

Tempera: Spatial Transformer Feature Pyramid Network for Cardiac MRI Segmentation

Christoforos Galazis^{1,2*}, Huiyi Wu², Zhuoyu Li³, Camille Petri², Anil A. Bharath⁴, and Marta Varela²

¹ Department of Computing, Imperial College London, UK

² National Heart & Lung Institute, Imperial College London, UK

³ Department of Metabolism, Digestion & Reproduction, Imperial College London, UK

⁴ Department of Bioengineering, Imperial College London, UK

*c.galazis20@imperial.ac.uk

Abstract. Assessing the structure and function of the right ventricle (RV) is important in the diagnosis of several cardiac pathologies. However, it remains more challenging to segment the RV than the left ventricle (LV). In this paper, we focus on segmenting the RV in both short (SA) and long-axis (LA) cardiac MR images simultaneously. For this task, we propose a new multi-input/output architecture, hybrid 2D/3D geometric spatial TransformEr Multi-Pass fEature pyRAmid (Tempera). Our feature pyramid extends current designs by allowing not only a multi-scale feature output but multi-scale SA and LA input images as well. Tempera transfers learned features between SA and LA images via layer weight sharing and incorporates a geometric target transformer to map the predicted SA segmentation to LA space. Our model achieves an average Dice score of 0.836 and 0.798 for the SA and LA, respectively, and 26.31mm and 31.19mm Hausdorff distances. This opens up the potential for the incorporation of RV segmentation models into clinical workflows.

Keywords: Cardiac MRI · Right Ventricle · Segmentation · Multi-View · 2D/3D Network · Spatial Transformer · Feature Pyramid.

1 Introduction

Cardiac Magnetic Resonance Imaging (MRI) is the most widely used imaging technique to quantify the structure and function of the heart [1]. As such, it can be used to assess the right ventricle (RV) and contribute to the diagnosis and monitoring of cardiac pathologies such as coronary heart disease, pulmonary hypertension, dysplasia and cardiomyopathies [5, 6].

Despite the RV’s importance, the left ventricle (LV) has traditionally been analysed in greater detail due to its pivotal role in a wider range of pathologies [17]. The RV is more challenging to accurately segment for both clinicians and (semi-)automated algorithms [3, 17]. This is due to the RV’s more complex crescent shape, thinner ventricular wall and heavier presence of wall trabeculations compared to the LV [8, 17].

In this paper, we focus on RV automatic segmentation as part of the ” *Multi-Disease, Multi-View & Multi-Center Right Ventricular Segmentation in Cardiac MRI* ” (M&Ms2) challenge [16]. We propose a novel hybrid 2D/3D deep neural network that takes both short-axis (SA) and long-axis (LA) images as inputs. It includes a novel multi-input/output feature pyramid that facilitates weight sharing across SA in-plane slices and between the SA and LA views. Additionally, we include a geometric target spatial transformer that utilizes the known spatial relationship between the different cardiac views.

2 Methods

Data The end-diastolic (ED) and end-systolic (ES) cardiac MRI image phases used for the experiments come from a 360-subject bSSFP CINE MRI dataset publicly available through the M&Ms-2 challenge [16]. The dataset is obtained across different sites, different vendors and has pathologies in the test set not included in the training. The available ground truth has manual segmentations of the LV and RV blood pool and the LV myocardium. From this set, 160 cases are dedicated for training. The remaining 40 and 160 cases are the development and test sets used for their respective phases in the competition. We further split the 160-case training set to 150 for model training and 10 for validation.

Preprocessing As the images were acquired across different centers and vendors, they need to be standardized before being passed to the model. The resolution of both SA and LA images was resampled to 1.25×1.25 mm in-plane using b-spline interpolation. The through-plane spacing was left unchanged at 10 mm. We first automatically identify a region of interest (ROI) containing the RV and LV by applying Canny edge detection on the ED and ES image difference, then using the circular Hough transform to identify the heart. This allows us to identify the regions in the image that have the largest movement across the cardiac cycle, which we assume to be the heart. The image size was standardized to $192 \times 192 \times 17$ through center ROI cropping. Finally, we standardize the cropped image such that it has a mean of zero and unit variance.

Additionally, we identify the affine transformation between the SA and LA images, which will subsequently be used in the segmentation network. The LA is treated as a 3D image (by adding a depth axis of size 1), thus allowing for a 3D/3D registration to take place. We pre-align the images based on the available file metadata of the images. For the registration, we follow a coarse-to-fine blurring approach to initially align with global features before finer ones. The registration is optimized to maximize the mutual information score. This is done using the SimpleITK package.

Model: Architecture To simultaneously segment the SA and LA images, we propose a new hybrid 2D/3D geometric spatial Transformer Multi-Pass Feature pyramid (Tempera) network, as shown in detail in *Figure 1*. Our architecture consists of 2D SA/LA hard-weight sharing layers, independent 3D (SA) and 2D (LA) branches and finally a geometric spatial transformer.

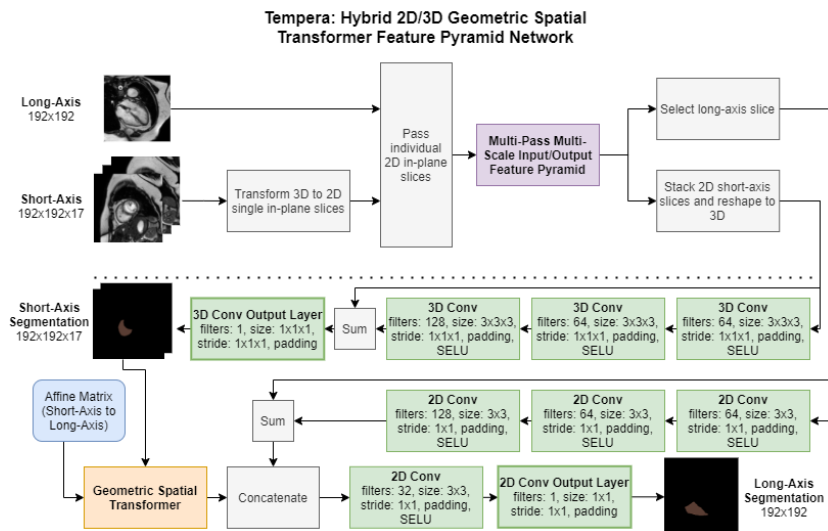


Fig. 1. An overview of the architecture of the hybrid 2D/3D Tempera network. It receives the individual 2D in-plane slices for both the SA and LA as inputs to the feature pyramid component. Then, the short-axis slices are used to reconstruct the 3D representation and each image view is passed to their respective 2D/3D convolutional branch. Finally, the predicted SA segmentation is transformed to the LA space to improve its prediction.

For the shared layers, we extend the Multi-scale Feature Pyramid network [14] to an architecture that we call Multi-Pass Feature Pyramid (MPFP). In the MPFP, the weight sharing, and thus internal (3x3) feature transfer, is achieved across the SA and LA in-plane slices and between the SA slices themselves. We are thus able to share features across all slices that relate to the RV shape, texture and contrast relative to the surroundings. The (partial) scale-invariant input/output of MPFP is designed to help Tempera generalize to unseen pathologies, such as dilated RV. An example is shown in *Figure 3* in which we base our assumption. The specific implementation details of our pyramid can be viewed in *Figure 2*.

Our second component is the Geometric target Spatial Transformer (GST), depicted in *Figure 4*. This is a differentiable non-trainable component, that we have built upon the spatial transformer from [2]. It applies transformations between different domain and target spaces, such as 3D to 2D and vice versa. The GST takes as input the pre-computed affine matrix to perform relevant coordinate transformations from the SA to LA. Specifically, we utilise it on the predicted SA segmentation to localize the RV in the LA.

Model: Optimizer We use the Adam optimizer [11] with: exponential decay rates for first moment estimates $\beta_1 = 0.9$, second moment estimates $\beta_2 = 0.999$ and $\epsilon = 10^{-8}$ to avoid division with 0. Additionally, we use an empirically-chosen initial learning rate of 5×10^{-4} , accompanied by a learning rate decay of 0.9 every 30 epochs.

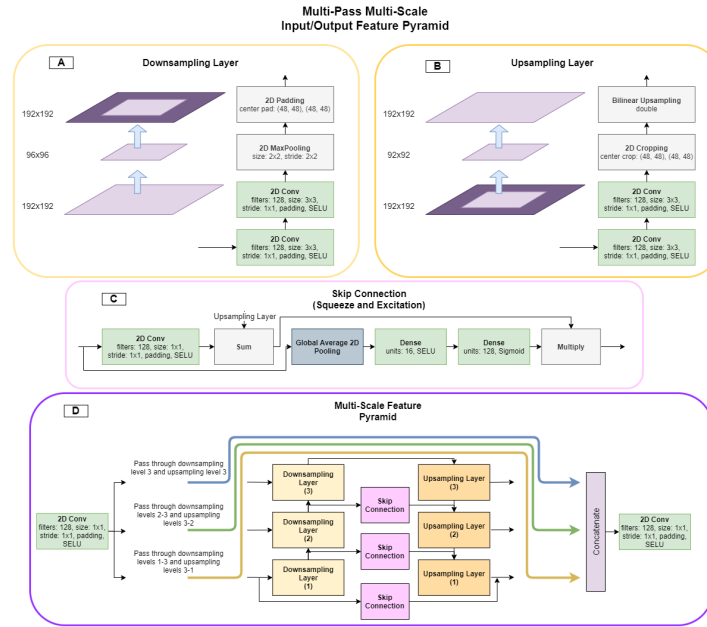


Fig. 2. Architecture of our proposed Multi-Pass and multi-scale input/output Feature Pyramid (MPFP). The downsampling component in block *A* consists of two 2D convolutional layers, followed by a pooling and padding layer. Similarly, the upsampling layer, block *B*, has two 2D convolutional layers, a cropping operation and finally bilinear upsampling. The skip connection in *C* contains a 2D convolutional layer followed by a summing operation and a squeeze-and-excitation component [10] to scale features. Finally, *D* illustrates how the components interact to form the MPFP. 2D convolutional layers are used to generate consistent feature sizes for the inputs and also to merge the outputs. After every layer we use a Scaled Exponential Linear Unit (SELU) [12] activation function, with exception the last layer of the skip connection uses a sigmoid. The input data to the pyramid is passed n times and generates n outputs, where $n = 3$ is the number of pyramid levels. Each subsequent pass goes through one less downsampling and upsampling layer.

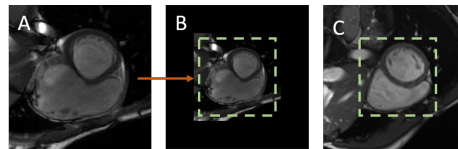


Fig. 3. Illustrative examples of how the feature pyramid may allow a hypothetical feature transfer to an unknown pathology for the network (dilated RV). Image *A* depicts SA basal slice of a patient with dilated RV. Image *B* is the downsampled version of image *A* by the MPFP in the first block. Finally, image *C* shows the SA basal slice of a healthy subject. The RV in images *B* and *C* are roughly of equal size. This allows feature transfer to untrained dilated RV cases.

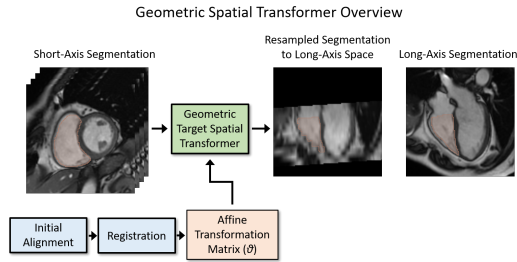


Fig. 4. The Geometric target Spatial Transformer (GST) affine transforms and resamples using linear interpolation the predicted segmentation from the SA to LA space, given an affine matrix as input. Their spatial relationship is stored in the image’s metadata, in which we further correct with an affine registration algorithm.

Model:Loss Function We minimized the combined [18] equal weighted Dice and focal loss (multiplicative weighted cross-entropy) [15] functions. We set the focal loss class weights $\alpha = 0.25$ and relaxation term $\gamma = 2.0$. Through our experimentation, we found that including the Dice loss improves the learning during initial epochs. The focal loss contributes during the later stages by dealing with class imbalance, in which background pixels outweigh foreground ones. Additionally, it helps with difficult to classify areas regardless if they are foreground or background. We set a λ of 75 for the SA loss and 1 for the LA. The equation is formed as bellow, where p is the probability of the ground truth label:

$$L_{Dice(p_t)} = 1 - \frac{2 \times true_positive_{p_t}}{2 \times true_positive_{p_t} + false_positive_{p_t} + false_negative_{p_t}} \quad (1)$$

$$L_{Focal}(p_t) = -a_t \times (1 - p_t)^\gamma \times \log(p_t) \quad (2)$$

$$L_{Total}(p_{SA}, p_{LA}) = \lambda_{SA} \times (L_{Dice}(p_{SA}) + L_{Focal}(p_{SA})) + \lambda_{LA} \times (L_{Dice}(p_{LA}) + L_{Focal}(p_{LA})) \quad (3)$$

Model:Hyperparameters We use LeCun weight initialization [12, 13] from a normal distribution and set the biases to zero. This is done to maintain the self-normalizing property when using Scaled Exponential Linear Units (SELUs) [12]. Also, we set the batch size to 1 due to memory limits and train for 300 epochs, which took nearly 16 hours to complete. The source code is available at https://github.com/cgalaz01/mnms2_challenge.

Model:Data Augmentation To prevent the model from overfitting and to expand the dataset, we perform data augmentations. These include: in-plane rotations, in-plane anisotropic scalings, blurring, Gaussian noise addition, mean intensity shift and removing square segments of the image (in-painting). We use randomly selected parameters at each data iteration for each epoch.

Model:Data Postprocessing We found that our model will occasionally identify surrounding tissues as the RV. To circumvent this, we added a post-processing step, in which only the largest connected region will be considered

as a valid prediction. Additionally, we apply a median filter to correct some of the interpolation errors caused by transforming back to the domain space using nearest neighbor interpolation.

3 Results

Tempera is able to accurately segment both the SA and LA views. On the test set we obtain 0.836 ± 0.23 Dice score and 26.31 ± 52.68 mm Hausdorff distance for the SA. The LA has a lower performance with 0.798 ± 0.28 Dice score and 31.19 ± 31.19 mm Hausdorff distance. From the predicted tests cases, there were a total of 11 SA and 20 LA that most likely failed during the ROI detection process, thus returning empty segmentations. Full details of the results are shown in *Table 1*.

A representative segmentation of the SA and LA images can be seen in *Figure 5*. Tempera’s segmentations closely match the ground truth labels, except in the LA basal slices, where Tempera occasionally under-segments the RV.

Dataset	Dice Score		Hausdorff Distance (mm)	
	SA	LA	SA	LA
Validation	$0.895 \pm \text{n/a}$	$0.829 \pm \text{n/a}$	$11.8 \pm \text{n/a}$	$15.118 \pm \text{n/a}$
Testing	0.836 ± 0.23	0.798 ± 0.28	26.31 ± 52.68	31.19 ± 71.29
Testing (no failures)	0.896 ± 0.05	0.899 ± 0.06	18.384 ± 24.99	14.186 ± 26.86

Table 1. Summary of Tempera’s performance on the (competition) validation and test sets using Dice score and Hausdorff distance. The test set’s evaluation was repeated (no failures), in which we removed cases where the preprocessing step failed.

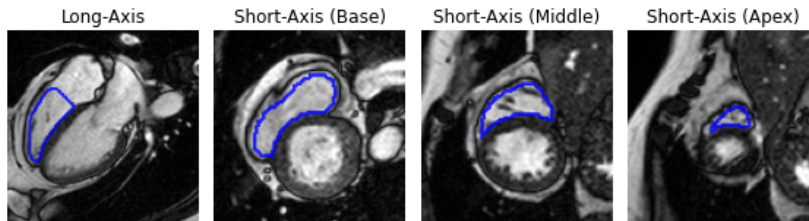


Fig. 5. A representative example from the development set of RV segmentation on the LA and SA images, showing how Tempera is able to delineate the RV successfully.

4 Discussion

Tempera network can perform RV segmentation well on a challenging dataset. By designing a model that interchanges between 2D and 3D space, we can reap several benefits. Using the SA as 2D slices rather than just a 3D volume, we are able to effectively increase the number of available training data. On the other

hand, when using 3D operations the network has the option to utilize structural information between slices to improve its performance. Finally, we exploit these geometric redundancies to improve the LA segmentations by using the GST.

The preprocessing ROI detection and standardization is an important step when dealing with data variation from different scan vendors. By detecting a small ROI which contains the RV, we can minimize erroneous segmentations, minimize the network learning irrelevant background features and allow for a smaller model. However, the detection may fail in the presence of pathology, artefacted or low SNR images and, in future work, we plan to improve this heart detection step. Short-term it can benefit from restricting the detection within a central region of the image. Additionally, we can take advantage of the spatial alignment of the LA and SA to project the centres from one image space to the other. In the long term, it can be improved by using a neural network-based detector.

There has been limited work done on incorporating the multiple views of cardiac MRI to improve the segmentation of cardiac chambers. [7] *et al* use the long axis views as shape priors for their short-axis cardiac segmentation. However, they restricted their investigations to the left myocardium, whose shape is much more regular than the RV's. Thus, their approach cannot be directly applied to the RV. Furthermore, the authors of [9] focus on preprocessing and data augmentation techniques to solve the data domain shift presented between different scan vendors and clinics, as part of the first M&Ms challenge [4]. For the RV segmentation, they obtained an average Dice score of 0.88.

5 Conclusions

For the M&Ms2 challenge, we propose the Tempera network to segment the right ventricle in the SA and LA cardiac MRI. Tempera uses both views and shares the learned features between them through hard-weight layer sharing in 2D space. A non-trainable spatial transformer, GST, allows for efficient sharing of information between the cardiac views.

Our Tempera network can seamlessly be extended to also segment the LV and the myocardium. Additionally, it can be updated to incorporate other chamber views, such as 2-Chamber and 3-Chamber. The only modification needed for this is the inclusion of additional convolutional branches.

Furthermore, we can improve the robustness of the network by utilising the temporal information (ED and ES). We can compute the deformation matrix across the cardiac cycle. Then, we can pass the deformation matrix to the GST to associate the temporal information between the ED and ES phases.

Acknowledgements

This work was supported by the UKRI CDT in AI for Healthcare <http://ai4health.io> (Grant No. EP/S023283/1) and the British Heart Foundation Centre of Research Excellence at Imperial College London (RE/18/4/34215).

References

1. Attili, A., Schuster, A., Nagel, E., et al.: Quantification in cardiac mri: Advances in image acquisition and processing. *The international journal of cardiovascular imaging* **26 Suppl 1**, 27–40 (02 2010). <https://doi.org/10.1007/s10554-009-9571-x>
2. Balakrishnan, G., Zhao, A., Sabuncu, M.R., others.: Voxelmorph: A learning framework for deformable medical image registration. *IEEE Transactions on Medical Imaging* **38**(8), 1788–1800 (2019)
3. Bonnemains, L., Mandry, D., Marie, P., et al.: Assessment of right ventricle volumes and function by cardiac mri: Quantification of the regional and global interobserver variability. *Magnetic resonance in medicine* **67**, 1740–6 (06 2012)
4. Campello, V.M., Gkontra, P., Izquierdo, C., et al.: Multi-centre, multi-vendor and multi-disease cardiac segmentation: The m amp;ms challenge. *IEEE Transactions on Medical Imaging* pp. 1–1 (2021). <https://doi.org/10.1109/TMI.2021.3090082>
5. Caudron, J., Fares, J., Vivier, P., et al.: Diagnostic accuracy and variability of three semi-quantitative methods for assessing right ventricular systolic function from cardiac mri in patients with acquired heart disease. *European Radiology* **21**, 2111–2120 (2011)
6. Caudron, J., Fares, J., Lefebvre, V., et al.: Cardiac mr assessment of right ventricular function in acquired heart disease: Factors of variability. *Acad Radiol.* (2012)
7. Chen, C., Biffi, C., Tarroni, G., et al.: Learning Shape Priors for Robust Cardiac MR Segmentation from Multi-view Images. In: *Lecture Notes in Computer Science*. vol. 11765 LNCS, pp. 523–531. Springer (oct 2019)
8. Friedberg, M., Redington, A.: Right versus left ventricular failure differences, similarities, and interactions. *Circulation* **129**, 1033–44 (03 2014)
9. Full, P.M., Isensee, F., Jäger, P.F., Maier-Hein, K.: Studying robustness of semantic segmentation under domain shift in cardiac mri. In: *Statistical Atlases and Computational Models of the Heart. M&Ms and EMIDEC Challenges*. pp. 238–249. Springer International Publishing, Cham (2021)
10. Hu, J., Shen, L., Sun, G.: Squeeze-and-excitation networks. In: *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 7132–7141 (2018). <https://doi.org/10.1109/CVPR.2018.00745>
11. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization (2017)
12. Klambauer, G., Unterthiner, T., Mayr, A., Hochreiter, S.: Self-normalizing neural networks. In: *31st International Conference on Neural Information Processing Systems*. p. 972–981. NIPS’17, Curran Associates Inc. (2017)
13. LeCun, Y., Bottou, L., Orr, G.B., Müller, K.R.: Efficient backprop. In: *Neural Networks: Tricks of the Trade, This Book is an Outgrowth of a 1996 NIPS Workshop*. p. 9–50. Springer-Verlag, Berlin, Heidelberg (1998)
14. Lin, T.Y., Dollár, P., Girshick, R., , et al.: Feature pyramid networks for object detection. In: *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. pp. 936–944 (2017). <https://doi.org/10.1109/CVPR.2017.106>
15. Lin, T.Y., Goyal, P., Girshick, R., et al.: Focal loss for dense object detection (2018)
16. Martín-Isla, C., Palomares, J.F.R., Guala, A., et al.: Multi-Disease, Multi-View & Multi-Center Right Ventricular Segmentation in Cardiac MRI (M&Ms-2) (Mar 2021). <https://doi.org/10.5281/zenodo.4573984>
17. Petitjean, C., Zuluaga, M.A., Bai, W., et al.: Right ventricle segmentation from cardiac mri: A collation study. *Medical Image Analysis* **19**(1), 187–202 (2015)
18. Zhu, W., Huang, Y., Zeng, L., et al.: Anatomynet: Deep learning for fast and fully automated whole-volume segmentation of head and neck anatomy. *Medical Physics* **46**(2), 576–589 (2019)