were utilized for training.

To analyze the learning behaviors of our model, we created four different configurations of our model: using only dense block (D-Net), using only DRBs (DR-Net), using DRBs and long connections (DRL-Net), using DRBs, long connections and attention-focused module (SIP-Net). We first analyzed
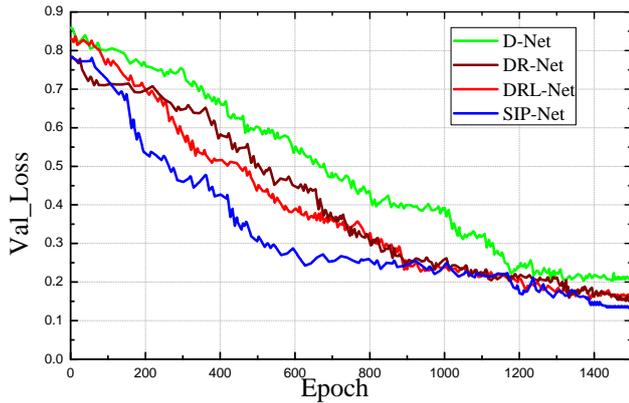
**Fig. 7** Validation loss of the networks with different structures.

the leaning behaviors of these models. Figs. 6 and 7 present the training and validation losses of different networks. It is observed that the models with either residual connections, long connections and attention-focused module converge faster and achieve lower validation loss than the one with only dense block, which demonstrates that the use of residual connections, long connections and attention-focused modules can improve the training efficiency and the performance of the models. Fig. 7 further shows that the long connections can accelerate the convergence speed and alleviate the risk of over-fitting on limited training data.

Table 5 shows the performance of our proposed model with different connections and blocks. It is can be seen that adding residual connections, long connections and attention-focused modules can achieve better Dice scores than the network with only dense blocks. The network with residual connections and dense block has marginally better performance than that with only dense block, which demonstrates that the enhanced information propagation inside each block can improve the performance of the model. The model with long connections obtained better performance than the one without. It is conceivable that enhancing information propagation both locally and globally inside the model and combining them together can further improve the performance. The network with attention-focused modules achieves the best performance in the ablation experiments, indicating that attention-focused module further improves the performance of model.

To demonstrate the efficiency of the proposed method in utilizing training data, we compare the performance of the model against that of FCN, which is indeed the version of our model without DRB and AM, using different amount of training data. In this experiment, we respectively used 40%, 50%, 60%, 70%, and 80% of data for training and reserved upto 20% of the data for testing. To avoid potential data distribution bias, in each setting, we randomly selected five different subsets from the entire dataset for training and testing. The average performances over the five runs under each setting are reported
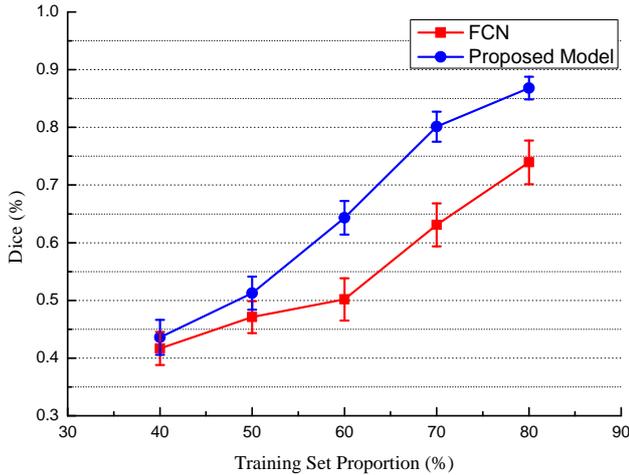
**Fig. 8** Performance of the proposed model and FCN under different training set proportions.

and shown in Fig. 8. It can be seen that, when only 40% of the training data were used, the proposed method and FCN achieved similar performance. The performances are poor as the training data is very limited in that case as we expected. As the size of the training dataset increases, both methods start to perform better. However, Fig. 8 shows that the proposed method improves in a much faster rate, with the contribution from the proposed DRB and AM modules. Eventually, the proposed model only needs less than 70% of the training data to outperform the FCN trained with the entire 80% of the data. The experiment demonstrates that the proposed structures can help deep CNNs get trained more efficiently with small number of training images, which is a very desired property for medical imaging applications where labeled data is usually in scarce.

## 5.2 Analysis of Attention-focused Modules

To further analyze the function of attention-focused modules, we visualized the generated attention masks in the up-sampling path. Four different types of input images were selected, which are selected from base, mid-gland, apex and also outside of the prostate. It can be seen that the attention masks have much higher weights in the prostate region than in the non-prostate region as shown in Fig. 9. And the shape of attention mask was very close to the ground truth. It is conceivable that higher weight was inside the attention masks, which helps to locate the region of prostate. The shape of the attention mask volume was again close to the ground truth. It suggests that the attention mask can help the network pay more attention to the region of prostate and suppress the features from the non-prostate region towards better image segmentation.
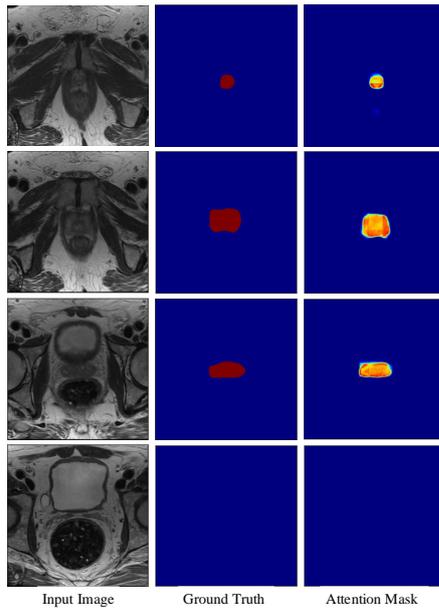
Input Image　　　　Ground Truth　　　　Attention Mask

**Fig. 9** Attention mask examples produced by the attention-focused module. The blue, red pixel represents background and prostate, respectively. And the attention mask is the corresponding heat map produced by attention-focused modules. For the heat map, the darker the color, the greater the weight value, the lighter the color, the smaller the weight value.

**Table 6** Performance of SIP-Net in different batch size.

| Batch-size | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| DSC [%] | 89.4 | 89.5 | 89.7 | 89.8 |

### 5.3 Effects of Batch Size

To evaluate the influence of batch size on the segmentation results, we compared the performance of our proposed model under various batch size. The prostate MR image dataset also was used, 10 patients were randomly selected for validation and the rest 40 patients were utilized for training. The segmentation performance is listed in Table 6. It can be seen that the size of batch has a slight effect on the segmentation results and the model performed the best when batch size is 4.

## 6 Conclusions

In this paper, we first prove that not all the feature maps transmitted by skip connections contribute positively to the network performance. And to adaptive select information passed through those skip connections, we propose a novel

network, named SIP-Net, which can adaptive select the information passed through those skip connections by our proposed attention-focused modules. Expect for making the skip connections between the down-sampling path and up-sampling path can further improve the context and gradient information propagation both forward and backward and address the vanishing-gradient problem, our proposed SIP-Net also makes the model focus on the region of interest. Extensive experiments on the publicly available MICCAI Prostate MR Image Segmentation 2012 Grant Challenge dataset, TCIA Pancreas CT-82 and MICCAI 2017 Liver Tumor Segmentation (LiTS) Challenge dataset demonstrate that our proposed method can get more accurate boundaries and achieve superior results compared with other state-of-the-art methods.

## Conflict of interest

The authors declare that they have no conflict of interest.

## References

1. Yongchun Zhu, Fuzhen Zhuang, Jindong Wang, Jingwu Chen, Zhiping Shi, Wenjuan Wu, and Qing He. Multi-representation adaptation network for cross-domain image classification. *Neural Networks*, 119:214–221, 2019.
2. Tobias Heimann, Bram Van Ginneken, Martin A Styner, Yulia Arzhaeva, Volker Aurich, Christian Bauer, Andreas Beck, Christoph Becker, Reinhard Beichel, György Bekes, et al. Comparison and evaluation of methods for liver segmentation from CT datasets. *IEEE Transactions on Medical Imaging*, 28(8):1251–1265, 2009.
3. Shu Liao, Yaozong Gao, Aytekin Oto, and Dinggang Shen. Representation learning: a unified deep learning framework for automatic prostate MR segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 254–261. Springer, 2013.
4. Qikui Zhu, Bo Du, and Pingkun Yan. Boundary-weighted domain adaptive neural network for prostate mr image segmentation. *IEEE transactions on medical imaging*, 39(3):753–763, 2019.
5. Qikui Zhu, Bo Du, Baris Turkbey, Peter Choyke, and Pingkun Yan. Exploiting interslice correlation for mri prostate image segmentation, from recursive neural networks aspect. *Complexity*, 2018, 2018.
6. Bo Du, Qiuci Wei, and Rong Liu. An improved quantum-behaved particle swarm optimization for endmember extraction. *IEEE Transactions on Geoscience and Remote Sensing*, 57(8):6003–6017, 2019.
7. Jia Wu, Zhibin Hong, Shirui Pan, Xingquan Zhu, Chengqi Zhang, and Zhihua Cai. Multi-graph learning with positive and unlabeled bags. In *Proceedings of the 2014 SIAM International Conference on Data Mining*, pages 217–225. SIAM, 2014.
8. Xiaojun Bi and Haibo Wang. Early alzheimers disease diagnosis based on eeg spectral images using deep learning. *Neural Networks*, 114:119–135, 2019.
9. Xue Li, Bo Du, Chang Xu, Yipeng Zhang, Lefei Zhang, and Dacheng Tao. Robust learning with imperfect privileged information. *Artificial Intelligence*, 282:103246, 2020.
10. Jia Wu, Xingquan Zhu, Chengqi Zhang, and S Yu Philip. Bag constrained structure pattern mining for multi-graph classification. *Ieee transactions on knowledge and data engineering*, 26(10):2382–2396, 2014.
11. Zengmao Wang, Bo Du, and Yuhong Guo. Domain adaptation with neural embedding matching. *IEEE Transactions on Neural Networks and Learning Systems*, 2019.

12. Fanzhen Liu, Shan Xue, Jia Wu, Chuan Zhou, Wenbin Hu, Cecile Paris, Surya Nepal, Jian Yang, and Philip S Yu. Deep learning for community detection: Progress, challenges and opportunities. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*, pages 4981–4987, 2020. Survey track.

13. Jia Wu, Zhihua Cai, Sanyou Zeng, and Xingquan Zhu. Artificial immune system for attribute weighted naive bayes classification. In *The 2013 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2013.

14. Xianjing Qin, Xuelong Li, Yang Liu, Hongbing Lu, and Pingkun Yan. Adaptive shape prior constrained level sets for bladder MR image segmentation. *IEEE journal of Biomedical and Health Informatics*, 18(5):1707–1716, 2014.

15. Y. Huo, J. Liu, Z. Xu, R. L. Harrigan, A. Assad, R. G. Abramson, and B. A. Landman. Robust multicontrast MRI spleen segmentation for splenomegaly using multi-atlas segmentation. *IEEE Transactions on Biomedical Engineering*, 65(2):336–343, 2018.

16. Chris McIntosh and Thomas G Purdie. Contextual atlas regression forests: multiple-atlas-based automated dose prediction in radiation therapy. *IEEE Transactions on Medical Imaging*, 35(4):1000–1012, 2016.

17. Hancan Zhu, Hewei Cheng, Xuesong Yang, and Yong Fan. Metric learning for label fusion in multi-atlas based image segmentation. In *13th International Symposium on Biomedical Imaging (ISBI)*, pages 1338–1341. IEEE, 2016.

18. Qinquan Gao, Akshay Asthana, Tong Tong, Daniel Rueckert, et al. Multi-scale feature learning on pixels and super-pixels for seminal vesicles MRI segmentation. In *Medical Imaging 2014: Image Processing International Society for Optics and Photonics*, volume 9034, page 903407, 2014.

19. Pingkun Yan, Sheng Xu, Baris Turkbey, and Jochen Kruecker. Adaptively learning local shape statistics for prostate segmentation in ultrasound. *IEEE Transactions on Biomedical Engineering*, 58(3):633–641, 2011.

20. O. Gloger, K. Tönnies, R. Laqua, and H. Vjlzke. Fully automated renal tissue volumetry in mr volume data using prior-shape-based segmentation in subject-specific probability maps. *IEEE Transactions on Biomedical Engineering*, 62(10):2338–2351, 2015.

21. Fulin Luo, Liangpei Zhang, Bo Du, and Lefei Zhang. Dimensionality reduction with enhanced hybrid-graph discriminant learning for hyperspectral image classification. *IEEE Transactions on Geoscience and Remote Sensing*, 2020.

22. Xiaoqiang Lu, Yaxiong Chen, and Xuelong Li. Hierarchical recurrent neural hashing for image retrieval with hierarchical convolutional features. *IEEE Transactions on Image Processing*, 27(1):106–120, 2018.

23. Jia Wu, Shirui Pan, Xingquan Zhu, and Zhihua Cai. Boosting for multi-graph classification. *IEEE transactions on cybernetics*, 45(3):416–429, 2014.

24. Xue Li, Bo Du, Yipeng Zhang, Chang Xu, and Dacheng Tao. Iterative privileged learning. *IEEE transactions on neural networks and learning systems*, 2019.

25. Fulin Luo, Liangpei Zhang, Xiaocheng Zhou, Tan Guo, Yanxiang Cheng, and Tailang Yin. Sparse-adaptive hypergraph discriminant analysis for hyperspectral image classification. *IEEE Geoscience and Remote Sensing Letters*, 2019.

26. Siwei Lai, Liheng Xu, Kang Liu, and Jun Zhao. Recurrent convolutional neural networks for text classification. In *AAAI*, volume 333, pages 2267–2273, 2015.

27. Yoav Goldberg. A primer on neural network models for natural language processing. *Journal of Artificial Intelligence Research*, 57:345–420, 2016.

28. Rie Johnson and Tong Zhang. Semi-supervised convolutional neural networks for text categorization via region embedding. In *Advances in Neural Information Processing Systems*, pages 919–927, 2015.

29. Bum-Chae Kim, Jee Seok Yoon, Jun-Sik Choi, and Heung-Il Suk. Multi-scale gradual integration CNN for false positive reduction in pulmonary nodule detection. *Neural Networks*, 115:1–10, 2019.

30. Jia Wu, Xingquan Zhu, Chengqi Zhang, and Zhihua Cai. Multi-instance multi-graph dual embedding learning. In *2013 IEEE 13th International Conference on Data Mining*, pages 827–836. IEEE, 2013.

31. Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.

32. Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 3431–3440, 2015.

33. Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. V-net: Fully convolutional neural networks for volumetric medical image segmentation. In *2016 Fourth International Conference on 3D Vision (3DV)*, pages 565–571. IEEE, 2016.

34. Lequan Yu, Xin Yang, Hao Chen, Jing Qin, and Pheng-Ann Heng. Volumetric convnets with mixed residual connections for automated prostate segmentation from 3D MR images. In *AAAI*, pages 66–72, 2017.

35. Eli Gibson, Francesco Giganti, Yipeng Hu, Ester Bonmati, Steve Bandula, Kurinchi Gurusamy, Brian Davidson, Stephen P Pereira, Matthew J Clarkson, and Dean C Barratt. Automatic multi-organ segmentation on abdominal CT with dense V-networks. *IEEE Transactions on Medical Imaging*, 2018.

36. Xiaomeng Li, Hao Chen, Xiaojuan Qi, Qi Dou, Chi-Wing Fu, and Pheng-Ann Heng. H-DenseUNet: Hybrid densely connected unet for liver and tumor segmentation from CT volumes. *IEEE Transactions on Medical Imaging*, 2018.

37. Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image Computing and Computer-Assisted Intervention*, pages 234–241. Springer, 2015.

38. Yi Li, Haozhi Qi, Jifeng Dai, Xiangyang Ji, and Yichen Wei. Fully convolutional instance-aware semantic segmentation. *arXiv preprint arXiv:1611.07709*, 2016.

39. Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.

40. Zifeng Wu, Chunhua Shen, and Anton van den Hengel. High-performance semantic segmentation using very deep fully convolutional networks. *arXiv preprint arXiv:1604.04339*, 2016.

41. Qikui Zhu, Bo Du, Baris Turkbey, Peter L Choyke, and Pingkun Yan. Deeply-supervised CNN for prostate segmentation. In *2017 International Joint Conference on Neural Networks (IJCNN)*, pages 178–184. IEEE, 2017.

42. Aliasghar Mortazi, Rashed Karim, Kawal Rhode, Jeremy Burt, and Ulas Bagci. Cardiacnet: Segmentation of left atrium and proximal pulmonary veins from MRI using multiview CNN. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 377–385. Springer, 2017.

43. Xiao Han. Automatic liver lesion segmentation using a deep convolutional neural network method. *arXiv preprint arXiv:1704.07239*, 2017.

44. Zhuowen Tu. Auto-context and its application to high-level vision tasks. In *IEEE Conference on Computer Vision and Pattern Recognition, 2008.*, pages 1–8. IEEE, 2008.

45. Hao Chen, Qi Dou, Lequan Yu, and Pheng-Ann Heng. Voxresnet: Deep voxelwise residual networks for volumetric brain segmentation. *arXiv preprint arXiv:1608.05895*, 2016.

46. Lequan Yu, Jie-Zhi Cheng, Qi Dou, Xin Yang, Hao Chen, Jing Qin, and Pheng-Ann Heng. Automatic 3D cardiovascular MR segmentation with densely-connected volumetric convnets. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 287–295, 2017.

47. Fei Wang, Mengqing Jiang, Chen Qian, Shuo Yang, Cheng Li, Honggang Zhang, Xiaogang Wang, and Xiaoou Tang. Residual attention network for image classification. *arXiv preprint arXiv:1704.06904*, 2017.

48. Chen-Yu Lee, Saining Xie, Patrick Gallagher, Zhengyou Zhang, and Zhuowen Tu. Deeply-supervised nets. In *Artificial Intelligence and Statistics*, pages 562–570, 2015.

49. Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *CVPR*, volume 1, page 3, 2017.

50. Holger R Roth, Le Lu, Nathan Lay, Adam P Harrison, Amal Farag, Andrew Sohn, and Ronald M Summers. Spatial aggregation of holistically-nested convolutional neural networks for automated pancreas localization and segmentation. *Medical Image Analysis*, 45:94–107, 2018.

51. Jinzheng Cai, Le Lu, Yuanpu Xie, Fuyong Xing, and Lin Yang. Improving deep pancreas segmentation in CT and MRI images via recurrent neural contextual learning and direct loss function. *arXiv preprint arXiv:1707.04912*, 2017.

52. Yuyin Zhou, Lingxi Xie, Wei Shen, Yan Wang, Elliot K Fishman, and Alan L Yuille. A fixed-point model for pancreas segmentation in abdominal CT scans. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 693–701. Springer, 2017.
53. Ozan Oktay, Jo Schlemper, Loic Le Folgoc, Matthew Lee, Mattias Heinrich, Kazunari Misawa, Kensaku Mori, Steven McDonagh, Nils Y Hammerla, Bernhard Kainz, et al. Attention U-Net: Learning where to look for the pancreas. *arXiv preprint arXiv:1804.03999*, 2018.
54. Holger R Roth, Amal Farag, Evrim B Turkbey, Le Lu, Jiamin Liu, and Ronald M Summers. Data from pancreas-CT. *The Cancer Imaging Archive. http://doi.org/10.7937/K9/TCIA.2016.tNB1kqBU*, 2016.
55. Lei Bi, Jinman Kim, Ashnil Kumar, and Dagan Feng. Automatic liver lesion detection using cascaded deep residual networks. *arXiv preprint arXiv:1704.02703*, 2017.