

CLINICAL EVALUATION OF ARTIFICIAL INTELLIGENCE-BASED SOFTWARE FOR AUTO-CONTOURING IN RADIATION THERAPY OF BREAST CANCER

Terzi G. Maria^{2,3}, Sakellaropoulos George^{1,3}, Spiridon N. Papatheodorou³, Despoina A. Spyropoulou^{2,4}

¹ Department of Medical Physics, School of Medicine, University of Patras, ² Department of Radiation Oncology, School of Medicine, University of Patras, ³ Department of Medical Physics, University Hospital of Patras, ⁴ Department of Radiation Oncology, University Hospital of Patras

Background: The aim of the study was to compare the auto-generated contours by the AI-based software “*Therapanacea*” with *manual - generated contours* by physicians. Can these software tools eliminate the quality evaluation of the experienced physician?

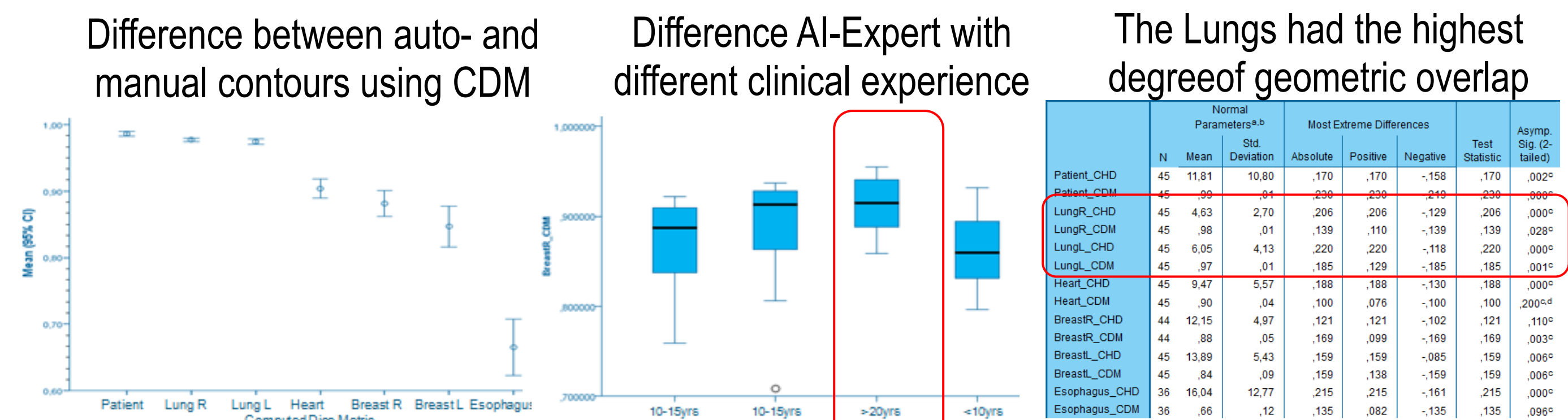
Materials and Methods: Four radiation oncologists with *different clinical experience* delineated the Clinical Target Volumes (CTVs) for Breast cancer radiotherapy and Organs at Risk (OARs) on 45 CT cases.

~Two volumetric overlap metrics, the Hausdorff distance (*CHD-95%*) and Coefficient Dice Metric (*CDM*), were used. The sample was divided into 14 groups according to the above metrics and 4 groups, one for each physician. Statistical analysis was performed.

~ In order to evaluate the dosimetric impact of AI: *Two plans* were created: one with contours by AI and the other using the manual-generated contours.

Results: Regarding the CHD-95%, statistically significant differences were found for the *Lung R* (p<0.05) and *Breast L* (p<0.05). No statistical difference was found for CDM (p>0.05).

Conclusion: This software could be used to *reduce* the *workflow time* and the *variability* among physicians. The quality evaluation of the experienced physician is *necessary*. The physician’s experience cannot be replaced by the metrics.



The AI software can *achieve similar* results with the *most experienced* physician. The differences between the auto- and manual-generated contours for *patient* and *esophagus* were due to a difference in *length*.

PTV Coverage	
PTV AI	93.83
PTV Expert	93.83

No difference for the PTV (Planning TargetVolume) coverage *despite the low degree* of geometric overlap

	CHD(95%)	CDM
Lung	3.75	0.98
Breast	15.00	0.86
Heart	25.25	0.81

PTV Coverage	
PTV AI	93.00
PTV Expert	95.15

Significant difference for the PTV coverage *and low degree* of geometric overlap

	CHD(95%)	CDM
Lung	9.21	0.98
Breast	10.23	0.94
Heart	3.95	0.93

There is a need to develop a *standardized method* for the validation of *any AI-based software*