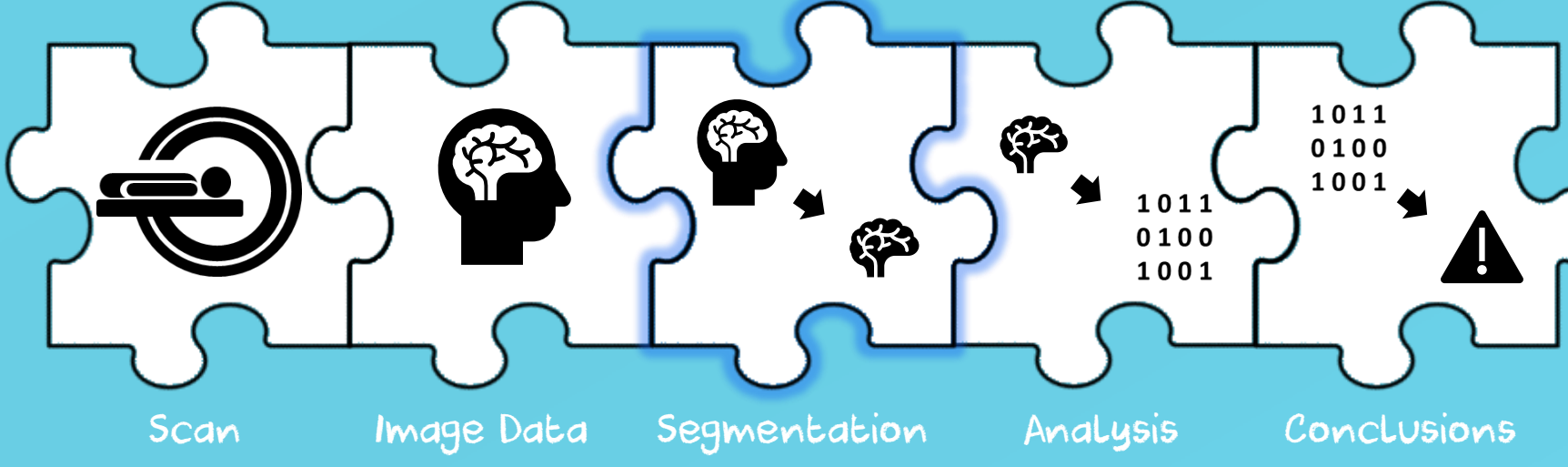


# Real-time Prediction of Segmentation Quality Swapping Expert Annotations for Generated Labels

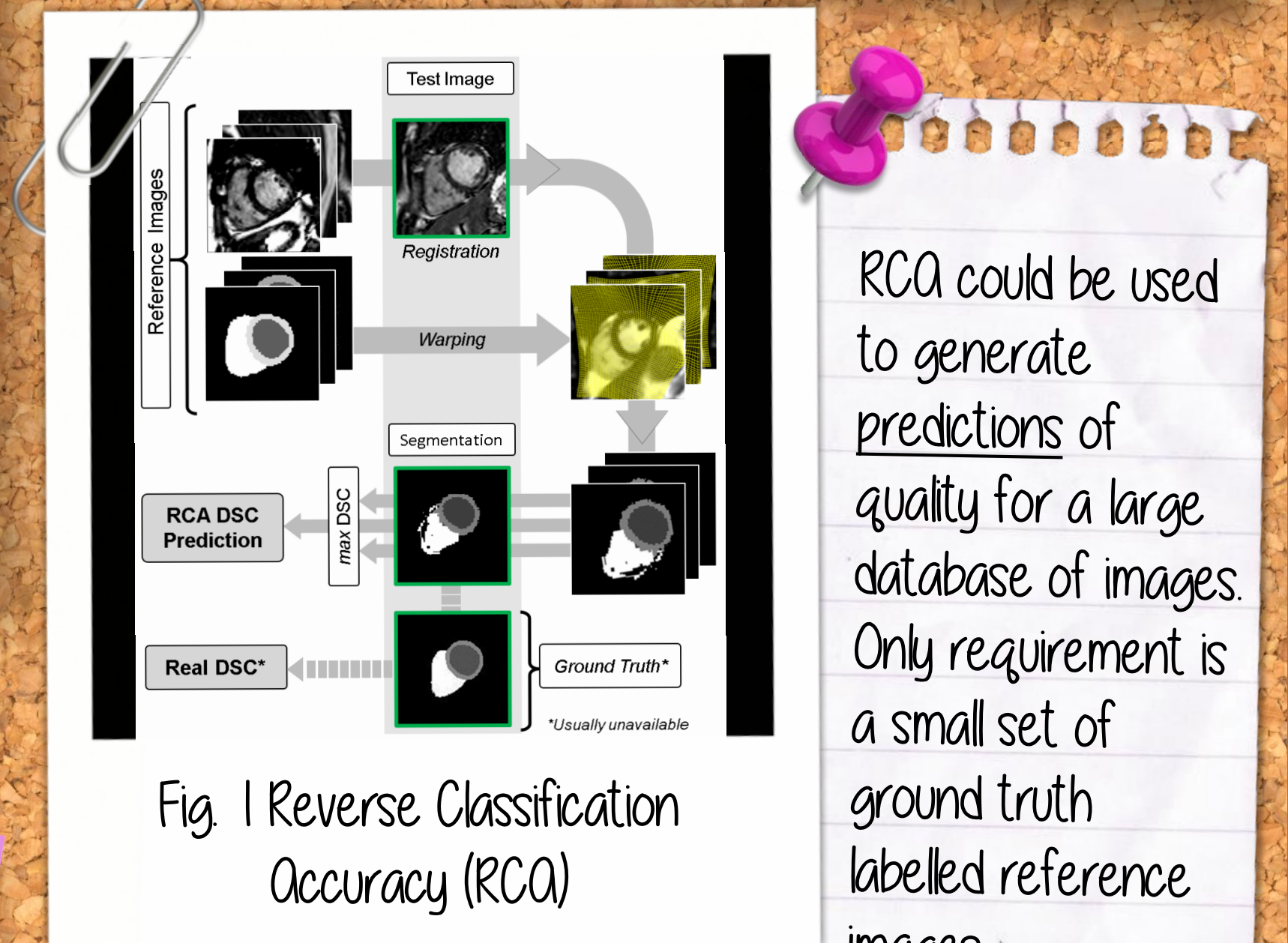
Rob Robinson<sup>1</sup>, O. Oktay<sup>1</sup>, W. Bai<sup>1</sup>, V.V. Valindria<sup>1</sup>, M.M. Sanghvi<sup>3,4</sup>, N. Oung<sup>3,4</sup>, J.M. Paiva<sup>3</sup>, F. Zemrak<sup>3,4</sup>, K. Fung<sup>3,4</sup>, E. Lukaschuk<sup>5</sup>, A.M. Lee<sup>3,4</sup>, V. Carapella<sup>5</sup>, Y.J. Kim<sup>5,6</sup>, B. Kainz<sup>1</sup>, S.K. Piechnik<sup>5</sup>, S. Neubauer<sup>5</sup>, S.E. Petersen<sup>3,4</sup>, C. Page<sup>2</sup>, D. Rueckert<sup>1</sup>, B. Glocker<sup>1</sup>

<sup>1</sup> BioMedIA Group, Imperial College London, UK  
<sup>2</sup> Research & Development, GlaxoSmithKline, UK  
<sup>3</sup> NIHR Barts, Queen Mary's University UK  
<sup>4</sup> Barts Heart Centre, Barts NHS Trust, London, UK  
<sup>5</sup> Radcliffe Dept. Medicine, University of Oxford, UK  
<sup>6</sup> Yonsei University College of Medicine, South Korea  
 ✉ rrobinson16@imperial.ac.uk

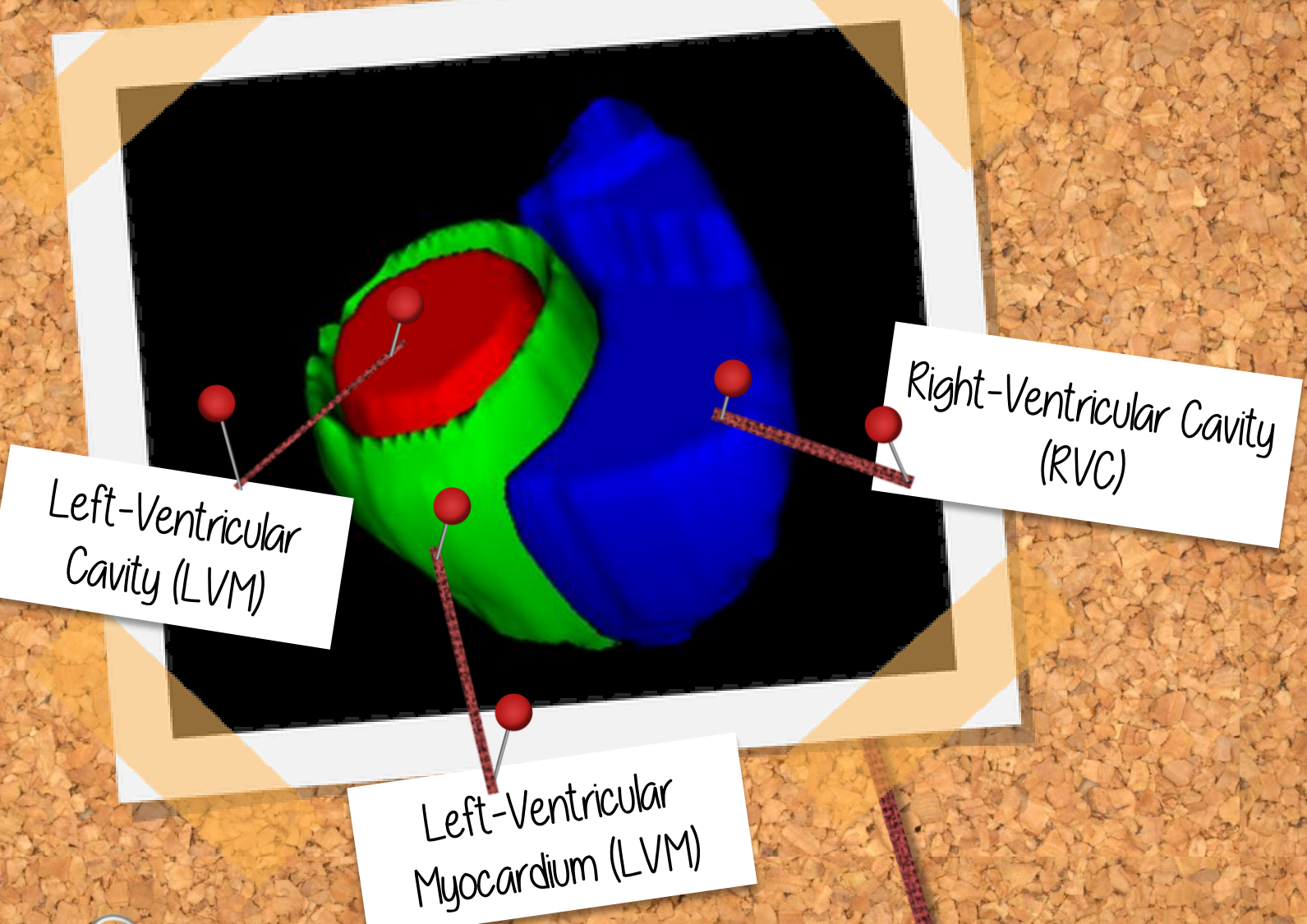
Can we do deep learning without the large manually-annotated dataset?



Desirable to perform quality control on-the-fly in real time. Deep Learning is ideal, but requires a large, labelled dataset. These don't exist in practice. So what can we do?



Advantages:  
 1. No need for expensive, expert manual labelling  
 2. Extend the training dataset by segmenting at different qualities (e.g. random forest depths)



In both experiments we use the same architecture.  
 Input: [224, 224, 85] array where the 5 channels are the 3D image and 4 one-hot-encoded segmentations, one for each of the classes: background, left-ventricular cavity, left-ventricular myocardium and right ventricular cavity. Network: A normal ResNet50 Residual Network  
 Output: [1, 5] array of predicted Dice Similarity Coefficients (DSC) for each class and one for the binary 'whole heart' case.

**Data** **biobank**  
 16,100 2D Stack CMR + Random Forest segmentations. Manual annotations available

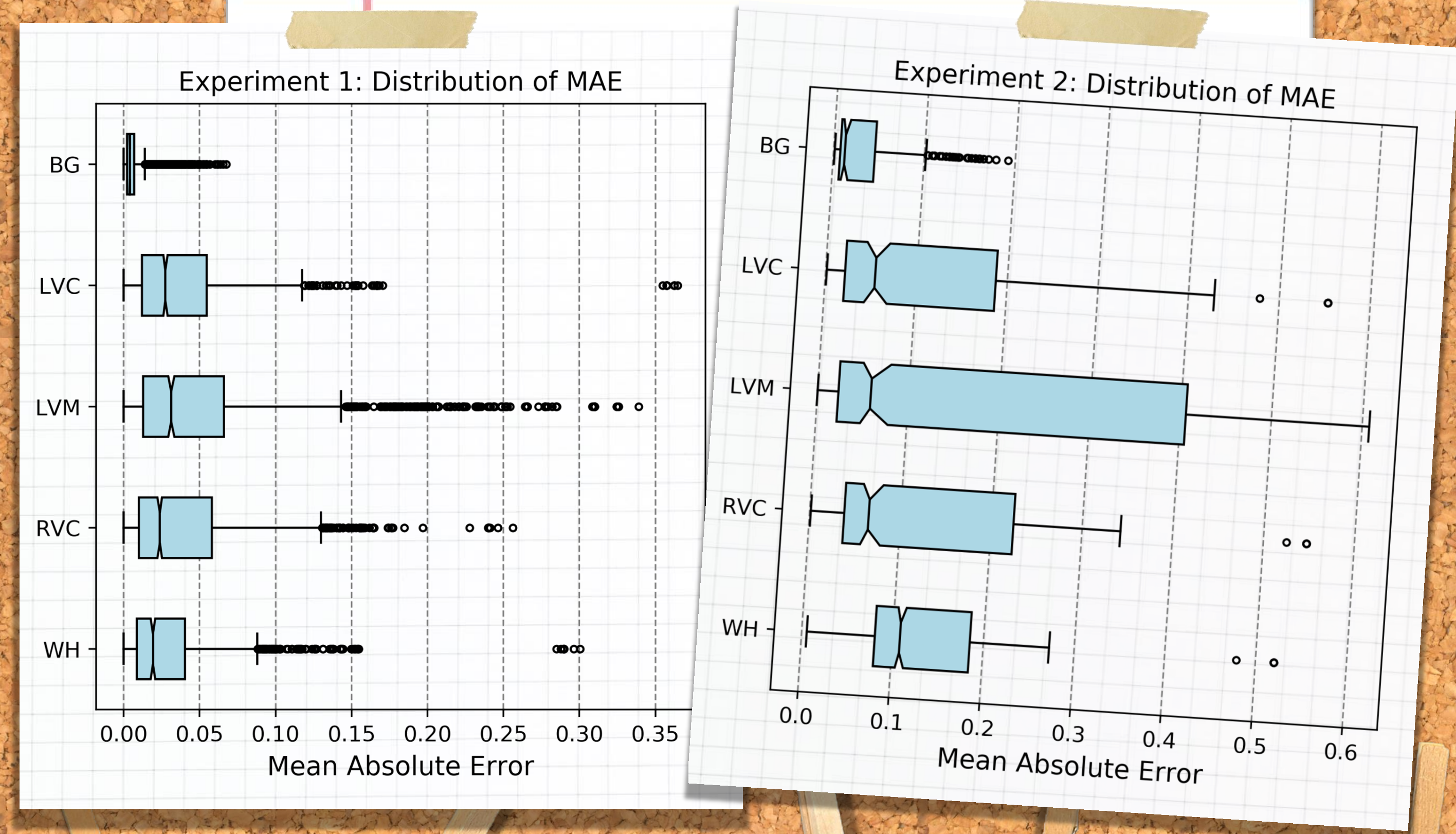
Use a 'reverse testing' strategy to generate a large set of predictions  
 Train a network to predict the predictions

Compare network prediction of RCA predictions (Experiment 2) with network predictions of real DSC from ground truth, manual labelmaps (Experiment 1)

Fig. 3: Results from Both Experiments  
 Table (top) showing mean absolute error for the different classes from Experiments 1 (left) and 2 (right). MOE for good (DSC ≥ 0.5) and poor (DSC < 0.5) ranges shown to be of similar order. True (TPR) and false (FPR) positive rates are shown below. Standard deviation in brackets.

Class	Experiment 1			Experiment 2		
	0 ≤ DSC ≤ 1 n = 1,610	DSC < 0.5 n = 817	DSC ≥ 0.5 n = 793	0 ≤ DSC ≤ 1 n = 288	DSC < 0.5 n = 160	DSC ≥ 0.5 n = 128
BG	0.008 (0.011)	0.012 (0.014)	0.004 (0.002)	0.034 (0.042)	0.048 (0.046)	0.074 (0.002)
LV	0.038 (0.040)	0.025 (0.024)	0.053 (0.047)	0.120 (0.128)	0.069 (0.125)	0.213 (0.065)
LVM	0.055 (0.064)	0.027 (0.027)	0.083 (0.078)	0.191 (0.218)	0.042 (0.041)	0.473 (0.111)
RVC	0.039 (0.041)	0.021 (0.020)	0.058 (0.047)	0.127 (0.126)	0.076 (0.109)	0.223 (0.098)
WH	<b>0.031 (0.035)</b>	<b>0.018 (0.018)</b>	<b>0.043 (0.043)</b>	<b>0.139 (0.091)</b>	<b>0.112 (0.093)</b>	<b>0.188 (0.060)</b>
	TPR 0.975	FPR 0.060	Acc. 0.965	TPR 0.879	FPR 0.000	Acc. 0.906

Graphs (bottom) of MOE over full DSC range. Low MOE for WH case in both experiments. Expected larger errors for Experiment 2 where predictions are used as GT labels.



What's the point?? There are still some things that need a little more consideration...

- ✓ Real-time feedback to technicians? 40 ms!
- ✓ Is the image analyzable?
- ✓ Reacquisition with patient in-situ? Reduced stress
- ✓ Cost/efficiency benefits?

- ✓ Network will learn how to assess Random Forest segmentations. So train the network with CNN/ensemble of segmentations
- ✓ Train the network on the logit transform of the DSC to give more fine-grain prediction at higher DSC

**References**  
 [1] Valindria, V.V., Laidis, I., Bai, W., Karamitsos, K., Odojny, E.O., Rookal, A.G., Rueckert, D., Glocker, B. 10, Feb. 2017.  
 [2] Petersen, S.E. et al. "Reference ranges for cardiac structure and function using cardiovascular magnetic resonance (CMR) in Caucasians from the UK Biobank population cohort." JCMR 19(1) (2017)

Fig. 4: Example Predictions  
 Predictions of network trained in Experiment 1. Showing (top) the midsagittal slice from 2D stack of CMR, (middle) random forest generated segmentation and (bottom) the manual ground truth labelmap. Examples segmentations are ordered by increasing quality from left to right. Real and predicted DSC shown with mean absolute error (MAE)

