

# Deep Generative Model-based Quality Control for Cardiac MRI Segmentation

✉<sup>1</sup>, ✉<sup>1</sup>, ✉<sup>2</sup>, ✉<sup>2</sup>, ✉<sup>2</sup>, ✉<sup>1</sup>, ✉<sup>1</sup>, ✉<sup>2</sup>, ✉<sup>2</sup>, ✉<sup>1,3</sup>

<sup>1</sup> Data Science Institute, Imperial College London, London, UK

<sup>2</sup> BioMedIA Group, Department of Computing, Imperial College London, UK

<sup>3</sup> Department of Brain Sciences, Imperial College London, UK

shuo.wang@imperial.ac.uk

**Abstract.** In recent years, convolutional neural networks have demonstrated promising performance in a variety of medical image segmentation tasks. However, when a trained segmentation model is deployed into the real clinical world, the model may not perform optimally. A major challenge is the potential poor-quality segmentations generated due to degraded image quality or domain shift issues. There is a timely need to develop an automated quality control method that can detect poor segmentations and feedback to clinicians. Here we propose a novel deep generative model-based framework for quality control of cardiac MRI segmentation. It first learns a manifold of good-quality image-segmentation pairs using a generative model. The quality of a given test segmentation is then assessed by evaluating the difference from its projection onto the good-quality manifold. In particular, the projection is refined through iterative search in the latent space. The proposed method achieves high prediction accuracy on two publicly available cardiac MRI datasets. Moreover, it shows better generalisation ability than traditional regression-based methods. Our approach provides a real-time and model-agnostic quality control for cardiac MRI segmentation, which has the potential to be integrated into clinical image analysis workflows.

**Keywords:** Cardiac segmentation · Quality control · Generative model.

## 1 Introduction

✉<sup>1</sup>, ✉<sup>1</sup>, ✉<sup>2</sup>, ✉<sup>2</sup>, ✉<sup>2</sup>, ✉<sup>1</sup>, ✉<sup>1</sup>, ✉<sup>2</sup>, ✉<sup>2</sup>, ✉<sup>1,3</sup>

$\mathbb{R}^n$   
 If  $\mathbf{x} \in \mathbb{R}^n$   
 $\mathbb{R}^n$  is a  
 is a  
 strictly  
 h(Q) is  
 a

**Related work:**

[6, 7, 8], [9, 10]  
 is  
 a

**Learning-based quality control:**

[11], [12]  
 85%  
 [2]  
 [3]  
 [4]  
 [5]  
 [6]

**Registration-based quality control:**

[7]  
 [8]  
 [9]  
 [10]  
 [11]  
 [12]  
 [13]

**Contributions:**

[14]  
 [15]  
 [16]  
 [17]  
 [18]

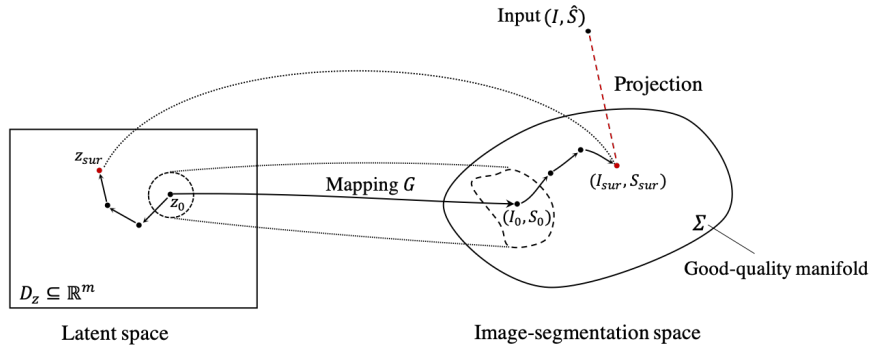
## 2 Methodology

### 2.1 Problem Formulation

Let  $F$  denote the generative model that maps an image  $I$  to a segmentation  $S$ , i.e.,  $S = F(I)$ . The quality of the segmentation is measured by a quality function  $q(S; I)$ . The goal is to find a segmentation  $\hat{S}$  that maximizes the quality function, i.e.,  $Q(\hat{S}; I) \approx q(S_{gt}, \hat{S})$ .

### 2.2 Deep Generative Model-based Quality Control

The deep generative model-based quality control framework is defined as follows. Let  $G$  denote the generative model that maps a latent variable  $z$  to a segmentation  $S$ , i.e.,  $S = G(z)$ . The quality of the segmentation is measured by a quality function  $q(S; I)$ . The goal is to find a latent variable  $z_{sur}$  that maximizes the quality function, i.e.,  $q(S_{sur}, \hat{S}) \approx q(S_{gt}, \hat{S})$ .



**Fig. 1.** Overview of the deep generative model-based quality control framework. The generative model  $G$  is trained to learn a mapping  $G(z)$  from the low-dimensional latent space  $D_z$  to the good-quality manifold  $\Sigma$ . The input image-segmentation pair  $(I, \hat{S})$  is projected to  $(S_{sur}, I_{sur})$  on the manifold through iterative search, which is in turn used as surrogate ground truth for quality prediction.  $z_0$  is the initial guess in the latent space and it converges to  $z_{sur}$ .

**Good-quality manifold:**  $\Sigma \subset D_I \times D_S$ , where  $D_I$  and  $D_S$  are the image and segmentation spaces, respectively. The manifold  $\Sigma$  is defined as the set of all image-segmentation pairs  $(I, S)$  that are of good quality, i.e.,  $(I, S) \in \Sigma$  if and only if  $q(S; I) \geq \tau$ , where  $\tau$  is a threshold.





ACDC dataset: 100 samples

split into  
training  
validation

### 3.2 Experimental Design

Model  
Architecture  
Implementation

VAE implementation and training:

Loss:  $ReLU$   
Activation:  $Sigmoid$   
Learning rate:  $\beta = 0.01$   
Batch size:  $m = 8$   
Optimizer: Adam  
Number of epochs: 100000  
Number of layers: 16  
Number of nodes per layer: 64  
Number of nodes in bottleneck layer: 16  
Number of nodes in output layer: 64  
Number of nodes in input layer: 64  
Number of nodes in hidden layer: 64

Baseline methods:

1)  $log(N)$   
2)  $log(N+1)$   
3)  $log(N+2)$

Experiment 1: UK Biobank

Model  
Architecture  
Implementation

Experiment 2: ACDC

Model  
Architecture

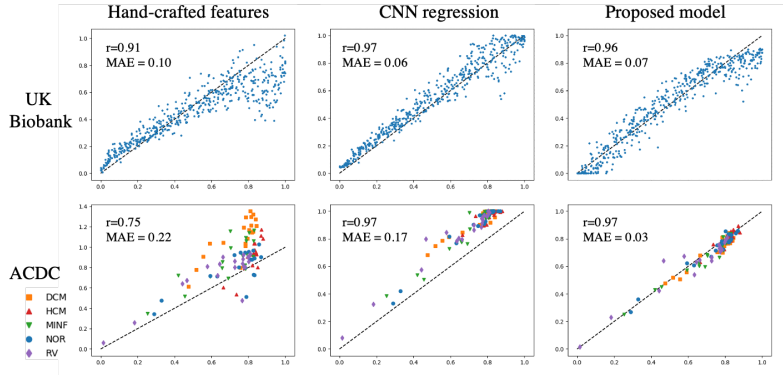
Ablation study:

Model  
Architecture

$m = 8$

### 4 Results and Discussion

The performance of different quality control methods is compared in Figure 3. The x-axis represents the real Dice of each subject, and the y-axis represents the predicted Dice by each method. A dashed line indicates the  $y = x$  reference. The top row shows UK Biobank data ( $n = 500$ ), and the bottom row shows ACDC data ( $n = 100$ ), with five subgroups plotted in different colors.



**Fig. 3.** Comparison of the performance of different quality control methods. The x-axis is the real Dice of each subject and the y-axis is the predicted Dice by each method. The dashed line is the  $y = x$ , plotted for reference. Top row: UK Biobank data ( $n = 500$ ). Bottom row: ACDC data ( $n = 100$ ), with five subgroups plotted in different colors.

The performance of different quality control methods is compared in Figure 3. The x-axis represents the real Dice of each subject, and the y-axis represents the predicted Dice by each method. A dashed line indicates the  $y = x$  reference. The top row shows UK Biobank data ( $n = 500$ ), and the bottom row shows ACDC data ( $n = 100$ ), with five subgroups plotted in different colors.

**Table 1.** Quality control performance of three models on two cardiac datasets. The Pearson correlation coefficient  $r$  and the mean absolute error (MAE) between predicted and true Dice metrics are reported. MAE is reported as mean (standard deviation). For ACDC dataset, the performance on five subgroups [4] are also reported.

Dataset	Hand-crafted features [11]		CNN regression [12]		Proposed model	
	$r$	MAE	$r$	MAE	$r$	MAE
UK Biobank	0.909	0.100(0.100)	<b>0.973</b>	<b>0.061(0.049)</b>	0.958	0.067(0.052)
ACDC	0.728	0.182(0.130)	0.968	0.165(0.044)	<b>0.969</b>	<b>0.033(0.028)</b>
DCM	0.802	0.353(0.123)	0.956	0.186(0.034)	<b>0.964</b>	<b>0.036(0.020)</b>
HCM	0.838	0.131(0.075)	0.836	0.155(0.027)	<b>0.896</b>	<b>0.023(0.015)</b>
MINF	0.815	0.184(0.123)	0.969	0.156(0.051)	<b>0.976</b>	<b>0.033(0.029)</b>
NOR	0.775	0.114(0.065)	0.985	0.158(0.040)	<b>0.985</b>	<b>0.026(0.017)</b>
RV	0.877	0.129(0.073)	0.974	0.169(0.053)	<b>0.960</b>	<b>0.045(0.042)</b>

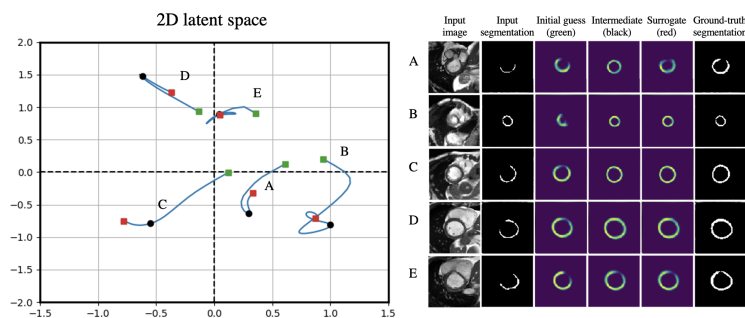
The proposed method is compared with the state-of-the-art methods in Table 1. The proposed method achieves the highest correlation coefficient  $r$  and the lowest MAE among all methods. The proposed method also achieves the highest correlation coefficient  $r$  and the lowest MAE among all methods. The proposed method also achieves the highest correlation coefficient  $r$  and the lowest MAE among all methods.

*prior to*

## 5 Conclusion

The proposed method is compared with the state-of-the-art methods in Table 1. The proposed method achieves the highest correlation coefficient  $r$  and the lowest MAE among all methods. The proposed method also achieves the highest correlation coefficient  $r$  and the lowest MAE among all methods. The proposed method also achieves the highest correlation coefficient  $r$  and the lowest MAE among all methods.





**Fig. 4.** Visualisation of searching path in a two-dimensional latent space. Left: searching paths for five exemplar samples (green point: initial guess from the VAE encoder; black point: intermediate state during iterative search; red point: convergence point for surrogate segmentation). Right: the input image and segmentation, reconstructed segmentations along the searching path and the ground-truth segmentation.

## References

1. WHO: Scale up prevention of heart attack and stroke. [https://www.who.int/cardiovascular\\_diseases/world-heart-day/en/](https://www.who.int/cardiovascular_diseases/world-heart-day/en/) accessed 16 March 2020.
2. Bai, W., Sinclair, M., Tarroni, G., Oktay, O., Rajchl, M., Vaillant, G., Lee, A.M., Aung, N., Lukaschuk, E., Sanghvi, M.M., et al.: Automated cardiovascular magnetic resonance image analysis with fully convolutional networks. *Journal of Cardiovascular Magnetic Resonance* **20**(1) (2018) 65
3. Tao, Q., Yan, W., Wang, Y., Paiman, E.H., Shamonin, D.P., Garg, P., Plein, S., Huang, L., Xia, L., Sramko, M., et al.: Deep learning-based method for fully automatic quantification of left ventricle function from cine mr images: a multivendor, multicenter study. *Radiology* **290**(1) (2019) 81–88
4. Bernard, O., Lalande, A., Zotti, C., Cervenansky, F., Yang, X., Heng, P.A., Cetin, I., Lekadir, K., Camara, O., Ballester, M.A.G., et al.: Deep learning techniques for automatic mri cardiac multi-structures segmentation and diagnosis: is the problem solved? *IEEE transactions on medical imaging* **37**(11) (2018) 2514–2525
5. Zheng, Q., Delingette, H., Duchateau, N., Ayache, N.: 3-d consistent and robust segmentation of cardiac images by deep learning with spatial propagation. *IEEE transactions on medical imaging* **37**(9) (2018) 2137–2148
6. Tarroni, G., Oktay, O., Bai, W., Schuh, A., Suzuki, H., Passerat-Palmbach, J., de Marvao, A., ORegan, D.P., Cook, S., Glocker, B., et al.: Learning-based quality control for cardiac mr images. *IEEE transactions on medical imaging* **38**(5) (2018) 1127–1138
7. Carapella, V., Jiménez-Ruiz, E., Lukaschuk, E., Aung, N., Fung, K., Paiva, J., Sanghvi, M., Neubauer, S., Petersen, S., Horrocks, I., et al.: Towards the semantic enrichment of free-text annotation of image quality assessment for uk biobank cardiac cine mri scans. In: *Deep Learning and Data Labeling for Medical Applications*. Springer (2016) 238–248
8. Zhang, L., Gooya, A., Dong, B., Hua, R., Petersen, S.E., Medrano-Gracia, P., Frangi, A.F.: Automated quality assessment of cardiac mr images using convolu-

- tional neural networks. In: International Workshop on Simulation and Synthesis in Medical Imaging, Springer (2016) 138–145
9. Robinson, R., Valindria, V.V., Bai, W., Oktay, O., Kainz, B., Suzuki, H., Sanghvi, M.M., Aung, N., Paiva, J.M., Zemrak, F., et al.: Automated quality control in image segmentation: application to the uk biobank cardiovascular magnetic resonance imaging study. *Journal of Cardiovascular Magnetic Resonance* **21**(1) (2019) 18
  10. Albà, X., Lekadir, K., Pereanez, M., Medrano-Gracia, P., Young, A.A., Frangi, A.F.: Automatic initialization and quality control of large-scale cardiac mri segmentations. *Medical image analysis* **43** (2018) 129–141
  11. Kohlberger, T., Singh, V., Alvino, C., Bahlmann, C., Grady, L.: Evaluating segmentation error without ground truth. In: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer (2012) 528–536
  12. Robinson, R., Oktay, O., Bai, W., Valindria, V.V., Sanghvi, M.M., Aung, N., Paiva, J.M., Zemrak, F., Fung, K., Lukaschuk, E., et al.: Real-time prediction of segmentation quality. In: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer (2018) 578–585
  13. Hann, E., Biasioli, L., Zhang, Q., Popescu, I.A., Werys, K., Lukaschuk, E., Carapella, V., Paiva, J.M., Aung, N., Rayner, J.J., et al.: Quality control-driven image segmentation towards reliable automatic image analysis in large-scale cardiovascular magnetic resonance aortic cine imaging. In: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer (2019) 750–758
  14. Liu, F., Xia, Y., Yang, D., Yuille, A.L., Xu, D.: An alarm system for segmentation algorithm based on shape model. In: Proceedings of the IEEE International Conference on Computer Vision. (2019) 10652–10661
  15. Valindria, V.V., Lavdas, I., Bai, W., Kamnitsas, K., Aboagye, E.O., Rockall, A.G., Rueckert, D., Glocker, B.: Reverse classification accuracy: predicting segmentation performance in the absence of ground truth. *IEEE transactions on medical imaging* **36**(8) (2017) 1597–1606
  16. Haskins, G., Kruger, U., Yan, P.: Deep learning in medical image registration: a survey. *Machine Vision and Applications* **31**(1) (2020) 1–18
  17. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. In: International Conference on Learning Representations. (2014)
  18. Higgins, I., Matthey, L., Pal, A., Burgess, C., Glorot, X., Botvinick, M., Mohamed, S., Lerchner, A.: beta-VAE: Learning basic visual concepts with a constrained variational framework. In: International Conference on Learning Representations. (2017)
  19. Bojanowski, P., Joulin, A., Lopez-Paz, D., Szlam, A.: Optimizing the latent space of generative networks. In: International Conference on Machine Learning. (2018)

## Appendix

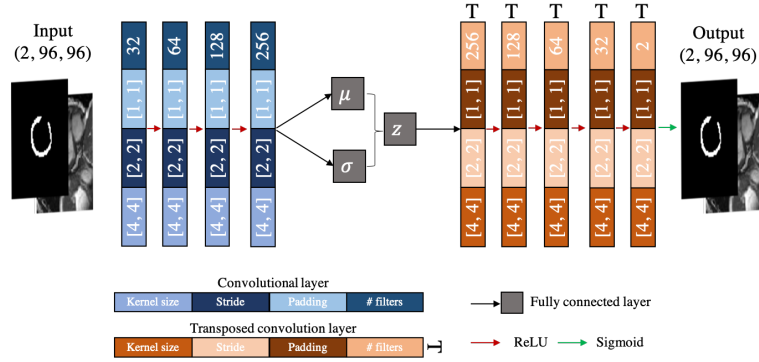


Fig. A1. Network architecture of the VAE.

Table A1. Ablation study of the latent space dimension  $m$  and the hyperparameter  $\beta$ . The mean absolute error (MAE) between predicted and true Dice metrics on UKBB validation set are reported.  $m=8$  and  $\beta = 0.01$  were selected as the optimal parameters according to the results.

$m$	$\beta$			
	0	1E-3	1E-2	1E-1
2	0.095	0.094	0.105	0.502
4	0.086	0.092	0.114	0.416
8	0.088	0.072	<b>0.068</b>	0.229
16	0.147	0.141	0.095	0.091