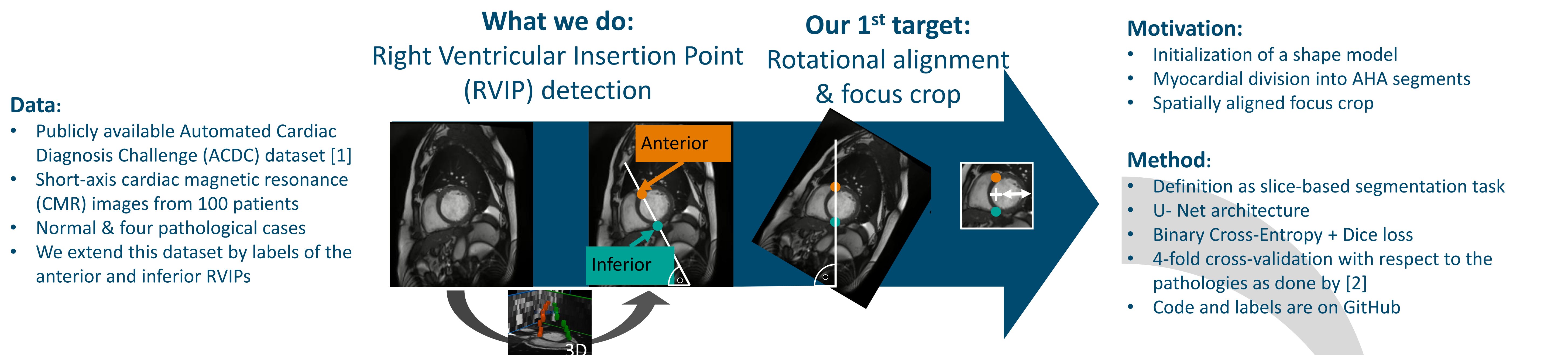


# Comparison of Evaluation Metrics for Landmark Detection in CMR Images

Sven Koehler<sup>1,2</sup>, Lalith Sharan<sup>1,2</sup>, Julian Kuhm<sup>1</sup>, Arman Ghanaat<sup>1</sup>, Jelizaveta Gordejeva<sup>1</sup>, Nike K. Simon<sup>1</sup>, Niko M. Grell<sup>1</sup>, Florian André<sup>1</sup>, Sandy Engelhardt<sup>1,2</sup>

<sup>1</sup> Department of Internal Medicine III, Heidelberg University Hospital, Heidelberg

<sup>2</sup> DZHK (German Centre for Cardiovascular Research), partner site Heidelberg/Mannheim



**Fig. 1.** Metric definition in terms of detection and localization for evaluation is not straight forward. Different definitions yield different results. (a) The total number of true positives (TP) and false positives (FP) are influential for accuracy measurements. (b) Each detection metric could be combined with a localization metric – with or without an upper boundary.

How to find a goal tailored TP/FP definition?

	Ground truth (g)	Prediction (p)	Threshold (t)	Ant	Inf	Euclidean distance
Base	+	+	+	+	+	+
Ant	+	+	+	+	+	+
Inf	+	+	+	+	+	+
Euclidean distance	+	+	+	+	+	+

	Base	Ant	Inf	Apex
I. Line (Septum approx.)	TN	FN	FN	TP
II. Points (single)	TN	FN	TP	TP
III. Threshold Points	TN	FN	TP	TP
if $d(g, p) < t, TP$ else FP	TN	TP	FN	TP

**a) Detection-metrics**

**b) Localization-metrics**

For each Detection-metric

**I) Volume-based**

1. average  $g$  and  $p$

2.  $d(g_{vol}, p_{vol})$

**II) Slice-based**

1.  $d(g_z, p_z); z \in \text{slices}$

2. average  $\bar{d}$

**III) with upper-bound**

farthest corner  $c$

Upper boundary for FP:  $|\bar{d}_{FP}| = d(g, c)$

**Tab. 1** Localization metric comparison for different experiments (Base: Baseline, Var.1 + hist. matching, Var.2 + Gauss  $\sigma = 2$ , Var.3: +  $\sigma = 4$ ).

Detection-Strategy	Exp.	(i) Line	(ii) Points
(i) Volume-based	Base	5.92 ± 4.83	3.86 ± 5.32
	Var.1	5.58 ± 6.25	4.16 ± 5.75
	Var.2	6.26 ± 7.08	3.54 ± 3.83
	Var.3	5.86 ± 4.95	3.33 ± 3.47
(ii) Slice-based	Base	4.42 ± 5.66	3.96 ± 7.07
	Var.1	3.79 ± 7.20	3.02 ± 4.39
	Var.2	3.88 ± 4.97	3.12 ± 7.10
	Var.3	4.42 ± 5.67	2.48 ± 2.20
(iii) Slice-based, $\uparrow$ -bound	Base	50.33 ± 65.01	49.68 ± 65.98
	Var.1	37.07 ± 46.70	36.83 ± 45.74
	Var.2	48.53 ± 64.66	47.58 ± 63.89
	Var.3	55.05 ± 74.88	53.76 ± 75.66

**Tab. 2** Detection metric comparison for different experiments (Base: Baseline, Var.1 + hist. matching, Var.2 + Gauss  $\sigma = 2$ , Var.3: +  $\sigma = 4$ ).

Detection-Strategy	Exp.	(i) Line	(ii) Points
Base & Thresh.	Base	0.84 ± 0.22	0.84 ± 0.22
	Var.1	0.88 ± 0.16	0.85 ± 0.19
	Var.2	0.85 ± 0.21	0.85 ± 0.23
	Var.3	0.82 ± 0.25	0.83 ± 0.26
Var.1 & Thresh.	Base	0.88 ± 0.16	0.85 ± 0.19
	Var.1	0.91 ± 0.15	0.96 ± 0.10
	Var.2	0.88 ± 0.19	0.92 ± 0.16
	Var.3	0.88 ± 0.21	0.89 ± 0.20
Var.2 & Thresh.	Base	0.88 ± 0.16	0.85 ± 0.19
	Var.1	0.91 ± 0.15	0.96 ± 0.10
	Var.2	0.88 ± 0.19	0.92 ± 0.16
	Var.3	0.88 ± 0.21	0.89 ± 0.20
Var.3 & Thresh.	Base	0.88 ± 0.16	0.85 ± 0.19
	Var.1	0.91 ± 0.15	0.96 ± 0.10
	Var.2	0.88 ± 0.19	0.92 ± 0.16
	Var.3	0.88 ± 0.21	0.89 ± 0.20

[1] Bernard O et al. Deep learning techniques for automatic mri cardiac multi-structures segmentation and diagnosis: Is the Problem Solved? IEEE Trans Med Imaging 37.11 (2018), pp. 2514–2525.

[2] Koehler S et al. How well do U-Net-based segmentation trained on adult cardiac magnetic resonance imaging data generalise to rare congenital heart diseases for surgical planning? Med Imaging Ed. by Fei B, Linte CA. SPIE, 2020, p. 55.