

Style Transfer and Self-Supervised Learning Powered Myocardium Infarction Super-Resolution Segmentation

1st Lichao Wang
*Department of Computing
Imperial College London*
l.wang22@imperial.ac.uk

2nd Jiahao Huang
*National Heart and Lung Institute
Imperial College London*
j.huang21@imperial.ac.uk

3rd Xiaodan Xing
*National Heart and Lung Institute
Imperial College London*
London, UK
x.xing@imperial.ac.uk

4rd Yinzhe Wu
*National Heart and Lung Institute
Imperial College London*
London, UK
yinzhe.wu18@imperial.ac.uk

5rd Ramyah Rajakulasingam
*National Heart and Lung Institute
Imperial College London*
London, UK
ramyah.rajakulasingam05@imperial.ac.uk

6rd Andrew D. Scott
*National Heart and Lung Institute
Imperial College London*
London, UK
a.scott07@imperial.ac.uk

7rd Pedro F Ferreira
*National Heart and Lung Institute
Imperial College London*
London, UK
p.f.ferreira05@imperial.ac.uk

8rd Ranil De Silva
*National Heart and Lung Institute
Imperial College London*
London, UK
r.desilva@imperial.ac.uk

9rd Sonia Nielles-Vallespin
*National Heart and Lung Institute
Imperial College London*
London, UK
s.nielles-vallespin@imperial.ac.uk

10rd Guang Yang
*National Heart and Lung Institute
Imperial College London*
London, UK
g.yang@imperial.ac.uk

Abstract—This study proposes a pipeline that incorporates a novel style transfer model and a simultaneous super-resolution and segmentation model. The proposed pipeline aims to enhance diffusion tensor imaging (DTI) images by translating them into the late gadolinium enhancement (LGE) domain, which offers a larger amount of data with high-resolution and distinct highlighting of myocardium infarction (MI) areas. Subsequently, the segmentation task is performed on the LGE style image. An end-to-end super-resolution segmentation model is introduced to generate high-resolution mask from low-resolution LGE style DTI image. Further, to enhance the performance of the model, a multi-task self-supervised learning strategy is employed to pre-train the super-resolution segmentation model, allowing it to acquire more representative knowledge and improve its segmentation performance after fine-tuning. https://github.com/wlc2424762917/Med_Img

Index Terms—Diffusion tensor imaging, late gadolinium enhancement, myocardium infarction segmentation, style transfer, self-supervised learning

I. INTRODUCTION

Diffusion tensor (DT) cardiovascular magnetic resonance (CMR) is a novel noninvasive tool that enables inference of sheetlet orientations, which are altered under pathological

conditions [1]. Preliminary studies have demonstrated the potential of DT CMR to detect microstructural abnormalities in myocardium infarction (MI) [2], suggesting the feasibility of MI segmentation on diffusion tensor imaging (DTI) images.

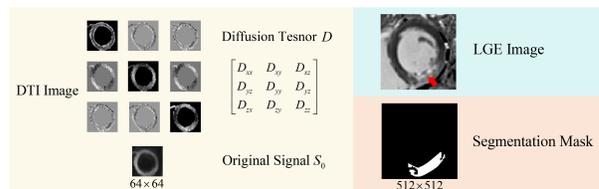


Fig. 1. Visualized data example of diffusion tensor imaging (DTI) and Late gadolinium enhancement (LGE) images. For DTI image, different components including DT D and the corresponding original signal S_0 are presented.

However, the domain of DTI has not been previously explored for MI segmentation. Current researches [3]–[6] were predominantly conducted using late gadolinium enhancement (LGE) images, due to certain limitations of DTI image compared to LGE image, including the lack of direct MI indication, a comparatively lower resolution, and a dearth of labeled data, as demonstrated in Fig. 1.

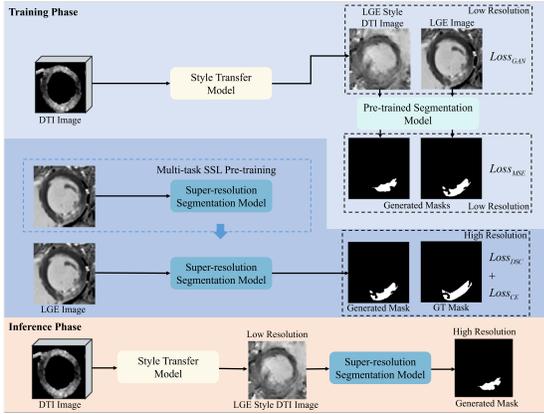


Fig. 2. The proposed pipeline. During the training phase, the style transfer model undergoes unsupervised training, while the super-resolution segmentation model is pre-trained using self-supervised pre-training and supervised fine-tuning strategy. In the inference phase, the Segmentor is excluded, and the pipeline is integrated. $Loss_{GAN}$, $Loss_{MSE}$, $Loss_{DSC}$, $Loss_{CE}$ stand for the CycleGAN loss [7], mean square error loss, dice loss, and cross-entropy loss respectively.

To this end, we propose a novel pipeline to leverage the advantages of LGE data for MI segmentation on DTI image. The pipeline consists of a style transfer model and a super-resolution segmentation (SSeg) model, as depicted in Fig.2. The style transfer model first converts the DTI image into LGE style, where MI areas are highlighted, and more labeled data is available. Then the SSeg model directly generates the high-resolution segmentation mask, facilitating improved segmentation of small foreground regions.

Our style transfer model is based on the CycleGAN [7]. To preserve the underlying segmentation mask throughout the style transfer, we develop the CycleGAN with a segmentation sub-network model, namely CycleGANSeg, aiming at keeping the integrity of the underlying MI segmentation mask.

Inspired by the Dual Super-Resolution Learning framework [8], we propose our SwinTransformer [9] based end-to-end SSeg model, namely SwinSSegNet. The SwinSSegNet aims to leverage more detailed information in the training process, and directly generates high-resolution segmentation masks in the inference process.

Moreover, despite the relatively higher availability of LGE images, the quantity remains limited. To address this constraint, we construct a hybrid dataset by incorporating LGE images obtained from publicly available datasets, ACDC [10] and LiVScar [11]. Based on the hybrid dataset, we adopt a multitask self-supervised learning (SSL) pre-training strategy, including contrastive learning [12], masked image modeling [13], and rotation prediction [14]. This approach enables our model to acquire a broad range of representation knowledge of LGE images, thereby enhancing its performance after fine-tuning. In summary, our contributions are as following:

1. We introduce a novel pipeline powered by style transfer and SSL for MI super-resolution segmentation on DTI image. This pipeline effectively harnesses the abundant LGE data, and enhance the accuracy of MI segmentation on DTI image.
2. We present the CycleGANSeg, which incorporates a segmentation sub-network within the CycleGAN framework.

This model is capable of converting DTI image to the LGE style, while preserving the integrity of the MI region. Notably, the CycleGANSeg can be trained in an unsupervised manner.

3. We propose the SwinSSegNet, an end-to-end super-resolution segmentation model that can directly generate high-resolution segmentation mask. This end-to-end paradigm has superior performance comparing to utilizing the 2-stage (first segmentation, then up-sample) paradigm.

4. We design a 2D multi-task self-supervised learning pre-training strategy on our curated hybrid LGE dataset to further enhance the performance of the SwinSSegNet model.

II. METHODS

As shown in Fig 2, our proposed pipeline contains two models, a style transfer model, i.e., CycleGANSeg, and an SSeg model, i.e., SwinSSegNet. The CycleGANSeg initially translates the DTI image into the LGE style, with a resolution of 64×64 . Subsequently, the SwinSSegNet generates the MI segmentation mask with an upsampled resolution of 512×512 .

A. CycleGAN with Segmentor

As shown in Fig 3, the style transfer model is based on the CycleGAN. To keep the underlying segmentation mask unchanged in the style transfer process, a pre-trained segmentation network, namely Segmentor, is incorporated, along with two specific loss functions.

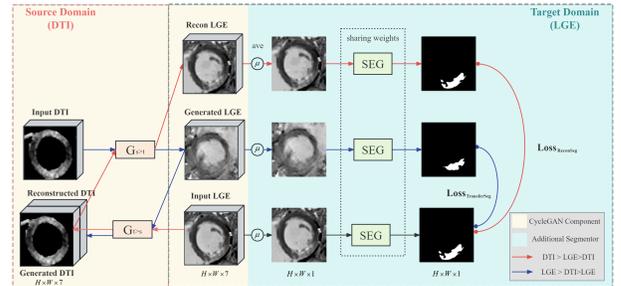


Fig. 3. The dataflow of CycleGANSeg. $G_{S>T}$, $G_{T>S}$, and SEG indicate the generator from DTI image to LGE image, the generator from LGE image to DTI image, and the Segmentor respectively. The black arrows, red arrows, and blue arrows depict the flow of data within the LGE domain, the bidirectional flow from LGE to DTI and back to LGE, and the bidirectional flow from DTI to LGE and back to DTI, respectively.

For the input and output data settings, we use the paired DTI images and LGE images as input, to ensure they have the same MI segmentation mask. Since our proposed CycleGANSeg comprises three sub-networks (a generator, a discriminator, and the Segmentor), the computational cost escalates significantly, as the resolution increases. Hence, to maintain efficiency during training, we standardize the resolution of all inputs and outputs to 64×64 .

The Segmentor outputs the segmentation mask of the original LGE image, the generated LGE image, and the reconstructed LGE image. To keep the integrity of the underlying segmentation mask, two mean square error (MSE) losses are adopted, as shown in (1) and (2). This approach only requires paired images from different domains, therefore the

CycleGANSeg offers the advantage of unsupervised training, obviating the need for extensive annotation.

$$Loss_{TransferSeg} = (Mask_{trans} - Mask_{ori})^2, \quad (1)$$

$$Loss_{ReconSeg} = (Mask_{recon} - Mask_{ori})^2, \quad (2)$$

where $Mask_{trans}$, $Mask_{recon}$, and $Mask_{ori}$ stand for the mask of the generated LGE style DTI image, the reconstructed LGE image, and the original LGE image respectively.

B. SwinTransformer Super-resolution Segmentation Model

The SwinSSegNet is designed to directly generate the high-resolution segmentation mask, by utilizing SwinTransformer as the encoder, employing a series of convolution-based blocks as the decoder. It is designed as a flexible framework, in which the upsample scale is adjustable and different backbones can be plugged in as the encoder. We set the downsampled LGE image (64×64) as the training input, and the original segmentation mask (512×512) as output. To keep more detailed information, the patch size is set as 2.

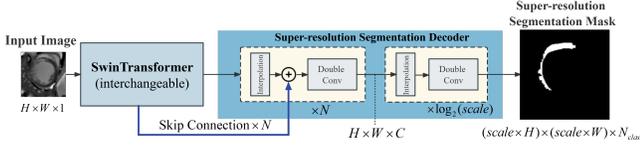


Fig. 4. The architecture of the SwinSSegNet. $scale$ stands for the upsampling scale. $H \times W$ and C represents the original spatial shape and the channel number of the feature map respectively. The decoder of the SwinSSegNet can be divided into two parts. The first part focuses on upsampling the feature map to match the spatial dimensions of the input image. The subsequent part is responsible for generating the upsampled output segmentation mask.

As shown in Fig. 4, within each layer of the first part, skip connection enables the transmission of detailed information from the encoder to the decoder. The second part has the flexibility to adjust the upsample scale. Each decoder block comprises an interpolation upsampling layer followed by a double convolution layer. The double convolution layer consists of two stacks, encompassing convolution, Rectified Linear Unit, and Batch Normalization layers.

C. Universal Multi-task Self-supervised Learning Pre-training

Motivated by the recent work on 3D patch-wise SSL pre-training framework for Swin UNETR [15], we adopt a 2D SSL pre-training strategy to fully harness the capabilities of the SwinTransformer encoder. This strategy encompasses various components, including contrastive learning, which aims to encourage the model to capture general semantic features; masked image modeling, which assists the model in learning detailed features; and rotation prediction, which promotes the acquisition of spatial features. Additionally, we collect and crop LGE images from the ACDC and LiVScar datasets, incorporating them alongside our private dataset to form a hybrid dataset. Subsequently, the SSL pre-training is performed on this hybrid dataset, leading to the formulation of a novel strategy termed universal multi-task self-supervised learning (U-SSL) pretraining strategy.

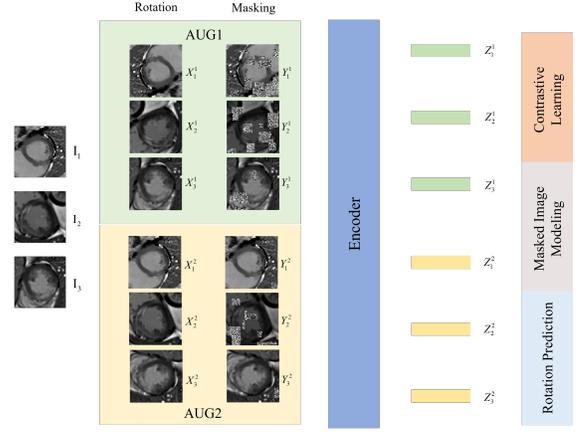


Fig. 5. The 2D multi-task self-supervised learning framework. Input late gadolinium enhancement images are augmented with rotation and random masking, subsequently fed to the encoder as input.

In the implementation, the same data augmentation method is utilized to generate similar/dissimilar pairs for contrastive learning, masked images for masked image modeling, and rotated images for the rotation prediction task. As shown in Fig. 5, the SSL data engine contains two sets of augmentation operations, AUG1 and AUG2. Each set encompasses rotation and masking. To exemplify, within AUG1, a batch of images, (I_1, I_2, \dots, I_n) , initially undergo rotation to yield $(X_1^1, X_2^1, \dots, X_n^1)$, and then are randomly masked out to produce $(Y_1^1, Y_2^1, \dots, Y_n^1)$. The rotation prediction task predicts \hat{y}_r , the rotation angle of (I_i, X_i) , and the associated loss, $Loss_{ROT}$, is designed as (3), where ground truth $y_r \in 0^\circ, 90^\circ, 180^\circ, 270^\circ$. The masked image modeling task aims at generating I_i from Z_i , and the loss, $Loss_{MIM}$, is designed as (4). Contrastive learning task maximizes the dot product similarity (sim) between positive embedding pairs (Z_i^1, Z_i^2) , while minimizing that between the other negative embedding pairs, and the loss, $Loss_{CL}$, is designed as (5).

$$Loss_{ROT} = \sum_{r=0}^3 y_r \log(\hat{y}_r), \quad (3)$$

$$Loss_{MIM} = \|X_i - I_i\|_1, \quad (4)$$

$$Loss_{CL} = -\log \frac{\exp(\sum_i^{2N} sim(Z_i^1, Z_i^2)/t)}{\sum_i^{2N} \sum_k^{2N} \mathbf{1}_{k \neq i} \exp(sim(Z_i^1, Z_k^2)/t)}, \quad (5)$$

where t is the measurement of normalized temperature scale. $\mathbf{1}_{k \neq i}$ is the indicator function evaluating to 1 if $k \neq i$.

III. EXPERIMENTS

A. Dataset

Our private dataset consists of 277 unlabeled DTI images and 271 labeled LGE images, with the resolution of 64×64 and 512×512 respectively. Among these, 77 DTI-LGE image pairs share identical segmentation masks. We denote the 271 labeled LGE images as the LGE sub-dataset, and the 77 DTI-LGE image pairs with same segmentation masks as paired DTI-LGE image sub-dataset. Both the LGE image and DTI image have the target subject, MI. Manual segmentation of

MI was undertaken by a CMR physicist with over 3 years of experience; these segmentations serve as the ground truth for the training and assessment of our proposed SwinSSegNet and pipeline. Additionally, We incorporate images from the ACDC and LiVScar datasets into our LGE subset to construct a hybrid LGE dataset for U-SSL pre-training. The ACDC dataset comprises 100 cine MRI scans (1,902 slices), and the LiVScar contains 30 images (200 slices). From each slice, the region of interest is cropped and upsampled to the resolution of 512×512 for our U-SSL pre-training process.

B. Implementation and Evaluation Details

We conducted our experiments on an NVIDIA RTX3090 GPU with 24GB GPU RAM. The masking rate was set as 45% for the SSL pretraining strategy. The AdamW optimizer was used. The CycleGANSeg was trained on our paired DTI-LGE image sub-dataset. The SwinSSegNet was pre-trained on the hybrid LGE dataset and fine-tuned on our LGE sub-dataset. All training procedures adopted the batch size of 24.

With respect to the SwinSSegNet, the Dice similarity coefficient [16] (DSC), indicative of the congruity between the segmentation outcome and the ground truth segmentation mask of the LGE image, was employed as the evaluation metric. In the context of the pipeline, the DSC between the segmentation result of the LGE style DTI image and the ground truth segmentation mask was utilized as the evaluation metric.

C. Results of the SwinSSegNet

We compared our proposed SwinSSegNet with state-of-the-art segmentation baselines (SwinUNet [17], UNet-2022 [18], and TransUNet [19]) on the LGE sub-dataset (Fig. 6.(A)). To ensure fairness, all models were trained from scratch. Compared with other models, the proposed SwinSSegNet achieves the best performance.

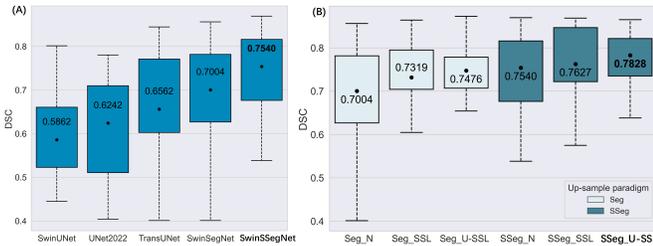


Fig. 6. (A). Comparison of the SwinSSegNet with baselines. SwinSegNet indicates adopting the upsample scale as 1 for the SwinSSegNet, and using interpolation to upsample the segmentation mask. (B). Comparison of the up-sample methods and the different pre-training strategies. SSeg stands for the end-to-end paradigm and Seg stands for the 2-stage paradigm (segmentation followed by bilinear interpolation upsampling). “_N”, “_SSL”, and “_U-SSL” represent trained from scratch, trained with multi-task self-supervised pre-training strategy, and trained with universal multi-task self-supervised pre-training strategy, respectively.

Furthermore, a series of ablation studies were conducted to investigate various aspects of the proposed methodology. The performance of different up-sample paradigms and the performance of different pre-training strategies are compared in Fig. 6.(B). The result indicates that the SSeg models are superior to the Seg models. Moreover, the employment of the SSL

leads to enhanced model performance, and the incorporation of the public data can further boost the enhancement. A noteworthy observation is the decreasing standard deviation across model trained from scratch, model pre-trained with SSL, and model pre-trained with U-SSL. This observation highlights the substantial impact of contrastive learning, which facilitates the acquisition of comprehensive representation knowledge. Consequently, the model becomes adept at effectively handling segmentation tasks, even for challenging samples that initially exhibited subpar performance. Moreover, we compared the effect of choosing different backbone models as the encoder, the result is shown in Table. I.

Encoder	ResNet [20]	ConvNext [21]	Biformer [22]	SwinTransformer
DSC	0.718±0.0094	0.732±0.0093	0.745±0.0093	0.755±0.0090

TABLE I

COMPARISON OF USING DIFFERENT BACKBONE MODEL AS ENCODER.

In the encoder comparison, we adopted the patch size as 2 for Biformer and SwinTransformer, and the downsample size as 2 for the first pooling layer in ResNet and ConvNext. The result indicates that SwinTransformer stands out as the most effective encoder.

D. Results of the Pipeline

To illustrate the superiority of our proposed pipeline (CycleGANSeg+SSeg), i.e., introducing the CycleGANSeg to convert DTI image to LGE style before using the SwinSSegNet trained with U-SSL pre-training strategy to perform segmentation, we compared the cases of not using style transfer model (DTI SSeg) and using the original CycleGAN as the style transfer model (CycleGAN+SSeg).

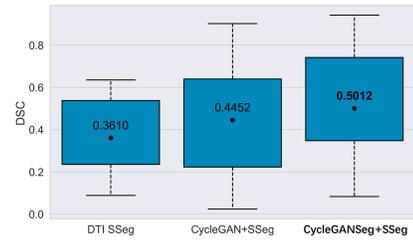


Fig. 7. Comparison of different pipelines. DTI SSeg stands for the SwinSSegNet trained from scratch on DTI image data. CycleGAN+SSeg stands for using the original CycleGAN as the style transfer model, and the SwinSSegNet trained with U-SSL pre-training strategy on LGE data as the SSeg model.

The quantitative result shown in Fig. 7 and the visualized results Fig. 8 indicate that our pipeline effectively highlights the MI area and achieves the best MI segmentation performance. The DTI SSeg model demonstrates poor performance, while CycleGAN+SSeg improves performance but introduces increased variance. This aligns with expectations, as the original CycleGAN lacks a dedicated mechanism to preserve the underlying segmentation mask. In contrast, the proposed CycleGANSeg+SSeg enhances the average performance while mitigating the range of both upper and lower bounds. This is achieved through the integration of the Segmentor in training process, which effectively helps the model to preserve the original segmentation mask, thereby stabilizing the style transfer process.

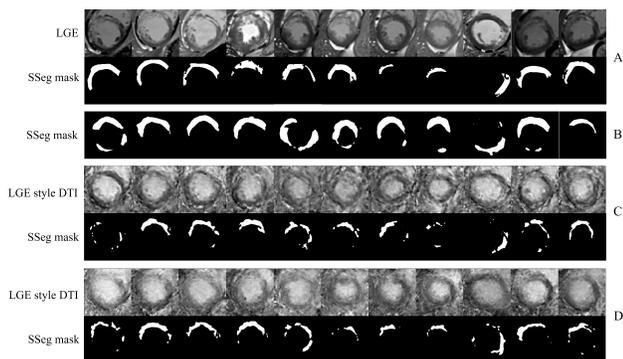


Fig. 8. Visualized results of the different pipelines. Row A illustrates the ground truth LGE image and SSeg mask. Row B illustrates the SSeg results of DTI SSeg model. Row C and Row D illustrate the LGE style DTI image, SSeg results of CycleGAN+SSeg pipeline and CycleGANSeg+SSeg pipeline respectively.

IV. CONCLUSIONS

We present a novel pipeline for MI super-resolution segmentation on DTI image, incorporating the CycleGANSeg and the SwinSSegNet. The CycleGANSeg transforms DTI image to LGE style, bridges domain gaps, and simplifies MI segmentation. The SwinSSegNet surpasses the two-step segmentation paradigm in generating high-resolution segmentation mask, and the integration of the U-SSL pre-training strategy further enhances the segmentation performance.

REFERENCES

- [1] I. Teh, D. McClymont, M.-C. Zdora, H. J. Whittington, V. Davidoiu, J. Lee, C. A. Lygate, C. Rau, I. Zanette, and J. E. Schneider, "Validation of diffusion tensor MRI measurements of cardiac microstructure with structure tensor synchrotron radiation imaging," *Journal of Cardiovascular Magnetic Resonance*, vol. 19, pp. 1–14, 2017.
- [2] Z. Khalique, P. F. Ferreira, A. D. Scott, S. NIELLES-Vallespin, D. N. Firmin, and D. J. Pennell, "Diffusion tensor cardiovascular magnetic resonance imaging: a clinical perspective," *Cardiovascular Imaging*, vol. 13, no. 5, pp. 1235–1255, 2020.
- [3] Z. Chen, A. Lalande, M. Salomon, T. Decourselle, T. Pommier, A. Qayyum, J. Shi, G. Perrot, and R. Couturier, "Automatic deep learning-based myocardial infarction segmentation from delayed enhancement MRI," *Computerized Medical Imaging and Graphics*, vol. 95, p. 102014, 2022.
- [4] C. Xu, Y. Wang, D. Zhang, L. Han, Y. Zhang, J. Chen, and S. Li, "BMAnet: Boundary Mining With Adversarial Learning for Semi-Supervised 2D Myocardial Infarction Segmentation," *IEEE Journal of Biomedical and Health Informatics*, vol. 27, no. 1, pp. 87–96, 2022.
- [5] S. Yang and X. Wang, "A Hybrid Network for Automatic Myocardial Infarction Segmentation in Delayed Enhancement-MRI," in *Statistical Atlases and Computational Models of the Heart. M&Ms and EMIDEC Challenges: 11th International Workshop, STACOM 2020, Held in Conjunction with MICCAI 2020, Lima, Peru, October 4, 2020, Revised Selected Papers 11*. Springer, 2021, pp. 351–358.
- [6] J. Wang, H. Huang, C. Chen, W. Ma, Y. Huang, and X. Ding, "Multi-sequence Cardiac MR Segmentation with Adversarial Domain Adaptation Network," in *Statistical Atlases and Computational Models of the Heart. Multi-Sequence CMR Segmentation, CRT-EPiggy and LV Full Quantification Challenges: 10th International Workshop, STACOM 2019, Held in Conjunction with MICCAI 2019, Shenzhen, China, October 13, 2019, Revised Selected Papers 10*. Springer, 2020, pp. 254–262.
- [7] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2223–2232.
- [8] L. Xilinx Wang, D. Li, Y. Zhu, L. Tian, and Y. Shan, "Dual Super-Resolution Learning for Semantic Segmentation," *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3773–3782, 2020.
- [9] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin Transformer: Hierarchical Vision Transformer using Shifted Windows," *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 9992–10002, 2021.
- [10] O. Bernard, A. Lalande, C. Zotti, F. Cervenansky, X. Yang, P.-A. Heng, I. Cetin, K. Lekadir, O. Camara, M. A. G. Ballester *et al.*, "Deep Learning Techniques for Automatic MRI Cardiac Multi-Structures Segmentation and Diagnosis: Is the Problem Solved?" *IEEE transactions on medical imaging*, vol. 37, no. 11, pp. 2514–2525, 2018.
- [11] T. Mansi, K. Mcleod, M. Pop, K. S. Rhode, M. Sermesant, and A. A. Young, "Statistical Atlases and Computational Models of the Heart. Imaging and Modelling Challenges," in *Lecture Notes in Computer Science*, vol. 8896, 2015, p. 296.
- [12] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A Simple Framework for Contrastive Learning of Visual Representations," in *International conference on machine learning*. PMLR, 2020, pp. 1597–1607.
- [13] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick, "Masked Autoencoders Are Scalable Vision Learners," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 16000–16009.
- [14] S. Gidaris, P. Singh, and N. Komodakis, "Unsupervised Representation Learning by Predicting Image Rotations," *ArXiv*, vol. abs/1803.07728, 2018.
- [15] Y. Tang, D. Yang, W. Li, H. R. Roth, B. A. Landman, D. Xu, V. Nath, and A. Hatamizadeh, "Self-Supervised Pre-Training of Swin Transformers for 3D Medical Image Analysis," *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 20698–20708, 2021.
- [16] L. R. Dice, "Measures of the Amount of Ecologic Association Between Species," *Ecology*, vol. 26, no. 3, pp. 297–302, 1945.
- [17] H. Cao, Y. Wang, J. Chen, D. Jiang, X. Zhang, Q. Tian, and M. Wang, "Swin-Unet: Unet-like Pure Transformer for Medical Image Segmentation," in *ECCV Workshops*, 2021.

- [18] J. Guo, H.-Y. Zhou, L. Wang, and Y. Yu, "UNet-2022: Exploring Dynamics in Non-isomorphic Architecture," *arXiv preprint arXiv:2210.15566*, 2022.
- [19] J. Chen, Y. Lu, Q. Yu, X. Luo, E. Adeli, Y. Wang, L. Lu, A. L. Yuille, and Y. Zhou, "TransUNet: Transformers Make Strong Encoders for Medical Image Segmentation," *ArXiv*, vol. abs/2102.04306, 2021.
- [20] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2015.
- [21] Z. Liu, H. Mao, C. Wu, C. Feichtenhofer, T. Darrell, and S. Xie, "A ConvNet for the 2020s," *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 11 966–11 976, 2022.
- [22] L. Zhu, X. Wang, Z. Ke, W. Zhang, and R. Lau, "BiFormer: Vision Transformer with Bi-Level Routing Attention," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.