

2D-3D Fully Convolutional Neural Networks for Cardiac MR Segmentation

Jay Patravali, Shubham Jain, and Sasank Chilamkurthy

Qure.ai
jay.patravali@qure.ai
shubham.jain@qure.ai
sasank.chilamkurthy@qure.ai
www.qure.ai

Abstract. In this paper, we develop a 2D and 3D segmentation pipelines for fully automated cardiac MR image segmentation using Deep Convolutional Neural Networks (CNN). Our models are trained end-to-end from scratch using the ACD Challenge 2017 dataset comprising of 100 studies, each containing Cardiac MR images in End Diastole and End Systole phase. We show that both our segmentation models achieve near state-of-the-art performance scores in terms of distance metrics and have convincing accuracy in terms of clinical parameters. A comparative analysis is provided by introducing a novel dice loss function and its combination with cross entropy loss. By exploring different network structures and comprehensive experiments, we discuss several key insights to obtain optimal model performance, which also is central to the theme of this challenge.

Keywords: Deep Learning, Medical Image Analysis, Computer Vision, MR Segmentation

1 Introduction

MR imaging is an effective non-invasive procedure for diagnosis and treatment of known or suspected Cardiac diseases. Cardiac MR images can produce highly detailed pictures of different structures within the heart. Delineation of these structures can provide relevant diagnostic information and evaluate the overall functioning of the heart. Segmentation of left ventricle, right ventricle and the myocardium can be used to calculate relevant diagnostics parameters such as ejection fraction and myocardial mass. Due to massive volumes of cardiac image data, relying on manual delineations can be a time-consuming process, often prone to error and rater variability. Hence there is a critical need for accurate, reproducible and fully-automated methods for cardiac segmentation.

In recent works, Deep Learning and Convolutional Neural Networks (CNNs) have shown tremendous progress in fully-automated segmentation tasks. The growing success of CNNs in solving computer vision problems such as image recognition and classification [1,2] can be attributed to its ability in learning a hierarchical representation of the input data, without relying on hand-crafted features. Deep learning techniques for segmentation have defined the state-of-the-art using Fully Convolutional Networks (FCN) [3]. The idea behind FCN is to use a contracting path to extract features at different spatial scales followed by an expanding path to upsample and increase the spatial resolution of learned features.

For segmentation in medical images, U-Net[6] is a well established 2D CNN architecture that builds upon the FCN. By adding skip connections between the contracting and expanding paths, the U-Net model showed reasonable segmentation accuracy with very few training samples. Cardiac segmentation based on original FCN have been proposed [4] with modifications to make it faster and memory efficient [5]. As raw 3D MR

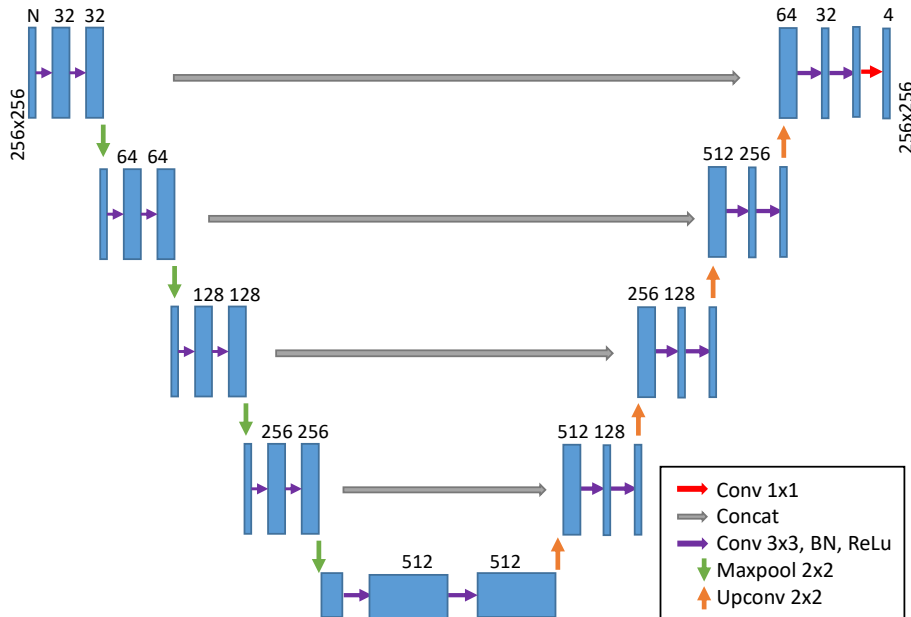


Fig. 1. 2D Model Architecture.

volumes are fed slice-by-slice as inputs to these 2D CNN models, they fail to capture the spatial contextual information required to segment the whole heart. To that end, U-Net3D[7] extends the 2D U-Net model by replacing its 2D convolutional operations with its 3D counterparts. In similar way, V-Net[8] employs a 3D CNN model with a novel dice loss function showing convincing results in medical image segmentation. To our best knowledge, there are very few methods that have applied 3D CNNs for Cardiac Segmentation and have obtained satisfactory performance.

In our work, we develop a fully-automated 2D and 3D CNN models designed to segment the Left Ventricle, Right Ventricle and Myocardium. This segmentation task is part of the Automatic Cardiac Detection Challenge 2017 [9]. The 2D segmentation model is trained slice-by-slice, whereas we compute volumetric segmentation for the 3D model. Our models are easy to implement, have modular architecture, and relatively short training and testing times. We introduce a new dice loss function, and compare its performance with traditional cross entropy loss and combined cross entropy-dice loss. Through our experiments we also compare and analyze the performance of our 2D and 3D models, both which achieve near state-of-the-art accuracy scores in terms of geometric metrics and clinical validity.

2 Method

2.1 Network Architecture

We develop a 2D segmentation model architecture that is adapted from U-Net [6] as illustrated in Fig. 1. On left side is the “contracting” stage and on the right side is “expanding” stage. At the bottom is a base layer. We provide an option to feed varying number (N) of image slices that can be passed as input channels to the model. Here, N can be 1 for single image slice or more. Every step on the contracting path consists of a series of a 3x3 convolutions (conv 3x3), batch normalization (bn) [10], rectified linear unit (ReLU) and conv 3x3 in a sequence that forms a *conv_bn_relu* block. Two of

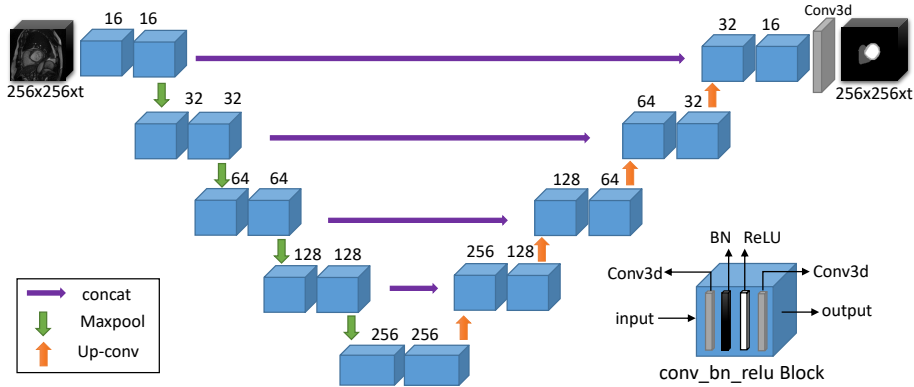


Fig. 2. 3D Model Architecture.

conv_bn_relu blocks in succession forms a *Conv* block that doubles the number of feature channels. The contracting path downsamples the image with a 2×2 maxpool operation of stride 2. Similar to contracting stage, every step in expanding stage has a sequence of conv 3×3 , bn, ReLU and conv 3×3 in a series of two consecutive blocks. The images are upsampled with a 2×2 up-convolution (upconv 2×2) with stride 2. For upsampling of images, a sequence of 2×2 up-convolution (upconv 2×2) with stride 2, concatenations and a *Conv* block forms a *deconv* block. A final 1×1 convolution layer maps the 32 feature channels to 4 classes.

Our 3D model is an extension to our 2D model with few modifications, and finds similarities with U-Net 3d [7]. First, we replace all 2D operations with its 3D counterparts. Second, due to memory constraints and less number of training examples we limit the maximum number of feature maps to 256. Overall since the number of slices are less across the dataset (9 avg.), we apply 3D $1 \times 2 \times 2$ maxpooling operation only in X and Y leaving out the Z dimension. This allows our 3D model to accept input volumes of varying slices at training or inference stage. Similar to our 2D model, every step in contracting and expanding stage consists of two repeating blocks, where each block is a sequence of conv $3 \times 3 \times 3$, bn, ReLU, conv $3 \times 3 \times 3$. Due to symmetric nature of the model, we can simultaneously add, remove or modify blocks across both paths. Additionally, the size of input data and final outputs images remains the same for both our 2D and 3D models. Thus, we are able to maintain modularity for faster experimentation and modification at the time of training and design. To prevent overfitting, we add dropout layers [11] with probability values of 0.5 in the last and 0.3 in the second last layer of the contracting stage, in both 2D and 3D models.

2.2 Dataset, Preprocessing and Augmentation

Our models are trained end-to-end from scratch using MICCAI’s ACD Challenge 2017 dataset. It contains 150 exams of fully-annotated cardiac MRI’s. Out of these 100 are used for training phase and 50 for testing phase of this challenge. These exams are obtained from multiple patients, each consisting of scans from End Diastole and End Systole phase taken in short axis orientation.

Since the data acquisition can bring inconsistencies in dataset, its necessary to carry preprocessing steps to ensure that the model receives uniform inputs. To remove noise and enhance contrast, we use contrast limited adaptive histogram localization (CLAHE) [12]. Then, we normalize the intensity values of all images between the range of 1-99 percentiles. Finally, we clip the image pixel values between 0 and 1.

To ensure both 2D and 3D model perceive heart features in similar proportion, we do a resampling operation across all input volumes to a common voxel spacing of 1.5x1.5x10 mm. For 2D segmentation we resize and crop the images to fixed size of 256x256. Due to maxpool operation applied only on height and width in our 3D segmentation model, when testing the 3D model, we can feed it image volumes of varying number of slices. At training, the 3D model is fed with raw input volumes that are resized and cropped to 256x256x12.

We apply a light data augmentation techniques on-the-fly to efficiently feed the input data volumes into our model. On a random basis, the data is rotated between -15 to +15 degree, and scaled between 0.9 - 1.1 range. This ensures slight robustness and variability in training the network.

2.3 Training

We train the 2D segmentation model by feeding it raw input images slice-by-slice. Whereas the 3D segmentation is trained by feeding it with entire 3D input volumes. Both the 2D and 3D models are trained with different optimization functions described more in detail in the following subsection.

During the training process, weights are updated using stochastic gradient descent with a momentum of 0.99. The initial learning rate is decayed by a factor of 10 every 30 epochs. For the training phase of this challenge, we do a 5-fold cross validation, which leaves 80 patients for training and 20 for validation. Training is completed after 300 epochs. Best models are checkpointed and stored for testing. For 2D segmentation we use a batch size of 8, whereas for 3D segmentation we use a batch size of 4.

Our models are implemented entirely using the PyTorch [13] framework, due to its flexibility in experimentation. We run all experiments on a standard workstation equipped with 64 GB of memory, Intel(R) Core(TM) i7-6700K CPU clocking at 4.00GHz, with a 12 GB NVidia Titan X Pascal GPU.

2.4 Optimization Function

In this section we describe three optimization functions that are used for training our 2D and 3D segmentation models. We use these functions to compare the performance of 2D and 3D models. At training, we apply a pixel-wise softmax activation in the final layer of the model to get the predicted probabilities $p(x, i)$ for each class i at each pixel x . The targets at location x is denoted by $t(x)$.

Cross Entropy Loss In segmentation tasks, the standard practice is to apply cross entropy loss function to measure the pixel-wise probability error between the predicted output and target and sum the errors across all the pixels. In addition to this, we apply weights to each class (w_i for class i) to offset the imbalance of pixel frequency across different classes. Concretely, the weighted cross entropy loss L_{CE} is defined as,

$$L_{CE} = - \sum_x w_{t(x)} \log(p(x, t(x))) \quad (1)$$

Dice Loss The Dice's Coefficient is a metric to measure the similarity between two given samples. Extending it as a loss function as shown in [8], improves the performance when dealing with situations where background pixels are higher than the labels. Here, we introduce a novel dice loss L_{dice} that is a logarithmic value of the dice score, making it easier to optimize. Similar to weighted cross entropy, we use weights to offset the class imbalance. Dice loss L_{dice} is weighted sum of dice losses l_i for each class i is given as,

$$L_{\text{dice}} = \sum_i w_i l_i \quad (2)$$

For class i , let's denote its binary map by t_i . i.e,

$$t_i(x) = \begin{cases} 1, & \text{if } t(x) = i \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

Then, dice loss for class i is given by,

$$l_i = \log \left(2 - \frac{\sum_x t_i(x)p(x, i) + \epsilon}{\sum_x t_i(x) + p(x, i) + \epsilon} \right) \quad (4)$$

Combined Dice Cross Entropy Loss While cross-entropy loss optimizes for pixel-level accuracy, the Dice loss function enhances the segmentation quality. Combining these two objective functions, we define a weighted average of cross entropy L_{CE} and dice loss function L_{dice} formulated as Cross Entropy-Dice Loss $L_{\text{CE} + \text{dice}}$ in Eqn 5. Here, λ_{CE} and λ_{dice} are weight parameters for cross entropy loss and Dice loss function respectively.

$$L_{\text{CE} + \text{dice}} = \lambda_{\text{CE}} * L_{\text{CE}} + \lambda_{\text{dice}} * L_{\text{dice}} \quad (5)$$

3 Results

In this section we evaluate the performance of the proposed 2D and 3D segmentation models in terms of geometric or distance metrics and clinical metric scores for all 100 studies provided in the training phase of ACDC 2017 contest. Table 1 and Table 2 presents the distance metric scores for our 2D model and 3D model respectively. For distance metric, we utilize the Dice Score and the Hausdorff Distance to measure the accuracy of segmented Left Ventricle (LV), Right Ventricle (RV) and Myocardium (MYO) in both end diastole (ED) and end systole (ES) phase. For each model, we compare the performance of the three optimization functions namely, the Cross Entropy Loss (CE Loss), Dice Loss and Combined Cross Entropy-Dice Loss (Dice-CE Loss). We observe that our proposed dice loss function outperforms CE Loss and CE-Dice Loss functions across all metrics in both 3D and 2D models. The Hausdorff distances in 3D models is observed to be much higher than 2D models, due to false positives as far-off speckles in 3D space. Illustration of results for both 2D and 3D models are provided in Fig. 3 and Fig. 4. Overall, we achieve near equal distance metric scores when compared to [14].

For clinical metrics, we use Correlation Coefficient (CC), Bias and Limits of Agreement (LOA). Our clinical metric results for 2D and 3D models are presented in Table 3 and 4 respectively. For both 2D and 3D models, the performance using cross-entropy and dice-loss functions for LV and RV is fairly similar, however the difference in performance is significant for MYO where dice-loss outperforms cross-entropy optimization. While distance metric scores are fairly similar for both 2D and 3D model, in clinical parameters we observe that the 3D model outperforms the 2D model.

Although accuracy scores are important when making clinical decisions, run-time efficiency and memory usage of the algorithm are also crucial to apply it in real-world applications. Our 2D model takes 2.9 hours to train using 4GB of GPU memory. At testing, it takes 0.3s and 1.2GB GPU memory to generate a single output(whole phase). Whereas our 3D model requires 2.6 hours for training and 4GB of GPU memory. At test time, it can generate output within 0.3s using 2GB GPU memory. This shows that our 2D and 3D models are efficient to train, are light-weight and relatively easy to deploy in clinical settings.

Table 1. 2D Segmentation: Distance Metric Results

	Dice Score						Hausdorff Distance					
	LV		RV		MYO		LV		RV		MYO	
	ED	ES	ED	ES	ED	ES	ED	ES	ED	ES	ED	ES
CE Loss	0.95	0.90	0.87	0.76	0.79	0.82	13.92	17.67	27.40	27.73	23.81	22.11
Dice Loss	0.95	0.90	0.90	0.79	0.86	0.88	9.51	12.29	16.1	20.38	13.45	14.88
Dice-CE Loss	0.95	0.90	0.89	0.81	0.83	0.84	9.15	11.7	16.0	18.22	13.87	15.35

Table 2. 3D Segmentation: Distance Metric Results

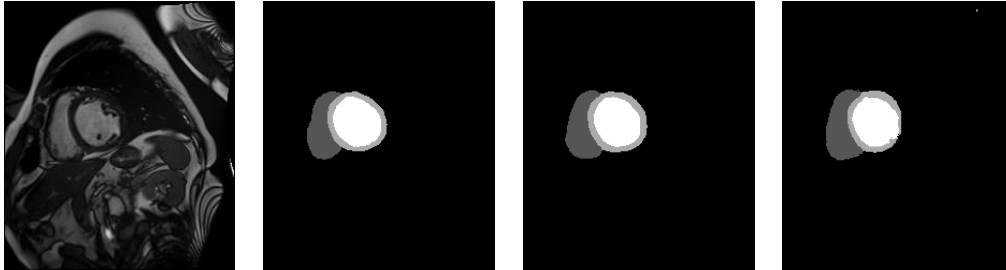
	Dice Score						Hausdorff Distance					
	LV		RV		MYO		LV		RV		MYO	
	ED	ES	ED	ES	ED	ES	ED	ES	ED	ES	ED	ES
CE Loss	0.94	0.89	0.86	0.73	0.76	0.81	12.36	14.41	25.85	29.57	43.47	43.82
Dice Loss	0.95	0.90	0.91	0.83	0.85	0.86	14.95	14.35	23.15	22.14	37.75	38.50
Dice-CE Loss	0.94	0.89	0.91	0.81	0.83	0.85	10.71	11.52	38.01	32.26	43.28	44.98

Table 3. 2D Segmentation: Clinical Metric Results

	Ejection Fraction						Myocardial Mass		
	LV			RV			MYO		
	CC	Bias	LOA	CC	Bias	LOA	CC	Bias	LOA
CE Loss	0.95	-0.74	-12.48,11.00	0.822	9.79	-10.89,30.47	0.93	-43.85	-92.12,4.42
Dice Loss	0.88	1.06	-18.10,20.22	0.822	9.35	-11.67,30.37	0.95	-6.32	-39.46,26.82
Dice-CE Loss	0.93	-0.46	-15.67,14.75	0.813	5.66	-16.59,27.91	0.94	-29.31	-68.20,9.58

Table 4. 3D Segmentation: Clinical Metric Results

	Ejection Fraction						Myocardial Mass		
	LV			RV			MYO		
	CC	Bias	LOA	CC	Bias	LOA	CC	Bias	LOA
CE Loss	0.975	1.04	-7.67,9.75	0.756	9.62	-15.11,34.35	0.922	-48.17	-97.51,1.17
Dice Loss	0.956	1.51	-10.09, 13.11	0.825	4.99	-17.17,27.15	0.958	3.77	-24.91,32.45
Dice-CE Loss	0.956	1.25	-10.37,12.87	0.867	6.19	-12.09,24.47	0.950	-10.08	-42.85,22.69

**Fig. 3.** Segmentation Results for 2D and 3D model. **From Left to Right:** Raw MR input image slice, Corresponding ground truth annotation, output predictions from 2D segmentation model and output predictions from 3D segmentation model.

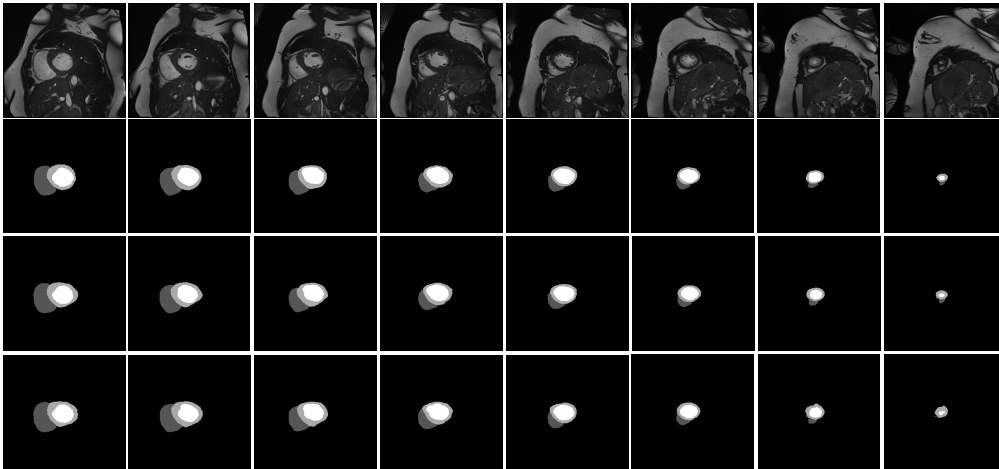


Fig. 4. Segmentation Results for a full MR image (Complete Phase) from slices 0-8 . **First Row:** Raw MR input images **Second Row:** Corresponding ground truth annotations. **Third Row:** Output predictions from 2D segmentation model. **Fourth Row:** Output predictions from 3D segmentation model.

4 Discussion

Model Structures Since 2D segmentation model is trained by splitting MR volumes into slices, it lacks the spatial context across the 3D volume. This reduces the model’s performance for slices at end of the phase, where the ratio of heart structure to background pixels is less. To solve this, we design our 2D model to accept a stack of image slices as input channels where the output is predicted for the middle slice. Given this design, we have 3 input options to pass as inputs: 1, 3 and 5. Among the three, we obtain the best performance with 3-input slices and report the scores for it in our results section.

Analyzing the 3D CNN segmentation model, we observe that the 3D model doesn’t meet the expectations in performance improvement over 2D, given its ability to exploit 3D structure from input volumes. We see that Dice Score for RV in 3D is better 2D model given the fact that it has complex shape and intensity inhomogeneities. Thus, predicting RV using single slice is much more difficult compared to looking at complete 3D context. Further improvements in performance can be brought about using post-processing techniques. Due to the modularity of 3D model architecture, we were able to quickly explore several designs and concepts. For example, we tried the recently introduced subpixel CNN’s [15] that proposes using subpixel layers as opposed to transposed convolutions. We executed this by replacing the *deconv* blocks in the expanding paths with a *subpixel* block that comprises of a *subpixel* layer between two *conv_bn_relu* blocks. However, no performance improvements were to be observed.

Data Augmentation At the initial stages of model design and training, we applied variety of data augmentation techniques demonstrated in [8,7], to incorporate randomness and robustness in the training the network. These include elastic deformations, random intensity jitter and affine transformations that include rotation, shearing, translation and flipping. Acquiring sub-par performances at testing, we found that applying heavy augmentations might misrepresent the anatomical structure of the heart. Instead, we opt to apply light data augmentations consisting of random rotations and scaling that can naturally match the variability of taking MR scans in real-world settings. With this

modification, we observed a jump in performance scores for both 2D and 3D segmentation model.

Optimization Functions Using our Dice loss as objective function improves the performance significantly for 2D and 3D models. As compared to pixel-level error optimization, dice loss is more robust and better at capturing the spatial context over the entire image. As shown in Table 1 and 2, the best dice scores are achieved by using dice loss function.

5 Conclusion

This paper introduces a 2D and 3D convolutional neural network for fully-automated cardiac MR segmentation. Our models have light-weight modular architecture, easy implementation and run-time efficiency. Using multiple loss criterion, we compare and analyze the performance of 2D and 3D model pipelines and show that both our models achieve near state-of-the-art accuracy scores in terms of distance metrics. With convincing performance in clinical accuracy metrics, we also prove our model's viability in real-world practical applications. Through our discussions, we derive several insights that can be used for optimizing overall performance of these segmentation models. For future work, we plan to utilize our segmentation models to learn and classify different cardiac diseases.

References

1. Girshick, R. (2015). Fast r-cnn. In Proceedings of the IEEE international conference on computer vision (pp. 1440-1448).
2. He K, Zhang X, Ren S, Sun J. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770-778 (2016).
3. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: Proc. CVPR. pp. 3431-3440 (2015).
4. P. V. Tran. A fully convolutional neural network for cardiac segmentation in short-axis mri. arXiv preprint arXiv:1604.00494, (2016).
5. Lieman-Sifry J, Le M, L au F, Sall S, Golden D. FastVentricle: Cardiac Segmentation with ENet. arXiv preprint arXiv:1704.04296. 2017 Apr 13.
6. O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In in Proc. of MICCAI, pages 234-241, (2015).
7. Cicek O, Abdulkadir A, Lienkamp SS, Brox T, Ronneberger O. 3D U-Net: learning dense volumetric segmentation from sparse annotation. In International Conference on Medical Image Computing and Computer-Assisted Intervention (pp. 424-432) 2016.
8. Milletari, F., Navab, N., Ahmadi, S.A.: V-Net: Fully convolutional neural networks for volumetric medical image segmentation. In: Proc. 3DV'16. pp. 565-571, (2016).
9. Automatic Cardiac Detection Challenge 2017. <http://www.creatis.insa-lyon.fr/Challenge/acdc/>
10. Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. CoRR abs/1502.03167 (2015)
11. Srivastava N, Hinton GE, Krizhevsky A, Sutskever I, Salakhutdinov R. Dropout: a simple way to prevent neural networks from overfitting. Journal of Machine Learning Research. 1;15(1):1929-58 (2014).
12. Pizer, Stephen M., et al. "Adaptive histogram equalization and its variations." Computer vision, graphics, and image processing 39.3 (1987): 355-368.
13. PyTorch. <http://pytorch.org/>
14. Zotti, C., Luo, Z., Lalande, A., Humbert, O., Jodoin, P. M. Novel Deep Convolution Neural Network Applied to MRI Cardiac Segmentation. arXiv preprint arXiv:1705.08943 (2017).
15. Dong C, Loy CC, He K, Tang X. Image super-resolution using deep convolutional networks. IEEE transactions on pattern analysis and machine intelligence. 1;38(2):295-307 (2016).