

# Multi-centre evaluation of atlas-based and deep learning contouring using a modified Turing Test

M. Gooding<sup>1</sup>, A. Smith<sup>1</sup>, D. Peressutti<sup>1</sup>, P. Aljabar<sup>1</sup>, E. Evans<sup>2</sup>, S. Gwynne<sup>3</sup>, C. Hammer<sup>4</sup>, H.J.M. Meijer<sup>5</sup>, R. Speight<sup>6</sup>, C. Welgemoed<sup>7</sup>, T. Lustberg<sup>8</sup>, J. van Soest<sup>9</sup>, A. Dekker<sup>8</sup>, W. van Elmpt<sup>8</sup>

<sup>1</sup>Mirada Medical Limited, Oxford, United Kingdom. <sup>2</sup>Velindre Cancer Centre, Clinical Oncology, Cardiff, United Kingdom. <sup>3</sup>South West Wales Cancer Centre, Clinical Oncology, Swansea, United Kingdom. <sup>4</sup>University Medical Center Groningen, Department of Radiation Oncology, Groningen, The Netherlands. <sup>5</sup>Radboud University Medical Center, Department of Radiation Oncology, Nijmegen, The Netherlands. <sup>6</sup>Leeds Cancer Centre, Leeds, United Kingdom. <sup>7</sup>Imperial College Healthcare NHS Trust, Radiotherapy Department, London, United Kingdom. <sup>8</sup>MAASTRO Clinic, Department of Radiation Oncology, Maastricht, The Netherlands.

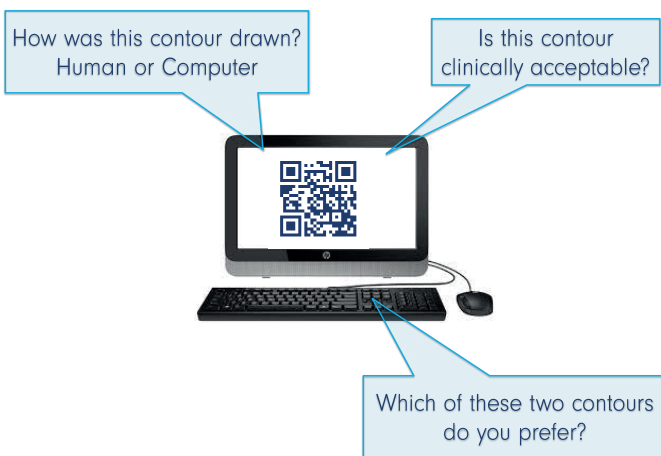
## Objective

While quantitative assessment of autocontouring quality is useful, frequently used measures do not necessarily indicate clinical acceptability or benefit. In contrast, clinical based assessment metrics, such as time saved with autocontouring or subjective evaluations, are both time consuming to perform and difficult to implement in a multi-centre evaluation.

Inspiration is taken from the Artificial Intelligence community to propose an assessment method based on the "Turing Test". The objective of this study was to perform a multi-centre evaluation of two autocontouring methods using this approach.

## Materials and Methods

A website was set up to facilitate multi-centre comparison, showing images and contours in a blinded fashion.

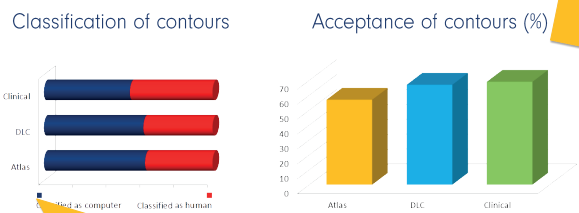


Use the QR-code to try it for yourself at <http://www.autocontouring.com>

- 60 Clinical test cases from a single institution (40 thoracic, 20 prostate)
- 3 Types of contour
  - Existing clinical
  - Multi-atlas contouring (Mirada Workflow Box 1.4)
  - Deep Learning Contouring (DLC) (Mirada DLCEXpert™ prototype)
- 15 Clinical participants (clinicians, dosimetrist or technicians)
- 5 Institutions participated
- 100 Randomly chosen questions answer by each participant

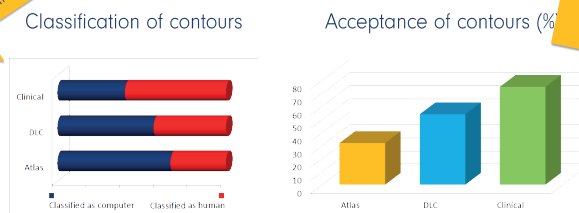
## Results

### Thoracic segmentation



DLC acceptance rate close that of clinical contouring. Note: Only 10% of clinical contours were accepted

### Prostate segmentation



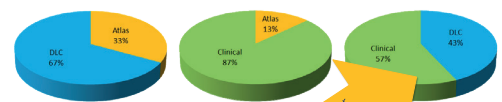
DLC shows increased clinical acceptance compared to atlas contouring

The source of the contours is often misidentified by the blinded participants

### Contour preference in blind side-by-side comparison



### Contour preference in blind side-by-side comparison



DLC has similar preference rate to the clinical contouring

## Conclusions

The Turing Test style assessment method provided an easy way to perform web-based multi-centre validation of autocontouring.

This study found that autocontours may be confused with clinical ones, when reviewed blindly. DLC showed increased clinical acceptance for prostate OAR contouring compared to atlas contouring. For thoracic imaging, DLC contours were accepted at a similar rate to clinical ones.