

# Novel methodology to assess the effect of contouring variation on treatment outcome

Alexander Jenkins\*

*School of Physics & Astronomy - Faculty of Science and Engineering, University of Manchester, Manchester M13 9PL, UK  
Division of Cancer Studies - School of Medical Sciences - Faculty of Biology- Medicine and Health, University of Manchester, Manchester M20 4BX, UK*

Thomas Soares Mullen\*

*School of Physics & Astronomy - Faculty of Science and Engineering, University of Manchester, Manchester M13 9PL, UK  
Champalimaud Research, Champalimaud Centre for the Unknown, Avenida Brasília, Doca de Pedrouços, Lisboa 1400-038, Portugal*

Corinne Johnson-Hart

*Division of Cancer Studies - School of Medical Sciences - Faculty of Biology- Medicine and Health, University of Manchester, Manchester M20 4BX, UK  
Christie Medical Physics and Engineering, The Christie NHS Foundation Trust, Wilmslow Road, Manchester M20 4BX, UK*

Andrew Green, Alan McWilliam, Marianne Aznar, Marcel van Herk and

Eliana Vasquez Osorio<sup>a)</sup>

*Division of Cancer Studies - School of Medical Sciences - Faculty of Biology- Medicine and Health, University of Manchester, Manchester M20 4BX, UK*

(Received 25 November 2020; revised 22 March 2021; accepted for publication 22 March 2021; published 24 April 2021)

**Purpose:** Contouring variation is one of the largest systematic uncertainties in radiotherapy, yet its effect on clinical outcome has never been analyzed quantitatively. We propose a novel, robust methodology to locally quantify target contour variation in a large patient cohort and find where this variation correlates with treatment outcome. We demonstrate its use on biochemical recurrence for prostate cancer patients.

**Method:** We propose to compare each patient's target contours to a *consistent* and *unbiased* reference. This reference was created by auto-contouring each patient's target using an externally trained deep learning algorithm. Local contour deviation measured from the reference to the manual contour was projected to a common frame of reference, creating *contour deviation maps* for each patient. By stacking the contour deviation maps, time to event was modeled pixel-wise using a multivariate Cox proportional hazards model (CPHM). Hazard ratio (HR) maps for each covariate were created, and regions of significance found using cluster-based permutation testing on the z-statistics. This methodology was applied to clinical target volume (CTV) contours, containing only the prostate gland, from 232 intermediate- and high-risk prostate cancer patients. The reference contours were created using ADMIRE<sup>®</sup> v3.4 (Elekta AB, Sweden). Local contour deviations were computed in a spherical coordinate frame, where differences between reference and clinical contours were projected in a 2D map corresponding to sampling across the coronal and transverse angles every 3°. Time to biochemical recurrence was modeled using the pixel-wise CPHM analysis accounting for contour deviation, patient age, Gleason score, and treated CTV volume.

**Results:** We successfully applied the proposed methodology to a large patient cohort containing data from 232 patients. In this patient cohort, our analysis highlighted regions where the contour variation was related to biochemical recurrence, producing expected and unexpected results: (a) the interface between prostate–bladder and prostate–seminal vesicle interfaces where increase in the manual contour relative to the reference was related to a *reduction* of risk of biochemical recurrence by 4–8% per mm and (b) the prostate's right, anterior and posterior regions where an increase in the manual contour relative to the reference contours was related to an *increase* in risk of biochemical recurrence by 8–24% per mm.

**Conclusion:** We proposed and successfully applied a novel methodology to explore the correlation between contour variation and treatment outcome. We analyzed the effect of contour deviation of the prostate CTV on biochemical recurrence for a cohort of more than 200 prostate cancer patients while taking basic clinical variables into account. Applying this methodology to a larger dataset including additional clinically important covariates and externally validating it can more robustly identify regions where contour variation directly relates to treatment outcome. For example, in the prostate case we use to demonstrate our novel methodology, external validation will help confirm or reject the counter-intuitive results (larger contours resulting in higher risk). Ultimately, the results of this methodology could inform contouring protocols based on actual patient outcomes.

© 2021 The Authors. Medical Physics published by Wiley Periodicals LLC on behalf of American Association of Physicists in Medicine. [https://doi.org/10.1002/mp.14865]

Key words: contour variation, CPHM, data mining, treatment outcome modeling

## 1. INTRODUCTION

Prostate cancer is the second most common cancer in men worldwide<sup>1</sup> and approximately 30-60% of cases will be treated with curative-intent radiotherapy.<sup>2,3</sup> Radiotherapy relies on accurate definition of the target volume, that is, the region where the prescribed dose of radiation is to be delivered. However, target definition is subjective, and is known to vary between observers; known as interobserver variability (IOV). Variation in contouring systematically deforms the planned dose distribution relative to the (unknown) true target, potentially delivering lower dose than required to the cancer cells and influencing the effectiveness of treatment.<sup>4</sup>

Despite the known presence of IOV in target contours, as far as the authors are aware, its effect on clinical outcomes has never been analyzed quantitatively. Typically, IOV is quantified in studies where multiple observers contour the same structure on patient images.<sup>5-7</sup> As drawing these contours takes a long time and requires expertise, most studies include only a limited number of observers on a limited number of patient cases.<sup>7</sup> Furthermore, variations are often reduced to a single number, where all spatial information is lost; typically Dice similarity coefficient or Hausdorff distance.<sup>7</sup> For example, computing the Dice similarity coefficient of two structures involves only their respective volumes and overlap, this results in a single scalar value that contains no spatial information and it is impossible to infer how each spatial region contributed to the metric. These two drawbacks, a small cohort and the simplification of the contour differences into a single number, makes it impossible to effectively analyze the spatial effects of IOV on clinical outcome.

In recent years, automatic contouring of structures has been made possible using deep learning (DL) techniques.<sup>8</sup> For some organs, DL auto-contoured structures are of comparable quality to those drawn by an observer.<sup>9-12</sup> Additionally, advances in computational tools are also being used to automatize and improve clinical target volumes (CTVs) beyond anatomical organs.<sup>13</sup> Contours produced by DL contouring tools can therefore be seen as being drawn by a virtual observer that is highly consistent and unbiased to data beyond the medical image being segmented.

In this study, we present a novel methodology to quantify the effect of contour deviations on clinical outcome. The methodology relies on quantifying local contour deviations by comparing the manually delineated contour with the DL generated contour, which is used as an arbitrary yet consistent reference. These local contour deviations are then analyzed statistically to define regions where observer deviation correlates with outcome, taking clinical variables into account.<sup>14</sup> A major advantage of this approach is that instead of being restricted to a limited IOV study setting, it can exploit the information contained in large quantities of routine clinical

data. As a first application of this novel methodology, we analyzed the effect of contour deviation of the prostate CTV on biochemical recurrence for a cohort of prostate cancer patients treated with radical radiotherapy.

## 2. MATERIALS AND METHODS

### 2.A. Patient dataset

Two hundred and forty-seven intermediate- and high-risk prostate cancer patients, treated between 2007 and 2013 at a single institution (The Christie Hospital NHS Foundation Trust) with 57 Gy in 19 fractions of intensity-modulated radiotherapy were included in this study. Patients were setup via an empty bladder and rectum protocol, both for planning and treatment. Patients were followed up for at least 4 yr as standard of care and biochemical recurrence status, defined as a rise in the blood level of prostate-specific antigen (PSA) of 2 ng/ml above nadir after treatment, was stored for all patients. The characteristics of this cohort are summarized in Table I. Each patient had one CTV contour encompassing the prostate gland only, defined by the treating oncologist. This contour shall be referred to as the *manual contour*. All data were collected from the ukCAT distributed learning database (ethics approval from the UK North West - Haydock Research Ethics Committee, reference number 17/NW/0060, local approval consent ukCAT ref. 2018-018).

For each patient, a DL auto-contour of the prostate gland was generated as a reference using the research version of ADMIRE<sup>®</sup> v3.4, (Elekta AB, Sweden).

TABLE I. Characteristics of the cohort of prostate cancer patients, before and after refinement.

Variable	Original (n = 247)		Refined (n = 232)	
	Nr	%	Nr	%
Gleason score				
6	30	12	25	11
7	134	54	128	55
8	42	17	39	17
9-10	41	17	40	17
Age (years)				
<65	82	33	75	33
65-75	147	60	140	60
>75	18	7	17	7
Recurrence				
Yes	83	34	74	32
No	164	66	158	68

The refined dataset includes only patients whose clinical target volume was limited to the prostate gland, and without anomalies in their contour deviation maps. For patients with a biochemical recurrence, the mean time to event was 4.80 yr (0.85–8.46 yr).

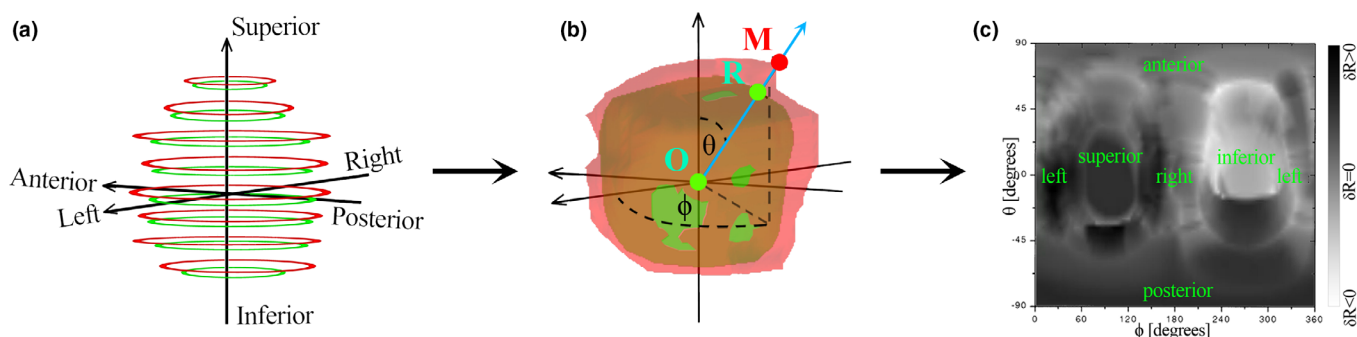


FIG. 1. Method to obtain a contour deviation map for a given patient from the manual and the DL reference contours of the CTV, as shown in (a). Contours are each triangulated to form a surface and a spherical polar coordinate system, centered at the center of mass of the auto-contour (O) is defined, as shown in (b). At each  $3^\circ \times 3^\circ$  angle, the difference between the distances from O to the manual contour and the DL reference contour  $\delta R(\theta, \phi)$  is calculated [Eq. (1)], to construct a single patient contour deviation map, shown in (c). Maps of hundreds of patients are next correlated with clinical outcome. [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

## 2.B. Quantifying local contour deviations

Figure 1 shows the steps for quantifying local contour deviations. The reference DL contour, and the manual contour are first triangulated into 3D surfaces using the marching cubes algorithm.<sup>15</sup> Then, using a similar approach as Remeijer *et al.*<sup>16</sup> local contouring deviation was measured in spherical polar coordinates. Contour deviation,  $\delta R$ , as a function of angles  $\Theta$  (for the coronal plane) and  $\phi$  (for the transverse plane) was defined as.

$$\delta R(\theta, \phi) = \left| \vec{OM} - \vec{OR} \right| \quad (1)$$

where distances  $|\vec{OM}|$  and  $|\vec{OR}|$  are quantified from the center of mass of the DL reference contour (O) in the direction determined by  $\Theta$  and  $\phi$  [blue arrow in Fig. 1(b)] to the point of intersection with the surface of the manual (M) and DL reference contour (R), respectively [Fig. 1(b)]. By sampling the coronal and transverse angles every  $3^\circ$ , contour deviation maps of  $60 \times 120$  pixels were created for each patient [Fig. 1(c)].

## 2.C. Cohort refinement

To ensure consistency in the input data, we curated the patient data. First, the manual contour of each patient was visually verified to ensure the contour only contained the prostate gland, that is, to ensure a consistent contouring protocol in the dataset. Second, the quality of the DL contour was visually verified. Patients with seminal vesicles included in the manual contour or where the DL contouring failed were removed from analysis [see Fig. 2]. Third, the contour deviation map of each patient was visually verified to ensure intuitive deviation, (i.e., in the order of mm), in all directions. Patients with local anomalies,  $\delta R(\theta, \phi)$ , often in the order of 10 cm, were removed. The cause of this problem was traced back to holes in the triangulation of contours into surfaces (see Fig. S1 in the supplementary materials). Visual verification was performed by AJ and TSM, to ensure consistency of contours rather than clinical correctness.

In addition to cohort refinement, we assessed the similarity between the DL contour and manual contour to identify if systematic differences exist across the cohort. We assessed this and report the following metrics: histogram of Dice similarity coefficient, scatter plot of DL vs manual contour volumes, three scatter plots of the DL vs manual contour x, y, and z center-of-mass coordinates, respectively, and boxplots of  $\delta R$  values in each region from Fig. 3.

## 2.D. Pixel-wise survival analysis

We assumed that each pixel in the contour deviation map referred to a consistent anatomical location.<sup>16</sup> For each pixel in the contour deviation maps, time to biochemical recurrence was modeled using a multivariable Cox proportional hazard model (CPHM) accounting for contour deviation, patient age, Gleason score, and manual CTV volume. Note that the last three variables were constant for all pixels for a single patient. This methodology has been developed by Green, *et al.*<sup>14</sup> and in their publication details the development of the statistical technique. It is important to note that Green's implementation is validated with respect to the "Survival" toolkit in R<sup>17</sup>; the R code for a single voxel is exemplified by Fig. S2 in the supplementary materials. By assembling the hazard ratios (HRs) of the 7200 CPHM in the  $60 \times 120$  grid, HR maps for each variable were created. Regions of significance were found using the cluster-based permutation test on the z-statistics<sup>18</sup> ( $10^4$  permutations). Ranges of HRs on each interface and within each prostate region are reported by extracting them using eight rectangular regions of interest (see Fig. 3).

## 3. RESULTS

### 3.A. Cohort refinement

From the original 247 patients, seven were removed as their manual contour included the seminal vesicles, one because of DL failure, and nine were removed due to local

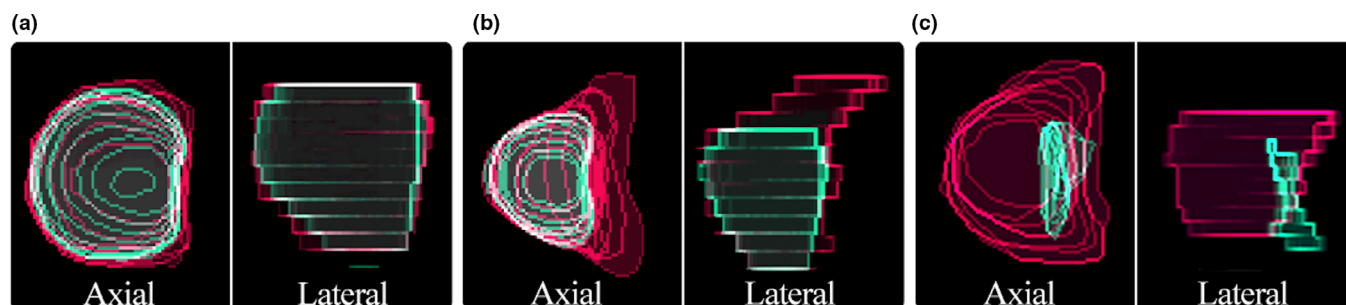


FIG. 2. Examples of the axial and lateral projections of the deep learning (DL) auto-contour (blue) and the manual contour (red) created for each patient. These images were used to identify and remove patients from the analysis if their seminal vesicles are included in the manual contour, or if the DL contouring has failed. (a) A patient kept in the analysis with only their prostate in the manual contour. (b) A patient removed from the analysis with their prostate and seminal vesicles in the manual contour. (c) The only patient where DL contouring failed (due to artifacts caused by bilateral metal hips) was also removed. [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

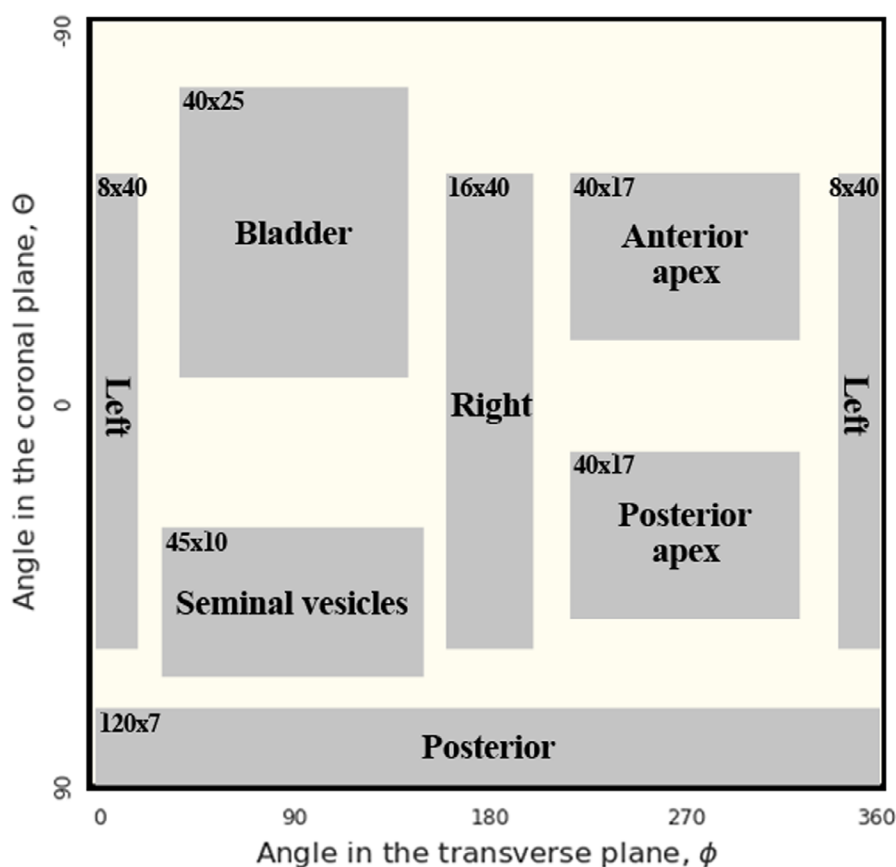


FIG. 3. The shaded rectangular kernels shown here represent regions and interfaces of interest on the prostate, which were used to calculate the range of Hazard ratios for the confounding variables. The dimensions (x,y) of each kernel is denoted in its upper left and is measured in pixels. [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

anomalies in surface reconstruction, leaving 232 patients for analysis (Table I, Fig. 2).

When assessing the similarity between the DL contour and manual contour, we found that the Dice similarity coefficient for the contours was on average 0.81 (range 0.41–0.92) (see supplementary materials, Fig. S3). A scatter plot of the DL vs manual contour volume revealed that the manual contour has a consistently larger volume on average than the DL contour (see supplementary materials, Fig. S4). Boxplots of

$\delta R$  values in each region from Fig. 3, revealed that median  $\delta R$  values were positive and different across all regions; this confirmed our observation from Fig. S2, but also identified the posterior apex region as that with greatest systematic over-contouring (see supplementary materials, Fig. S5). Scatter plots of the x, y, and z center-of-mass coordinates for the DL and manual contour revealed no systematic shifts in coordinates across the cohort (see supplementary materials, Fig. S6).

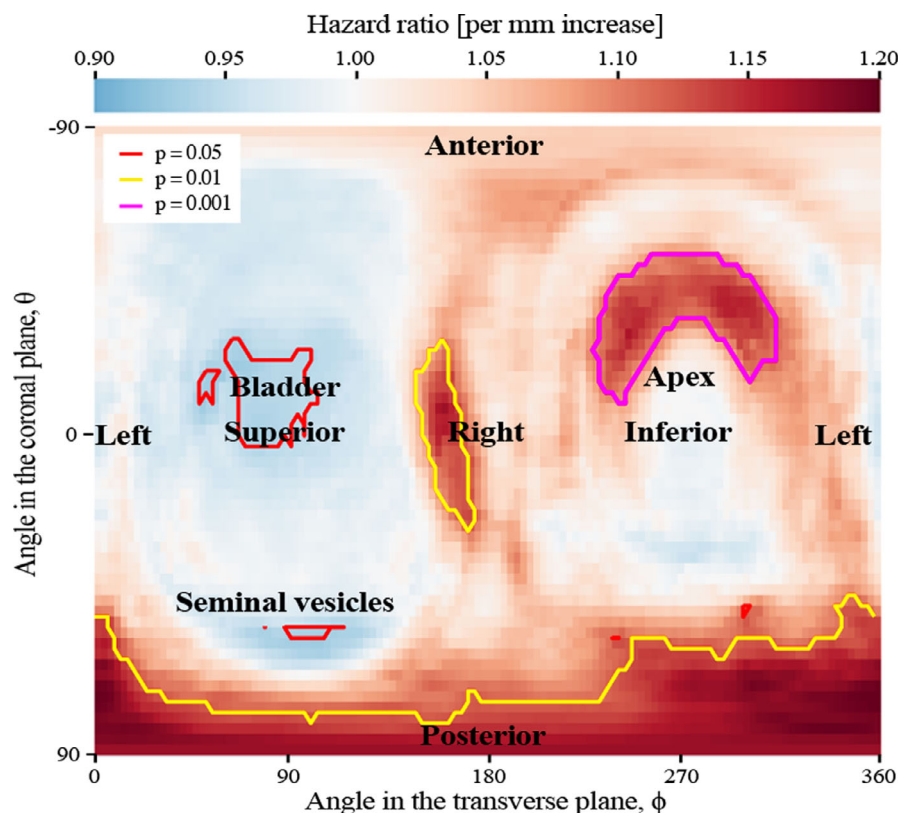


FIG. 4. Hazard ratio (HR) map for contouring deviation. Regions of significance are contoured. The map suggests that contouring larger volumes in the left region, bladder, and seminal vesicle interfaces could lead to better biochemical control. This HR map is modulated by all variables included in the analysis. [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

### 3.B. Pixel-wise survival analysis

Figures 4 and 5 show the HR maps for the considered variables in the CPHM. The effect of contouring deviation,  $\delta R(\theta, \phi)$ , against biochemical recurrence, controlling for Gleason score, age, and manual CTV volume is shown in Fig. 5. The contours encapsulate varying statistically significant regions of  $HR < 1$  and  $HR > 1$ . This result shows that per mm increase in the manual contour relative to the DL reference in the prostate–bladder and prostate–seminal vesicle interfaces, *reduces* the risk of biochemical recurrence by 4–8% ( $P < 0.05$ ). Conversely, per mm increase in the manual contour relative to the DL reference in the prostate’s right, anterior and posterior regions, *increases* the risk of biochemical recurrence by 8–24% ( $P < 0.01$ ). Figure 5 shows HR maps for the controlled confounding variables. Patient age showed a significant relationship in the bladder and seminal vesicles interfaces, and the posterior and apical regions ( $P < 0.05$ ) as shown in Fig. 5(a). This implies an interaction between contour variation, age, and biochemical recurrence. Manual CTV volume shows a significant relationship with biochemical recurrence throughout the prostate’s superior ( $P < 0.001$ ), as shown in Fig. 5(b). When Gleason scores 7 and 8 are compared to Gleason score 9–10, as shown in Figs 5(c)–5(d), all values in the HR maps are less than unity, as expected. Throughout Fig. 5(c) all values are statistically significant ( $P < 0.05$ ), however, only the region contoured is

statistically significant for Fig. 5(d). This shows that patients with a Gleason score lower than 9–10 have a reduced relative risk of biochemical recurrence, which is an intuitive result, with some interaction with contour variation. Table II displays the range of HRs on each interface and within each prostate region, extracted using the rectangular regions defined in Fig. 3.

### 4. DISCUSSION

In this study we proposed a novel methodology to analyze the effect of contouring uncertainty on clinical outcome for a large cohort of patients. Here, we used our method on the prostate site and demonstrated how to identify regions where contour deviation and recurrence are correlated, by measuring the deviation of each patient’s clinical contour relative to a highly consistent reference contour and applying pixel-wise CPHM, followed by permutation testing. It is important to notice that this methodology can be used to explore relationships to other outcomes as well.

After applying our proposed methodology to a cohort of 232 prostate cancer patients, we found regions where contour deviations were correlated with biochemical recurrence. In detail, we observed:

- A per mm increase in the manual contour relative to the DL reference in the prostate’s bladder and seminal

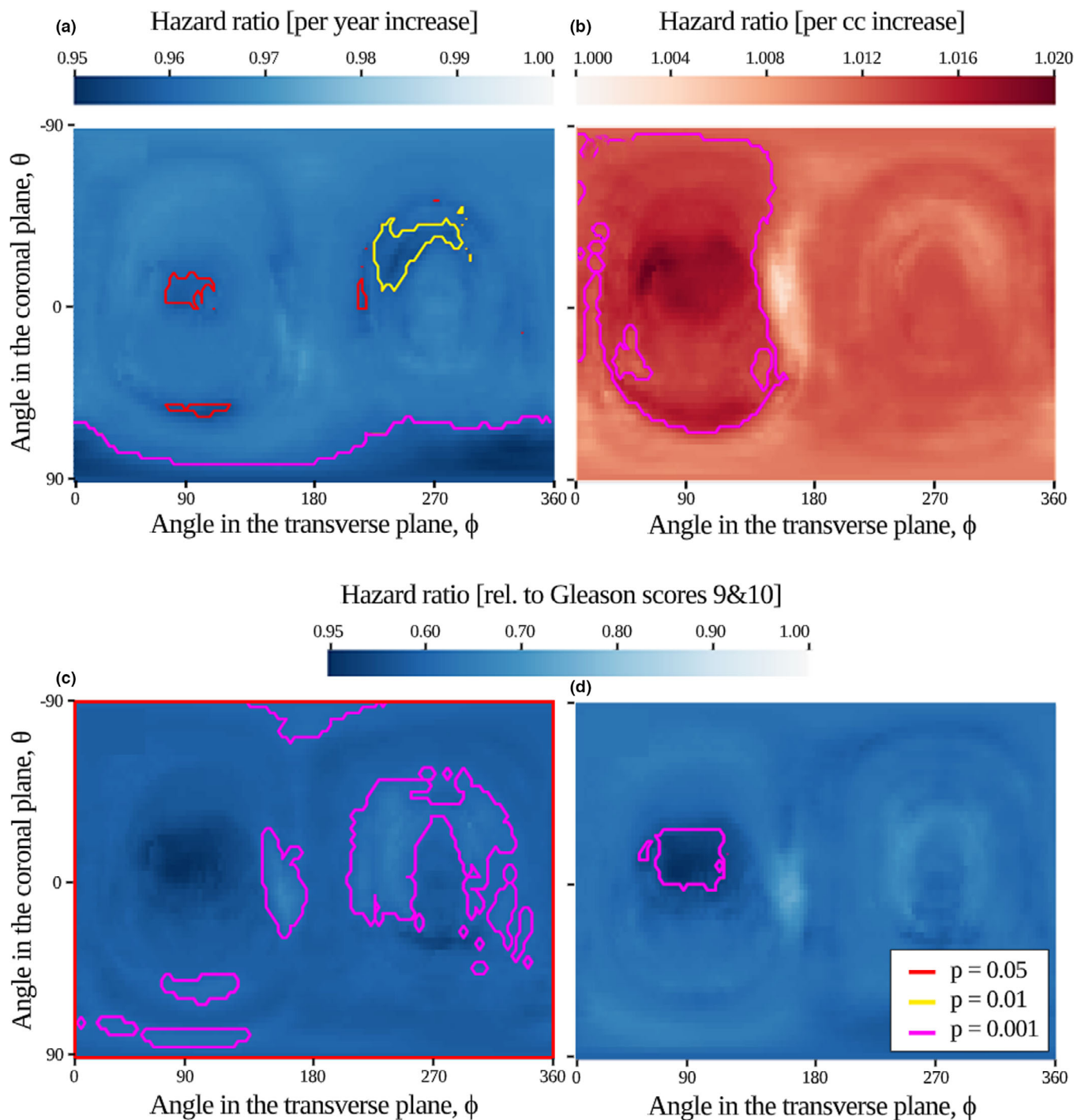


FIG. 5. The Hazard ratio (HR) maps of the risk of recurrence for the confounding variables modulated by the spatially varying contour deviation. (a) Age per year increase, (b) Manual CTV volume per  $\text{cm}^3$  increase, (c) Gleason score 7 relative to 9–10, and (d) Gleason score 8 relative to 9–10. Notice that every HR map is based on modeling all variables included in the analysis. [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

vesicle interfaces is related to a *reduction* of the risk of biochemical recurrence by 4–8% ( $P < 0.05$ ). This can be interpreted as having larger manual contours in these regions is associated with improved control.

- A per mm increase in the manual contour relative to the DL reference in the prostate's right, anterior, and posterior regions, is related to an *increase* in the risk of biochemical recurrence by 8–24% ( $P < 0.01$ ). This

means that larger manual contours in these regions are associated with poorer control.

To the best of our knowledge, this is the first investigation showing a direct effect of contouring variation on biochemical recurrence following prostate radiotherapy. For the pixel-wise CPHM analysis, we included the Gleason score at diagnosis, patient age at the start of treatment, manual CTV

TABLE II. Range of Hazard ratio (HR) for each covariate at different regions across the prostate, extracted from the HR maps shown in Figs. 4 and 5, using the regions of interest defined in Fig. 3.

Region	Prostate Volume (per cc)	Age (per year)	Gleason score 7 (rel. 9–10)	Gleason score 8 (rel. 9–10)	Contour variation (per mm)
Posterior	1.007–1.011	0.954–0.968	0.546–0.592	0.514–0.626	1.002–1.236
Right	1.003–1.011	0.963–0.974	0.544–0.645	0.538–0.717	0.994–1.201
Anterior Apex	1.007–1.011	0.958–0.986	0.552–0.636	0.542–0.666	0.987–1.190
Posterior Apex	1.009–1.011	0.962–0.970	0.504–0.613	0.527–0.649	0.948–1.151
Seminal vesicles	1.008–1.013	0.958–0.969	0.547–0.600	0.560–0.608	0.916–1.078
Left	1.008–1.013	0.962–0.969	0.548–0.591	0.550–0.610	0.952–1.188
Bladder	1.004–1.016	0.962–0.971	0.476–0.598	0.454–0.628	0.915–1.088

volume, and contour deviation. However, the presence of the counter-intuitive observation made, that larger tumor coverage may increase the risk of recurrence, potentially invalidates the more logical one. This clearly points to the need to correct for additional confounding variables. Confounding variables missing from the analysis may influence the interpretation and significance of our findings. Such additional variables include the PSA level on diagnosis, spatial variation of the planned target dose, and the patient's rectum volume upon planning, which are all known to affect the risk of biochemical recurrence.<sup>19–21</sup> This analysis should therefore be repeated in a cohort, preferably larger, where these covariates are available. In addition, the pixel-wise CPHM needs to be internally and externally validated to ensure the observed HRs will generalize well to new patient data, from a variety of different populations. These additional steps will help turn our observations into conclusions, and ultimately translate these results into the clinic.

HR maps of the other covariates also had significant regions (Fig. 5), indicating that contour deviation is not the only variable affecting survival, which is not surprising. Interpretation of the HR maps is less intuitive than the contour deviation maps. The HR maps of the confounding variables will be approximately constant when there is no interaction between contouring deviation and the confounding variable. This is clearly the case for age. However, the HR maps of prostate volume shows different HR for different regions, which is logical because delineation deviation and volume change have the same effect — a motion of the delineated contours. The most interesting interaction is for Gleason scores (7 and 8 relative to 9–10) which both suggest an increase in risk of delineation deviation for higher Gleason scores. We also checked whether there was a correlation between Gleason score and manual CTV volume, however, this was not observed (see supplementary materials, Fig. S7).

Classical contouring variation studies require the contours of multiple observers, often oncologists: an expensive resource. In a recent review of IOV, the number of patients range between 1 and 26 for prostate cancer cases.<sup>7</sup> Such a limited number of patients makes meaningful analysis of clinical outcomes impossible. In our novel approach, we use a DL auto-contour as a reference to quantify contour deviation for more than 200 patients. Despite the clinical correctness of DL auto-contours being debatable, it provides a consistent

baseline to compare the manual (clinically used) contour. As DL only uses the image dataset as input, clinical circumstances cannot influence the results of the DL model, meaning that detected deviations are not confounded by such clinical variables. However, it is important to note that the DL model will consistently reproduce any bias present in the training data, which may be the reason of the DL contours are smaller than the manual on average. Recent developments are starting to handle CTV definition automatically<sup>13</sup> which, if implemented clinically, may improve consistency for future patients. For our work, we visually inspected the contours and removed failed DL contours, therefore minimizing the impact of image quality on DL contouring. Retraining the DL model is beyond our reach as we used commercial tools. As discussed in the results section, we observed that the manual contour is consistently larger on average than the DL contour (see supplementary materials, Fig. S4), and also found that the manual contour is systematically larger in all regions of interest from Fig. 3 (see supplementary materials, Fig. S5). This consistent difference does not affect the estimates resulting from our analysis: the estimate for  $\delta R$  for each Cox model (at each pixel) is determined from the spread of  $\delta R$ , rather than its absolute value. Having a reference structure that is consistently smaller will only affect the intercept term at each pixel from the Cox model, which is not of interest here.

For the current analysis, contour deviation maps were extracted using a spherical coordinate system centered on the DL reference contour. This assumes that the same spherical direction (i.e.,  $\Theta$  and  $\varphi$ ) will capture the same anatomical prostate region for all patients. This assumption is likely valid for convex organs which shape and orientation is similar among patients, such as the prostate. A similar assumption was made by Witte *et al.*<sup>19</sup> for image-based data mining. However, further work on improving interpatient mapping of structures to a common frame of reference could further improve the results and allow this methodology to be extended to nonconvex structures.

For our current analysis, we explored the magnitude of  $\delta R$  and its relation to biochemical recurrence. As the contours of the prostate shape the region of high dose, we assumed that radial variations on these contours would indirectly relate to treatment failure. Our method could be adapted if directional variability is of interest, where instead of looking at the

magnitude of  $\delta R$ , we could look at its vector components separately (e.g., right–left, anterior–posterior, or inferior–superior components).

For the studied patient cohort, our method produced both expected and unexpected results when relating contour deviations to biochemical recurrence, for example larger contours around the seminal vesicles predict better control, but larger contours around the posterior region predict worse control. As such, these results should be interpreted with caution and extra analysis on an external and larger dataset should be performed before translating them into clinical practice. These results also highlight the need to include deviation cross-correlations introduced by observer contouring “styles” into the analysis. The effect of these cross-correlations could be minimized by including a large number of observers to reduce the bias on contouring style. Alternatively, the observer identification could be added as a confounding categorical variable in the CPHM. In our case, individual observers could not be identified from the retrospective data and therefore, this effect could not be accounted for in our analysis. Other deviation correlations may be introduced by the observer’s level of expertise and the individual interpretation of the local contouring protocol factors that could be included in future studies.<sup>22–24</sup>

From the HR map for contouring deviation, as presented here, regions are identified where clinical contours could be altered in order to limit a patient’s risk of biochemical recurrence following prostate radiotherapy. Again, we highlight that it is important to notice our methodology can be used to explore relationships between contouring and any outcome. As a result, competing risk models could be built on top of our methodology and used to highlight the regions where contouring deviation should be reduced to find the optimal therapeutic ratio. Thus, our developed methodology could be used to better define protocols for contouring, and potentially improve patient outcomes, once applied to a dataset where the aforementioned additional covariates are available, and an internal and external validation has been conducted. Furthermore, the methodology proposed here could be adapted to other radiotherapy treatment sites.

The primary goal of this manuscript was to introduce a methodology to explore the correlations between contour variation and outcome. The general framework followed by our method, that is, image-based data-mining or voxel-based analysis, been used in neuroimaging for over a decade<sup>25,26</sup> and it has been successfully used to explore radiotherapy dose and outcome in several sites.<sup>19,27–29</sup> We refer the interested reader to the recent article by Palma *et al.* where a “Cookbook” dedicated to this method for use in radiation oncology is presented<sup>30</sup> and a critical editorial on the importance of correct statistical analysis.<sup>31</sup>

## 5. CONCLUSION

We have proposed a novel method to analyze the effect of contouring variation on clinical outcome for a large cohort of patients using deviations to a consistent virtual

observer. We exemplify our methodology on a cohort of prostate cancer patients, and use time to biochemical recurrence following treatment as our outcome. Regions were identified in which contour deviations of the prostate relate to biochemical recurrence, with some expected and unexpected results (produced both expected and unexpected results). After including relevant covariates and validating results with an external dataset, results from this methodology could inform contouring protocols based on actual patient outcomes.

## ACKNOWLEDGMENTS

M.V.H. was supported by NIHR Manchester Biomedical Research Centre. Prostate Cancer UK (RIA15-ST2-031) and Cancer Research UK via funding to the Cancer Research Manchester Centre (C147/A25254) supported this work. M.A. is also supported by Cancer Research UK (C8225/A21133).

## CONFLICT OF INTEREST

The author(s) declare no conflict of interest.

## DATA AVAILABILITY STATEMENT

The data that support the findings of this study are available from the corresponding author upon reasonable request.

\*These authors contributed equally to this work

<sup>a)</sup>Author to whom correspondence should be addressed. Electronic mail: eliana.vasquezosorio@manchester.ac.uk; Telephone: (+44) 161 918 7480.

## REFERENCES

1. Bray F, Ferlay J, Soerjomataram I, Siegel RL, Torre LA, Jemal A. Global cancer statistics 2018: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries. *CA Cancer J Clin.* 2018;68:394–424.
2. National Cancer Registration & Analysis Service, Cancer Research UK. National Cancer Registration and Analysis Service Short Report: Chemotherapy, Radiotherapy and Surgical Tumour Resections in England: 2013–2014; 2018; [http://www.ncin.org.uk/cancer\\_type\\_and\\_topic\\_specific\\_work/topic\\_specific\\_work/main\\_cancer\\_treatments](http://www.ncin.org.uk/cancer_type_and_topic_specific_work/topic_specific_work/main_cancer_treatments).
3. Pascale M, Azinwi CN, Marongiu B, Pesce G, Stoffel F, Roggero E. The outcome of prostate cancer patients treated with curative intent strongly depends on survival after metastatic progression. *BMC Cancer.* 2017;17:651.
4. Nyholm T, Jonsson J, Söderström K, *et al.* Variability in prostate and seminal vesicle delineations defined on magnetic resonance images, a multi-observer, -center and -sequence study. *Radiat Oncol.* 2013;8:1–12.
5. Fiorino C, Reni M, Bolognesi A, Cattaneo GM, Calandrino R. Intra- and inter-observer variability in contouring prostate and seminal vesicles: Implications for conformal treatment planning. *Radiother Oncol.* 1998;47:285–292.
6. Villeirs GM, Van Vaerenbergh K, Vakaet L, *et al.* Interobserver delineation variation using CT versus combined CT + MRI in intensity-modulated radiotherapy for prostate cancer. *Strahlentherapie und Onkol.* 2005;181:424–430.
7. Vinod SK, Jameson MG, Min M, Holloway LC. Uncertainties in volume delineation in radiation oncology: A systematic review and recommendations for future studies. *Radiother Oncol.* 2016;121:169–179.

8. Cardenas CE, Yang J, Anderson BM, Court LE, Brock KB. Advances in auto-segmentation. *Semin Radiat Oncol.* 2019;29:185–197.
9. McCarroll RE, Beadle BM, Balter PA, et al. Retrospective validation and clinical implementation of automated contouring of organs at risk in the head and neck: A step toward automated radiation treatment planning for low- And middle-income countries. *J Glob Oncol.* 2018;2018.
10. Zabel WJ, Conway JL, Gladwish A, et al. Clinical evaluation of deep learning and atlas-based auto-contouring of bladder and rectum for prostate radiation therapy. *Pract Radiat Oncol.* 2021;11:e80–e89.
11. Gooding MJ, Smith AJ, Tariq M, et al. Comparative evaluation of auto-contouring in clinical practice: A practical method using the Turing test. *Med Phys.* 2018;45:5105–5115.
12. Almeida G, Tavares JMRS. Deep learning in radiation oncology treatment planning for prostate cancer: A systematic review. *J Med Syst.* 2020;44:1–15.
13. Unkelbach J, Bortfeld T, Cardenas CE, et al. The role of computational methods for automating and improving clinical target, vol. definition. *Radiother Oncol.* 2020;153:15–25.
14. Green A, Vasquez Osorio E, Aznar MC, McWilliam A, van Herk M. Image based data mining using per-voxel cox regression. *Front Oncol.* 2020;10:1178.
15. Lorensen WE, Cline HE. Marching cubes: A high resolution 3D surface construction algorithm. In: Proceedings of the 14th Annual Conference on Computer Graphics and Interactive Techniques, SIGGRAPH 1987. Vol 21. Association for Computing Machinery. Inc; 1987:163–169.
16. Remeijer P, Rasch C, Lebesque JV, Van Herk M. A general methodology for three-dimensional analysis of variation in target volume delineation. *Med Phys.* 1999;26:931–940.
17. Therneau TM, Grambsch PM. *Modeling Survival Data: Extending the Cox Model.* Berlin: Springer; 2000.
18. Bullmore ET, Suckling J, Overmeyer S, Rabe-Hesketh S, Taylor E, Brammer MJ. Global, voxel, and cluster tests, by theory and permutation, for a difference between two groups of structural MR images of the brain. *IEEE Trans Med Imaging.* 1999;18:32–42.
19. Witte MG, Heemsbergen WD, Bohoslavsky R, et al. Relating dose outside the prostate with freedom from failure in the dutch trial 68 Gy vs. 78 Gy. *Int J Radiat Oncol Biol Phys.* 2010;77:131–138.
20. Ang M, Rajcic B, Foreman D, Moretti K, O'Callaghan ME. Men presenting with prostate-specific antigen (PSA) values of over 100 ng/mL. *BJU Int.* 2016;117:68–75.
21. Heemsbergen WD, Hoogeman MS, Witte MG, Peeters STH, Incrocci L, Lebesque JV. Increased risk of biochemical and clinical failure for prostate patients with a large rectum at radiotherapy planning: Results from the dutch trial of 68 GY versus 78 Gy. *Int J Radiat Oncol Biol Phys.* 2007;67:1418–1424.
22. Rasch C, Barillot I, Remeijer P, Touw A, Van Herk M, Lebesque JV. Definition of the prostate in CT and MRI: A multi-observer study. *Int J Radiat Oncol Biol Phys.* 1999;43:57–66.
23. Altorjai G, Fotina I, Lütgendorf-Caucig C, et al. Cone-beam CT-based delineation of stereotactic lung targets: The influence of image modality and target size on interobserver variability. *Int J Radiat Oncol Biol Phys.* 2012;82.
24. Nicholls L, Gorayski P, Poulsen M, et al. Maintaining prostate contouring consistency following an educational intervention. *J Med Radiat Sci.* 2016;63:155–160.
25. Ashburner J, Friston KJ. Voxel-based morphometry - the methods. *NeuroImage.* 2000;11:805–821.
26. Whitwell JL. Voxel-based morphometry: An automated technique for assessing structural changes in the brain. *J Neurosci.* 2009;29:9661–9664.
27. Palma G, Monti S, D'Avino V, et al. A voxel-based approach to explore local dose differences associated with radiation-induced lung damage. *Int J Radiat Oncol Biol Phys.* 2016;96:127–133.
28. McWilliam A, Kennedy J, Hodgson C, Vasquez Osorio E, Faivre-Finn C, van Herk M. Radiation dose to heart base linked with poorer survival in lung cancer patients. *Eur J Cancer.* 2017;85:106–113.
29. Guo Y, Jiang W, Lakshminarayanan P, et al. Spatial radiation dose influence on xerostomia recovery and its comparison to acute incidence in patients with head and neck cancer. *Adv Radiat Oncol.* 2020;5:221–230.
30. Palma G, Monti S, Cella L. Voxel-based analysis in radiation oncology: A methodological cookbook. *Phys Medica.* 2020;69:192–204.

31. Shortall J, Palma G, Mistry H, et al. Flogging a dead salmon? Reduced dose posterior to prostate correlates with increased PSA progression in voxel-based analysis of 3 randomised phase 3 trials – Marcello et al. *Int J Radiat Oncol Biol Phys.* 2021.

## SUPPORTING INFORMATION

Additional supporting information may be found online in the Supporting Information section at the end of the article.

**Fig S1.** Each patient's deviation map was inspected to identify and remove patients with local contouring anomalies induced by triangulating each contour into a surface. (a) The deviation map of a patient, kept in the analysis. (b) The deviation map of a patient with a local anomaly of order 10 cm, removed from the analysis.

**Fig S2.** Example code snippet demonstrating the equivalent R code to perform pixel-wise cox proportional hazard model toolbox using the “survival” package.

**Fig S3.** Distribution of Dice similarity coefficient across the entire patient cohort. Dice similarity coefficient measures the spatial overlap between the volume of the manual contours and the DL algorithm contours. The majority of Dice coefficients were greater than 0.8, suggesting high similarity between the DL and manual contours.

**Fig S4.** Scatter plot showing the volume of the manual contour against the volume of the deep learning algorithm. We observe a tendency of the datapoints to be below the line of unity (dashed line), meaning that the DL CTVs were smaller than the manual CTVs. The  $y = 0.659x + 5.333$  line is the fitted linear model ( $R^2 = 0.867$ ).

**Fig S5.** The distribution of  $\partial R$  pixel values for regions and interfaces of interest on the prostate, defined in Figure 3 in the main manuscript, indicating systematic differences across all regions. All median  $\partial R$  values are  $\geq 0$  cm agreeing with Figure S4 showing that the DL segmentations are systematically smaller than the manual contours. One-way ANOVA test showed significant systematic differences over all regions of interest ( $F\text{-score} = 1.08e4$ ,  $P < 10e-4$ ). This systematic variation is expected as contour variation has been reported to be different in different regions.

**Fig S6.** Scatter plots of the centre-of-mass of the manual contours (y-axis) against the DL contours (x-axis) in the x plane (right-left), y-plane and z-plane shown in (a), (b) and (c) respectively. The null hypothesis that distributions of x, y and z centre-of-mass coordinates between manual and DL contours have equal means, has been tested using a two-sample t-test. We found no significant evidence to suggest the means centre-of-mass coordinates are different in x ( $P = 0.64$ ), y ( $P = 0.64$ ), and z ( $P = 0.98$ ).

**Fig S7.** Distribution of the clinical target volume (CTV) for patients with different Gleason scores. We performed a one-way analysis of variation (ANOVA) to test if there were differences for the CTV patients with different Gleason score groups ( $F\text{-score} = 1.9196$ ,  $P = 0.127$ ). The insignificant score suggests that the Gleason score of the patient has an insignificant interaction with their defined CTV.