

Article

# Cross-domain data augmentation for deep-learning-based male pelvic organ segmentation in cone beam CT

Jean Léger <sup>1,†,\*</sup> , Elliott Brion <sup>1,†</sup> , Paul Desbordes <sup>1</sup> , Christophe De Vleeschouwer <sup>1</sup> , John A. Lee <sup>1,2</sup>  and Benoit Macq <sup>1,\*</sup> 

<sup>1</sup> ICTEAM, UCLouvain, 1348 Louvain-la-Neuve, Walloon Brabant, Belgium

<sup>2</sup> IREC/MIRO, UCLouvain, 1200 Woluwe-Saint-Lambert, Brussels, Belgium

\* Correspondence: jean.leger@uclouvain.be, benoit.macq@uclouvain.be

† These authors contributed equally to this work.

Version February 5, 2020 submitted to Appl. Sci.

**Abstract:** For prostate cancer patients, large organ deformations occurring between radiotherapy treatment sessions create uncertainty about the doses delivered to the tumor and surrounding healthy organs. Segmenting those regions on cone beam CT (CBCT) scans acquired on treatment day would reduce such uncertainties. In this work, a 3D U-net deep-learning architecture was trained to segment the bladder, rectum, and prostate on CBCT scans. Due to the scarcity of contoured CBCT scans, the training set was augmented with CT scans already contoured in the current clinical workflow. Our network was then tested on 63 CBCT scans. The Dice similarity coefficient (DSC) increases significantly with the number of CBCT and CT scans in the training set, reaching  $0.874 \pm 0.096$ ,  $0.814 \pm 0.055$ , and  $0.758 \pm 0.101$  for the bladder, rectum, and prostate respectively. This is about 10% better than conventional approaches based on deformable image registration between planning CT and treatment CBCT scans, except for the prostate. Interestingly, adding 74 CT scans to the CBCT training set allowed to maintain high DSCs, while halving the number of CBCT scans. Hence, our work shows that although CBCT scans include artifacts, cross-domain augmentation of the training set is effective and can rely on large datasets available for planning CT scans.

**Keywords:** segmentation; deep learning; deformable image registration; cone beam CT; pelvis; prostate cancer; radiotherapy; CNN; U-net

## 1. Introduction

Fractionated external beam radiotherapy (EBRT) cancer treatment relies on two steps. In the treatment planning phase, clinicians delineate the tumor and surrounding healthy organs' volumes on a computed tomography (CT) scan and compute the dose distribution. In the treatment delivery phase, the patient is aligned with a specific treatment planning position and the dose fraction is delivered. Patient positioning relies on a daily cone beam computed tomography (CBCT) scan acquired in the treatment position before each treatment fraction is delivered.

CT and CBCT are both based on X-ray propagation through the patient's body. However, the CBCT scans are of lower quality than CT scans due to different types of artifact, including noise, beam hardening, and scattering, as shown in Figure 1. In particular, scattering is an important limitation that could rule out the use of CBCT for radiotherapy treatment planning [1]. However, CBCT scans are currently used in order to detect daily variations in patient anatomy, which are particularly large in the pelvic region due to physiological function (e.g., bladder and rectal filling and voiding). Detecting

31 such variations is important since they can impair treatment dose conformity, which means delivering  
32 too large a dose to the healthy organs (e.g., the bladder and rectum in the case of prostate cancer)  
33 and too low a dose to the clinical target volume (which simply corresponds to the prostate itself for  
34 a significant proportion of patients) [2]. To improve treatment dose conformity in the pelvic region  
35 further, proposals have been made to change treatment plan delivery as a function of time based on  
36 observed anatomic variations [3,4].

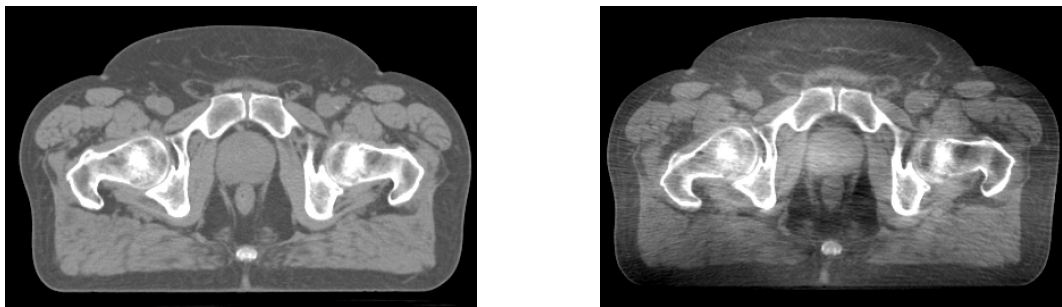
37  
38 However, a step towards better adaptive radiotherapy would require automatic segmentation  
39 of the pelvic organs on daily CBCT scans in order measure the anatomical variations accurately.  
40 Automating this segmentation is necessary to be able to integrate it in the clinical workflow, as  
41 delineating the organs manually on daily scans is excessively time-consuming. Measuring anatomical  
42 variations is particularly important in proton therapy because the proton dose distribution is highly  
43 sensitive to changes in patient geometry [5,6].

44  
45 Currently, organ segmentation is classically performed by deformable image registration (DIR)  
46 algorithms between the planning CT and daily CBCT scans [7,8]. These algorithms include such  
47 clinical software packages as MIM [9] and RayStation [10]. Although the results are better than those  
48 of rigid registration, these intensity-based DIR algorithms fail in the presence of large deformations  
49 between the registered scans, as is the case in the pelvic region [11,12]. Zambrano et al. [11] and  
50 Thor et al. [12] implemented a featurelet-based algorithm [13] and the demons DIR algorithm [14],  
51 respectively. As a result, more complex DIR approaches, such as a B-spline DIR algorithm relying  
52 on mutual information, have been proposed [15]. This last approach implements a 6-pass DIR with  
53 progressively finer resolution and, after visual inspection, an optional final pass using a narrow  
54 region around the region of interest. Another approach uses a DIR framework where a locally rigid  
55 deformation is enforced for bone and/or the prostate, while surrounding tissue is still allowed to  
56 deform elastically [16]. Alternatively, statistical shape models can capture shape variations and have  
57 also been considered for bladder segmentation on CBCT scans [17,18]. However, those methods  
58 require the definition of landmarks or meshes. Moreover, several delineated CBCT scans must be  
59 available to build a patient-specific shape model. That thwarts the application of such methods at  
60 the start of treatment. So, none of these methods accomplishes the challenging task of pelvic organ  
61 segmentation on CBCT scans. In parallel, recent advances in computing capabilities, the availability of  
62 representative datasets, and the great versatility of deep-learning (DL) approaches have enabled DL  
63 algorithms to achieve impressive segmentation performance. Unlike the aforementioned techniques,  
64 DL algorithms are supposed to be robust to variations in shape and appearance if those variations  
65 are captured in the training database and do not require landmark definition. DL algorithms have  
66 already been used successfully to segment pelvic organs on CT scans [19,20]. The 3D U-net fully  
67 convolutional neural network [21] has been used to segment female pelvic organs on CBCT scans  
68 [22,23]. Concurrently, we showed that adding annotated CT scans to the training set improved  
69 bladder segmentation on CBCT scans [24]. This approach was motivated by the scarcity of annotated  
70 CBCT scans, compared with annotated CT scans, and the fact that CBCT scans can be roughly  
71 considered to be noisy, distorted CT scans from a segmentation perspective, hence sharing shape and  
72 contextual information with the CT scans. The current paper extends our previous conference paper  
73 [24] in that it considers additional male pelvic organs (the rectum and prostate), and presents more  
74 comparative results (including the morphons deformable registration algorithm). It also involves  
75 data from an additional hospital and provides a more detailed discussion. Segmentation of male  
76 pelvic organs (bladder, rectum, prostate, and seminal vesicles) on CBCT and CT scans using a DL  
77 approach was the subject of a recent paper [25]. These authors' contribution consists mainly of the use  
78 of artificially-generated pseudo CBCT scans in the training set along with a high segmentation quality.  
79 Our approach adds training on real CBCT scans and provides a new and larger test set as well as more

80 extensive comparison with clinically-used registration tools.

81

82 The main contributions of this work are to provide (i) a DL-based segmentation method for male  
83 pelvic organs on CBCT scans and (ii) a detailed comparison of state-of-the-art segmentation tools in  
84 order to guide the choice of method in clinical practice. The impacts of the number of training scans  
85 and addition of CT scans to the training database were studied in order to provide detailed information  
86 on the amount of annotations required for use in clinical practice.



(a) Slice of a CT scan.

(b) Slice of a CBCT scan.

**Figure 1.** Comparison of CT and CBCT scans.

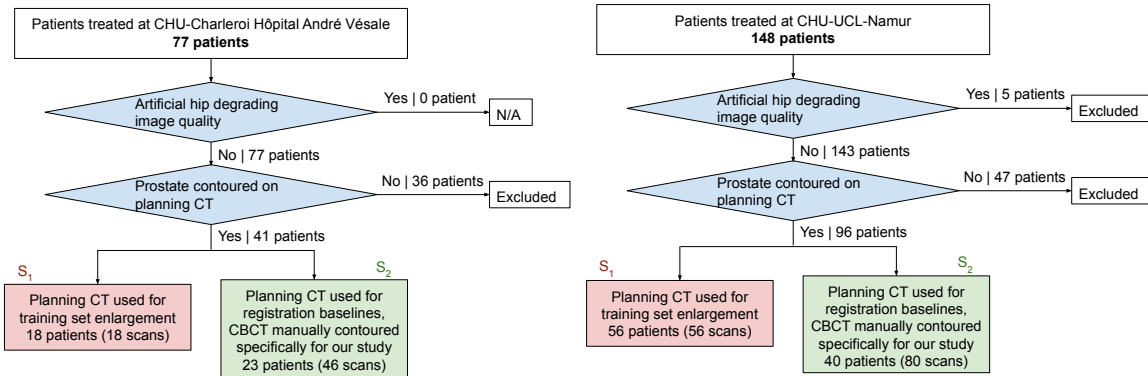
## 87 2. Materials and Methods

### 88 2.1. Data and preprocessing

89 Our data consist of (i) a set  $S_1$  of 74 patients for whom we have delineated CT scans and (ii) a  
90 set  $S_2$  of 63 patients (different from the 74 patients mentioned above) for whom we have delineated  
91 planning CT scans and delineated daily CBCT scans. The contours of the bladder, rectum, and prostate  
92 were delineated on the CT scans during the clinical workflow. The contours on the CBCT scans were  
93 delineated by a trained expert specifically for this study. Within set  $S_1$ , 18 and 56 patients underwent  
94 EBRT for prostate cancer at two teaching hospitals, CHU-Charleroi Hôpital André Vésale and  
95 CHU-UCL-Namur, respectively. Within set  $S_2$ , 23 and 40 patients underwent EBRT for prostate cancer  
96 at CHU-Charleroi Hôpital André Vésale (CBCT scans acquired with a Varian TrueBeam STx version  
97 1.5) and CHU-UCL-Namur (CBCT scans acquired with a Varian OBI cone beam CT), respectively.  
98 The use of these retrospective, anonymized data for this study has been approved by each hospital's  
99 ethics committee (dates of approval: May 24, 2017 for CHU-Charleroi Hôpital André Vésale and May  
100 12, 2017 for CHU-UCL-Namur). In order to ensure data uniformity across the entire dataset, all the  
101 3D CT and CBCT scans (as well as the 3D binary masks representing the manual segmentations)  
102 were re-sampled on a 1.2x1.2x1.5 mm regular grid. All re-sampled image volumes and binary mask  
103 volumes were cropped to volumes of 160x160x128 voxels containing the bladder, rectum, and prostate.

104

105 The case selection procedure is described in Figure 2. Patients with an artificial hip were excluded  
106 from this study because the presence of an artificial hip degrades the image too much for the organs  
107 to be segmented accurately by a human expert. Patients for whom the prostate was not contoured  
108 on the planning CT scan were also excluded. This corresponds to patients for whom the clinical  
109 target volume (CTV) differed from that of the prostate, either because this organ had been surgically  
110 removed or because the CTV included other areas in addition to the prostate. Note that it is common in  
111 radiotherapy to inject contrast media into the bladder. Different inter-subject levels of contrast product  
112 increased the variability of this organ's appearance, making its automatic contouring more challenging.  
113 Since our case selection procedure includes all patients regardless of the use of contrast media, our  
114 method is supposed to be robust to such variability.



**Figure 2.** Case selection from CHU-Charleroi Hôpital André Vésale and CHU-UCL-Namur.

## 2.2. Model architecture and learning strategy

The bladder, rectum, and prostate were segmented on CBCT scans using the 3D U-net fully convolutional neural network [21,26]. The 3D input goes through a contracting path to capture context and an expanding path to enable precise localization. In the last layer, a softmax is applied and the network outputs the probability of each voxel's belonging to the bladder, rectum, prostate, or none of these organs. The network architecture is shown in Figure 3. To obtain a binary mask for each organ, the most probable class label was assigned to each pixel individually. In practice, each organ was segmented as a single region of connected voxels. No disconnected region of the same organ was observed. The main advantage of fully convolutional neural networks is that they output predictions at the same resolution as the input. One output channel was considered per organ. The network was trained with the Dice loss. The Adam optimization algorithm was used with a learning rate of  $10^{-4}$ . The number of epochs was chosen such that convergence was reached. The hyper-parameters mentioned here are the same as in Brion et al. [24] and proved satisfactory on the data used in this work. For this reason and to keep data available for training and testing, no validation set was considered here. Training data were augmented online using rotation (between  $-5^\circ$  and  $5^\circ$  along each of the three axes), shift (between -5 and 5 pixels along each axes), and shear (reasonable values for the affine transformation matrix). The batch size was set to two, which is the maximum size affordable on our 11 Gb graphical processing units (GPU).

We performed 3-fold cross-validation with the 63 CBCT scans of set  $\mathcal{S}_2$ , where 2 folds ( $n_{CBCT} \leq 42$  volumes in total) were used as the training set and one fold (21 volumes) as the test set, as shown in Table 1. The number of training CBCT scans  $n_{CBCT}$  was varied such that  $n_{CBCT} \in \{0, 6, 10, 20, 30, 42\}$ . The training set was augmented with  $n_{CT}$  annotated CT scans from set  $\mathcal{S}_1$  such that  $n_{CT} \in \{0, 20, 74\}$ . The same CT scans were added to the training CBCT scans independently on the considered training folds. Hence, the training set contains  $n_{CBCT} + n_{CT}$  volumes in total. Note that the test set contains no CT scans (since our goal was to segment CBCT scans only). The source code is publicly available on [https://github.com/eliottbrion/pelvis\\_segmentation](https://github.com/eliottbrion/pelvis_segmentation).

**Table 1.** Three-fold cross-validation. To train the model, we used  $n_{CT}$  CT scans from  $\mathcal{S}_1$  and the  $n_{CBCT}$  first volumes from the CBCT folds labeled "train." To test the model, we used all 21 volumes from the CBCT fold labeled "test."

$\mathcal{S}_1$ (CT)	$\mathcal{S}_2$ (CBCT)		
	fold1	fold2	fold3
train	train	train	test
train	train	test	train
train	test	train	train

### 142 2.3. Validation and comparison baselines

143 In order to evaluate our contouring results, we used four metrics comparing the predicted and  
 144 manual segmentations. The Dice similarity coefficient (DSC) and the Jaccard index (JI) measure the  
 145 overlap between two binary masks, while the symmetric mean boundary distance (SMBD) assesses  
 146 the distance between the contours (i.e., the sets of points located at the boundary of the binary masks)  
 147 delineating those binary masks. We also computed the difference between the manual and predicted  
 148 volumes for all the organs considered. More specifically,

$$\text{DSC} = \frac{2|M \cap P|}{|M| + |P|}, \quad (1)$$

$$\text{JI} = \frac{|M \cap P|}{|M \cup P|}, \quad (2)$$

$$\text{SMBD} = \frac{\bar{D}(M, P) + \bar{D}(P, M)}{2}, \quad (3)$$

149 where  $M$  and  $P$  are the sets containing the matricial indices of the manual and predicted segmentation  
 150 3D binary masks, respectively;  $\bar{D}(M, P)$  is the mean of  $D(M, P)$  over the voxels of  $\Omega_M$ ; and  $D(M, P) =$   
 151  $\{\min_{x \in \Omega_P} \|s \odot (x - y)\|, y \in \Omega_M\}$ , where  $\Omega_M, \Omega_P$  are the boundaries extracted from  $M$  and  $P$ ,  
 152 respectively, and  $s^\top = (1.2, 1.2, 1.5)$  is the pixel spacing in mm. Comparing the manual and predicted  
 153 organ volumes was motivated by the field of application of this study. Indeed, from the perspective  
 154 of adaptive radiotherapy, the organs' volumes are needed in order to compare the initial CT plan  
 155 dose-volume histograms for the bladder, rectum, and prostate with the doses actually delivered as  
 156 determined from CBCT scans acquired during the image-guided treatment [27]. The manual and  
 157 predicted organ volumes were compared using a Bland-Altman plot, which allows quantification of  
 158 the agreement between two quantitative measurements (i.e., the manual and predicted organ volumes)  
 159 by studying their mean difference and constructing limits of agreement [28]. We computed the bias as

$$\text{Bias} = \frac{1}{n} \sum_{i=1}^n (p_i - m_i), \quad (4)$$

160 where  $n$  is the number of patients in the test set and  $p_i = s_1 \times s_2 \times s_3 \times |M_i|$ ,  $m_i = s_1 \times s_2 \times s_3 \times |P_i|$  are  
 161 the volumes of the manual and predicted segmentations of the  $i$ -th patient. It provides the systematic  
 162 under- or overestimation of the predicted volumes. We also computed the precision,

$$\text{Precision} = \frac{1}{n} \sum_{i=1}^n |p_i - m_i|, \quad (5)$$

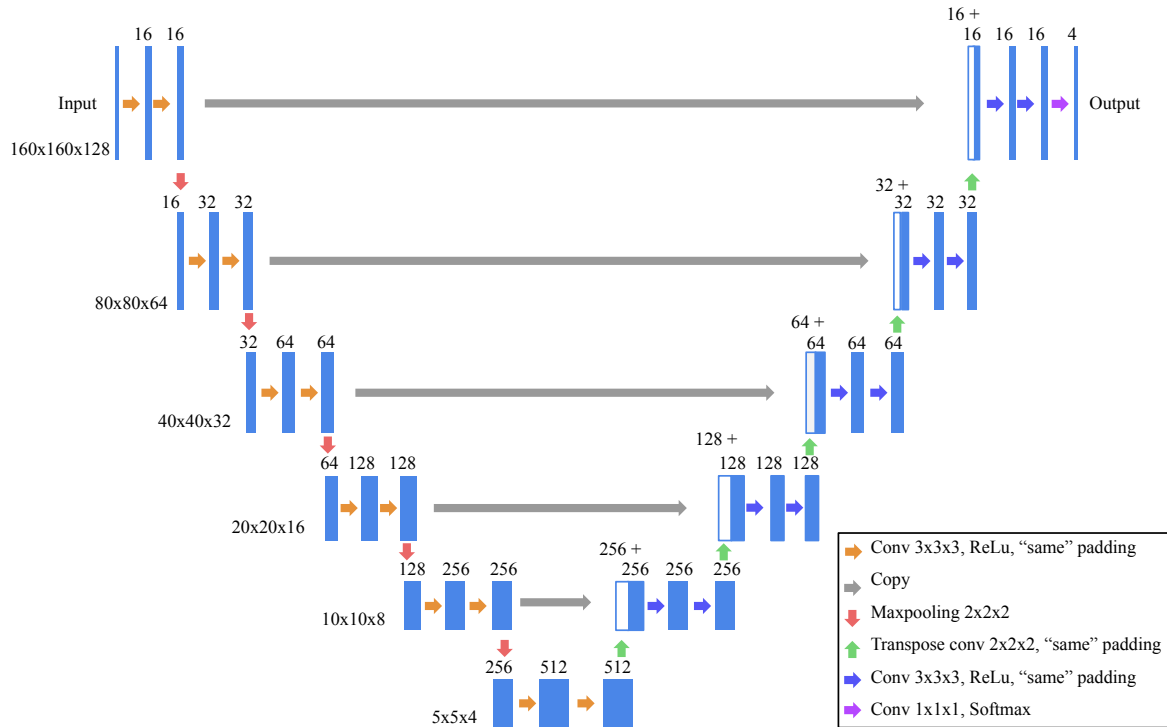
163 which measures the difference between manual and predicted volume (in absolute value).

164  
 165 The DL-based segmentation was compared with different alternative approaches as summarized  
 166 in Table 2. Two segmentation methods based on deformable image registration (denoted DIR in Table  
 167 2, second column) were applied to our dataset. First, the contours from the planning CT scans of set  $\mathcal{S}_2$   
 168 were mapped to the follow-up CBCT scans of the same patient by using a rigid registration followed  
 169 by DIR with the ANACONDA algorithm without controlling regions of interest (ROIs) in RayStation<sup>1</sup>  
 170 (version 5.99.50.22) [29]. This algorithm adopts an intensity-based registration. Second, the contour  
 171 was mapped from the planning CT scan to the follow-up CBCT scan using the diffeomorphic  
 172 morphons DIR algorithm implemented in OpenReggui<sup>2</sup> [30]. This method exploits the local phase of

<sup>1</sup> <https://www.raysearchlabs.com/raystation/>.

<sup>2</sup> <https://openreggui.org/>.

173 the image volumes to perform the registration. Therefore, it is suited for registering image volumes  
 174 with different contrast enhancement, such as CT and CBCT scans. The diffeomorphic version of the  
 175 algorithm forces anatomically plausible deformations. We also compared our DL method with the  
 176 Mattes mutual information rigid registration algorithm [31], as implemented in OpenReggui.  
 177



**Figure 3.** 3D U-Net model architecture. Each blue rectangle represents the feature maps resulting from a convolution operation, while white rectangles represent copied feature maps. For the convolutions, the zero padding was chosen such that the volume size was preserved ("same" padding). The output size is 4: one per segmentation (bladder, rectum, and prostate) and one for the background.

### 178 3. Results

179 In this section, we assess the performance of our algorithm in terms of overlap (i.e., DSC and  
 180 JI), distance (i.e., SMBD), and volume agreement measurements. In the first part, we compare the  
 181 overlaps and distances measured between our algorithm in different settings and the considered  
 182 DIR-based segmentation approaches. In the second part, we further evaluate the performance  
 183 of our best algorithm (i.e., 3D U-net trained with all available CT and CBCT scans) by assessing  
 184 whether the predicted organ volumes are in good agreement with the volumes determined by manual  
 185 segmentation. This is done by Bland-Altman analysis.  
 186

187 In Figure 4, the DSC between the segmentation output of the fully convolutional neural network  
 188 (FCN) and the ground truth segmentation were computed and averaged over all 63 CBCT scans from  
 189 the three test folds. This was done for different numbers of training CBCT and CT scans. The results  
 190 were then compared with the RayStation DIR algorithm, diffeomorphic morphons algorithm, and rigid  
 191 registration. Table 2 completes the plots in Figure 4 by providing the means and standard deviations  
 192 of the DSC, JI, and SMBD for different numbers of training CBCT scans and different numbers of  
 193 training CT scans. The statistical model used for comparing the performances is a mixed model with a  
 194 random intercept on the patient. It showed significant differences between algorithms' performance  
 195 for all organs regarding their DSC (bladder, rectum, prostate  $p < 10^{-3}$ ), JI (bladder, rectum, prostate

196  $p < 10^{-3}$ ), and SMBD (bladder, rectum, prostate  $p < 10^{-3}$ ). In the following paragraphs, the notation  
197 Ours( $n_{CBCT}$ ,  $n_{CT}$ ) stands for the 3D U-net proposed in this study with  $n_{CBCT}$  CBCT scans and  $n_{CT}$  CT  
198 scans in the training set. The  $P$ -values provided below were obtained by performing a Tukey's range  
199 test on the DSCs. The following observations can be made based on Figure 4 and Table 2.

200  
201 First, CBCT scans are more valuable than CT scans to train a CBCT segmentation model. This is  
202 not surprising, and supported by the observation that a model trained on 40 CBCT and 0 CT scans  
203 performed significantly better than a model trained on 0 CBCT and 40 CT scans for all organs (bladder,  
204 rectum, prostate  $p < 10^{-3}$ ). The DSCs reached 0.634, 0.286, and 0.525 with Ours(0, 40) and 0.845, 0.754,  
205 and 0.722 with Ours(40, 0), for the bladder, rectum, and prostate, respectively. Also, a model trained  
206 only on 74 CT scans reached approximately the same performance as a network trained on only 6 to 10  
207 CBCT scans for all the organs. Moreover, the more CBCT scans there were in the training set, the higher  
208 the DSCs on the test set were. This result makes sense since adding new CBCT scans to the training set  
209 allows the network to generalize on the test set (exclusively composed of CBCT scans) better. More  
210 surprisingly, we observed that once a sufficient number (typically 20) of CBCT scans were part of the  
211 training set, the benefit of adding CBCT or CT scans was practically the same. Indeed, compared with a  
212 model trained on 20 CBCT and 20 CT scans, the model trained on 40 CBCT and 0 CT scans did not lead  
213 to a significant improvement in performance (bladder  $p = 0.877$ , rectum  $p = 0.700$ , prostate  $p = 0.629$ ).  
214 The DSCs reached 0.815, 0.731, and 0.682 with Ours(20, 20) for the bladder, rectum, and prostate,  
215 respectively. This confirms that augmenting a CBCT training set with CT scans might be quite valuable.

216  
217 Second, from the CT perspective, we clearly observed that the more CT scans there were in the  
218 training set, the higher the mean DSC became. Indeed, Ours(20, 74) is significantly better than Ours(20,  
219 0) for all organs (bladder, rectum  $p < 10^{-3}$ , prostate  $p < 10^{-2}$ ). We explain this improvement by the  
220 learning of more generic features, leading to better generalization. However, we observed that the  
221 difference in the average DSC between Ours(20, 0) and Ours(20, 20) was approximately equal to the  
222 difference in the average DSC between Ours(20, 20) and Ours(20, 74), whereas 20 new CT scans were  
223 added to the training set in the first case, and 54 new CT scans, in the second case. This may indicate  
224 saturation of the performance improvement produced by adding CT scans to the training set. Moreover,  
225 when the number of training CBCT scans was large, adding training CT scans improved performance  
226 for the rectum only ( $p < 0.01$ ): no statistically significant incremental change in performance was  
227 observed for the bladder or prostate ( $p = 0.780$  and  $p = 0.630$ , respectively) when Ours(42, 74) and  
228 Ours(42, 0) were compared. A plausible interpretation is that most of the useful information present in  
229 the CT scans was already captured in the relatively large CBCT training set. More importantly, in  
230 line with our objective of limiting the annotation of CBCT scans, we observed that the performance  
231 obtained with 42 CBCT and 0 CT scans could be reached with 20 CBCT and 74 CT scans for all organs  
232 (bladder  $p = 0.940$ , rectum  $p = 0.882$ , prostate  $p = 0.994$ ). Hence, the availability of 74 annotated  
233 CT scans reduced the number of annotated CBCT scans significantly (by a factor of approximately two).

234  
235 Third, when all available CT and CBCT scans (42 CBCT and 74 CT scans) were used for training,  
236 our approach significantly outperformed the rigid registration, RayStation DIR algorithm, and  
237 diffeomorphic morphons algorithm for the bladder and rectum ( $p < 10^{-3}$ ) but not for the prostate  
238 ( $p = 0.911$ ). These conclusions are illustrated on a representative patient in Figure 5. The results also  
239 show that the rigid registration is outperformed by the ANACONDA algorithm, which is in turn  
240 outperformed by the diffeomorphic morphons algorithm for the bladder and rectum. As mentioned  
241 above, both DIR methods are statistically similar to the rigid registration approach when it comes to  
242 segmenting the prostate. This supports the hypothesis that the prostate undergoes less deformation  
243 than the bladder and rectum, which are subject to regular influxes and voiding of matter. Although  
244 our analysis is based on the DSC, both the JI and the SMBD lead to the same conclusions.

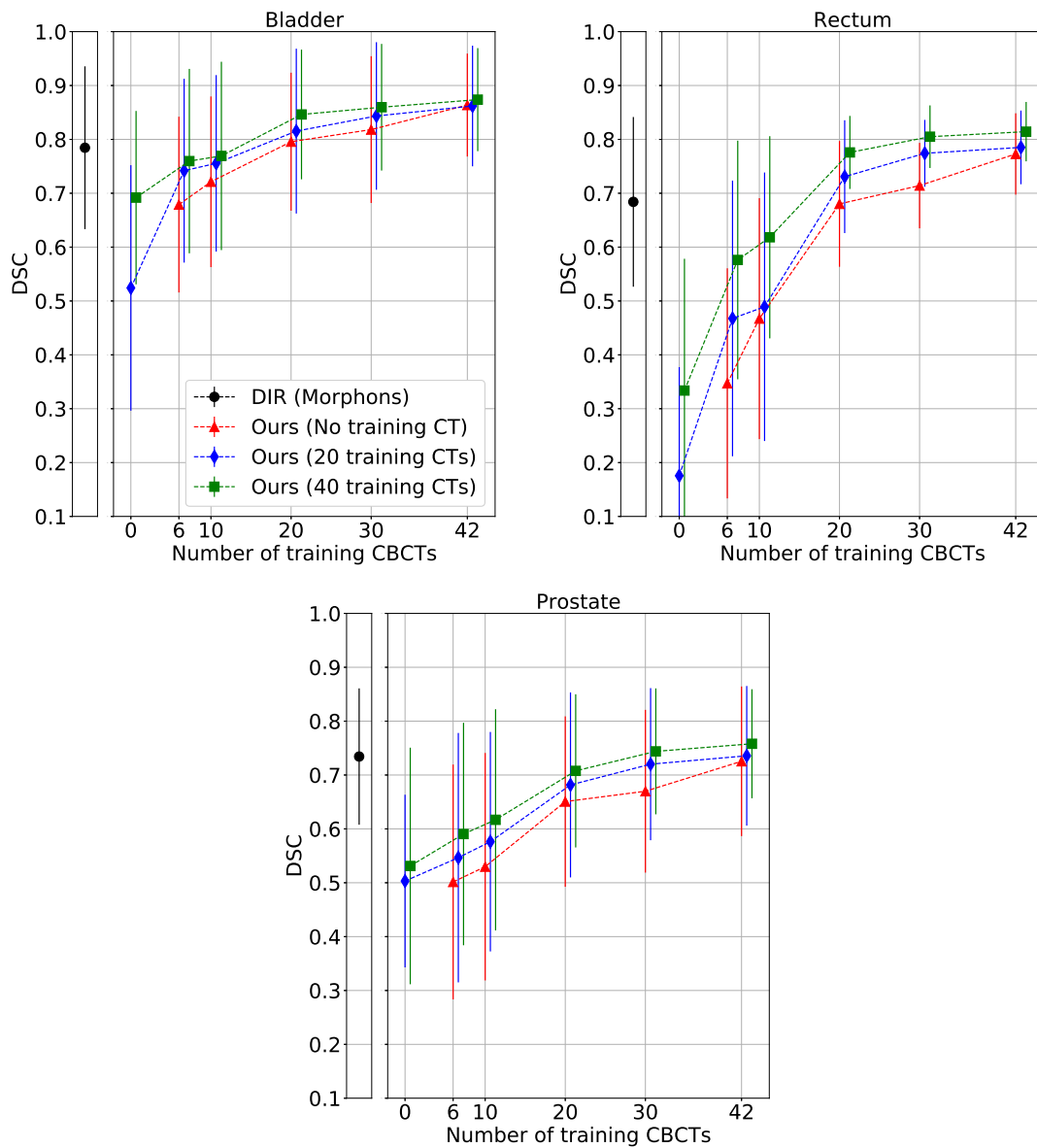
245

246 Figure 6 presents Bland-Altman plots comparing the organ volumes reached manually and by  
247 our DL-based predictions (obtained with Ours(42, 74)), using the bias, precision, and 95% limits of  
248 agreements (LoA). The bias normalized by the manual volume is below 5 % for all organs (bladder  
249 4.78%, rectum 1.21%, prostate 2.51%). The precision normalized by the manual volume is similar for  
250 the bladder and the rectum (bladder 13.3%, rectum 13.9%) and larger for the prostate (27.9%). The LoA  
251 of the bladder are also close to the LoA of the rectum (-32% and 41% for the bladder and -33% and  
252 35% for the rectum), whereas they are larger for the prostate (-65% and 70%). Table 3 completes the  
253 Bland-Altman plots by providing the means and standard deviations for the manual and predicted  
254 organ volumes. Moreover, a one-sample *t*-test was performed on the differences between the manual  
255 and predicted volumes normalized by the manual volume for each organ. The resulting *P*-values  
256 for all organs are presented in Table 3 and are not significantly different (bladder  $p = 0.285$ , rectum  
257  $p = 0.897$ , prostate  $p = 0.438$ ). This means that the predicted and manual contours are similar in  
258 means according to the *t*-test.

259

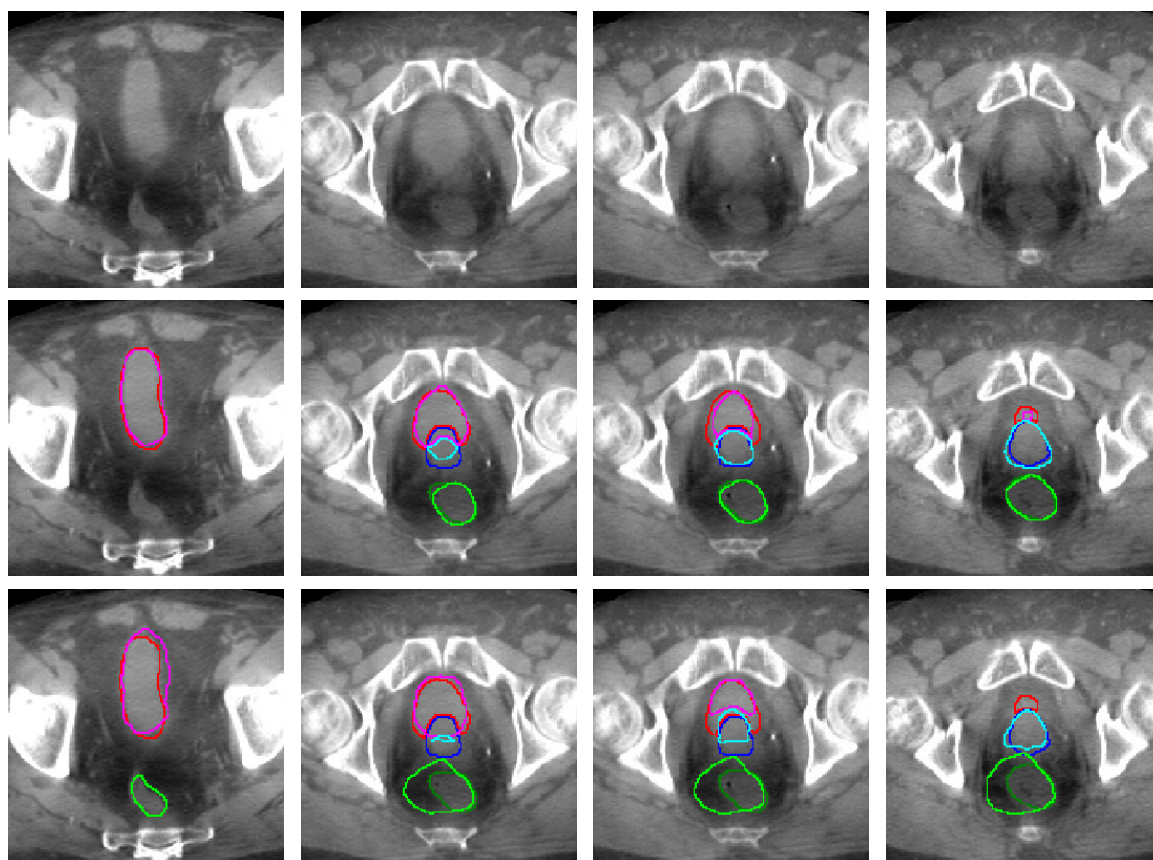
260 Computational cost analysis was performed by measuring the running time on our machine  
261 equipped with a 11Gb GeForce GTX 1080 Ti graphics card. The rigid registration of one image ran in  
262 1.05 min. The deformable image registration with the ANACONDA and morphons algorithms ran in  
263 1.92 min and 8.33 min, respectively. The inference time for one image with the DL approaches was  
264 much lower. It reached 0.15 s independently of the learning strategy. Indeed, the number of images in  
265 the training set has no impact on the inference time. The training time needed to reach convergence  
266 depends on the size of the training set. Hence, Ours(20, 0), Ours(20, 74), Ours(42, 0), and Ours(42, 74)  
267 were trained in 17.3, 224, 167, and 220 min, respectively.



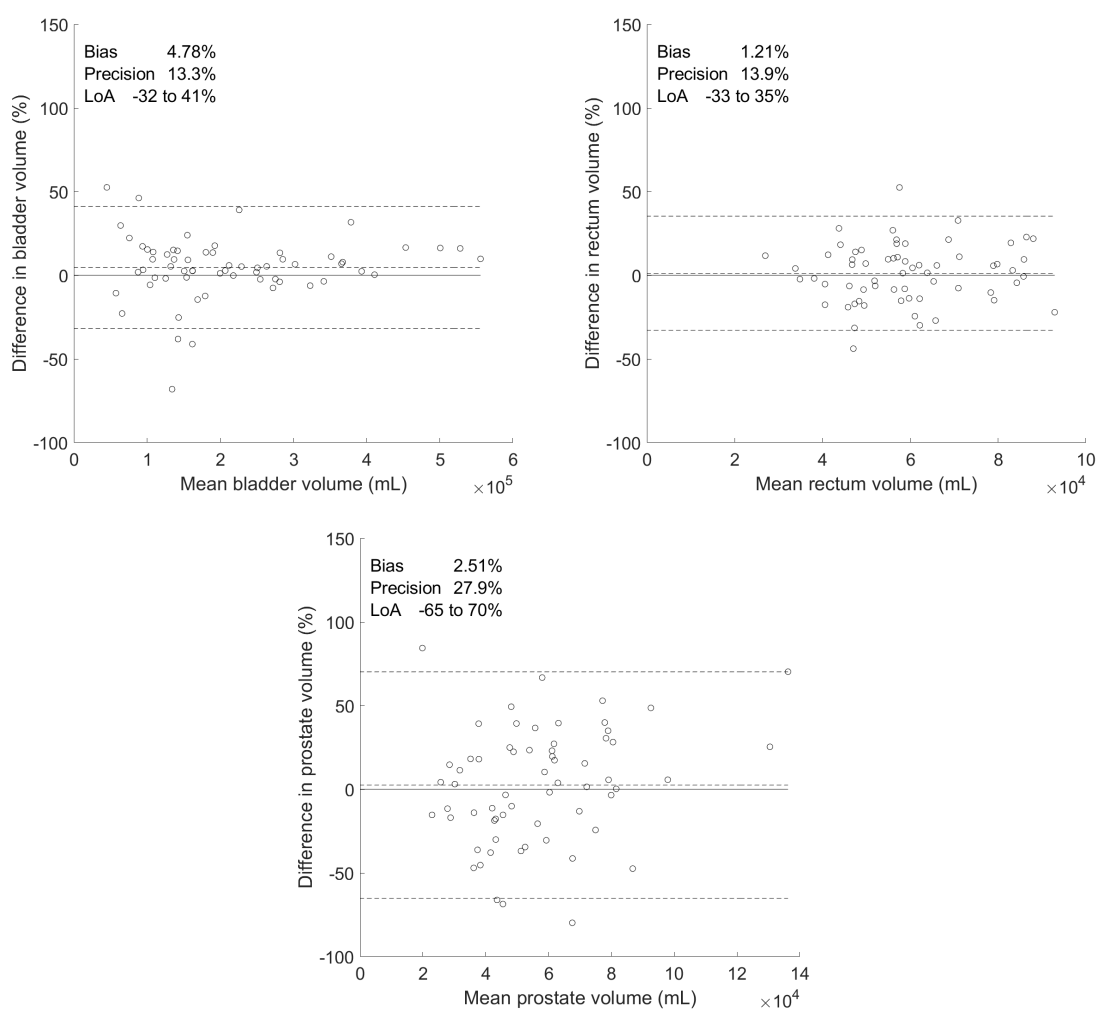


**Figure 4.** Influence of the number of training CBCT and CT scans on the DSC. Bars indicate one standard deviation for the group of 63 patients. DSC: Dice similarity coefficient.





**Figure 5.** Comparison of manual, 3D U-net, and morphons DIR-based segmentation for a representative patient. Each column corresponds to a slice of the same CBCT scan. Dark colors represent reference segmentations (both second and third rows), while light colors show 3D U-net segmentation (second row) and morphons DIR-based segmentation (third row). The predicted bladder, in pink, has a DSC of 0.940 (U-net) and 0.864 (morphons); the rectum, in light green, has a DSC of 0.791 and 0.759; the prostate, in light blue, has a DSC of 0.780 and 0.730.



**Figure 6.** Bland–Altman plots for the bladder, rectum, and prostate derived from the differences between the predicted and manual segmentations. The solid lines represent no difference; the dotted lines depict the mean difference (bias) and 95% limits of agreements (LoA).

**Table 3.** Absolute and relative differences between manual and predicted organs volumes. *P*-values are calculated using a one-sample *t*-test on percentage differences.

Organ	Volumes ( $\times 10^4$ mL)		Differences between manual and predicted volumes				
			Absolute ( $\times 10^4$ mL)		Percentage (%)		
	Manual	Predicted	Bias	Precision	Bias	Precision	<i>P</i> -value
Bladder	$21.9 \pm 12.9$	$20.7 \pm 11.4$	1.18	2.46	4.78	13.3	.285
Rectum	$5.96 \pm 1.66$	$5.87 \pm 1.55$	.094	.826	1.21	13.9	.897
Prostate	$5.87 \pm 2.98$	$5.53 \pm 2.07$	.340	1.64	2.51	27.9	.438

#### 268 4. Discussion

269 Based on Table 2 (first part) and Figure 4, 3D U-net approach is the most satisfactory approach for  
 270 automatic segmentation of the bladder and rectum on CBCT scans. This supports the initial hypothesis  
 271 that registration-based approaches fail in the case of large deformation and alternative approaches  
 272 using the statistics of the target image (i.e., the CBCT scan) are more suitable. This observation is also  
 273 consistent with the state-of-the-art algorithms shown in Table 2 (second part), where DL approaches

274 outperform alternative approaches for the bladder and rectum.

275

276 Still based on Table 2 (first part) and Figure 4, the 3D U-net slightly outperforms the  
277 registration-based approaches for the prostate, but this improvement is not statistically significant. 3D  
278 U-net's lower performances for the prostate than for the bladder and rectum is further supported by  
279 the Bland-Altman analysis of the manual and predicted volumes. Indeed, this analysis provides less  
280 than 5% bias for all organs but higher precision (i.e., a larger spread of the predictions, as defined  
281 in (5)) for the prostate than for the bladder and rectum. Also, most other state-of-the-art DIR-based  
282 algorithms outperform our approach for the prostate. This shows that DIR-based approaches are still  
283 valuable in situations with limited organ deformation and where poor contrast makes the use of  
284 vanilla DL models challenging. A first way to improve the segmentation results for the prostate and  
285 outperform DIR-based approaches without annotating more CBCT scans might be to generate pseudo  
286 CBCT scans as in Schreier et al., but our study shows that increasing the number of already annotated  
287 CT scans further is a valuable alternative, albeit with a risk of saturation. If few data are available, a  
288 second option could be to promote a desired shape or structure in the deep model prediction [32,33].  
289 A third option could be to perform unsupervised domain adaptation [34]. This approach requires  
290 annotations in a source domain (CT) but not in the target domain (CBCT). This will be the subject of  
291 future research.

292

293 From an application point of view, the study shows that the more CBCT scans are contoured, the  
294 better the DSC on the predicted contours. However, contouring CBCT scans is not part of the clinical  
295 workflow, is time-consuming, and is not easy because of the poor contrast between the different  
296 regions of interest. Hence, we have shown that expanding the training set with CT scans improves the  
297 segmentation performances for all considered organs, especially when few contoured CBCT scans are  
298 available. Our 3D U-net that reached the best segmentation performances was trained with 42 CBCT  
299 and 74 CT scans.

300

301 Most cases of failure have been observed for the prostate, which has the lowest DSC of the organs.  
302 This may be due to the fact that the prostate is hard to see on CBCT scans and often pushes on the  
303 bladder as we can see in Figure 5. Hence, some upper parts of the prostate are often wrongly classified  
304 as bladder, which decreases the DSC for the prostate. Since the bladder is larger than the prostate,  
305 misclassification at the boundary between the two organs has less impact on the DSC of the bladder. A  
306 second case of failure occurs at the top and bottom slices of the rectum, which is wrongly classified  
307 as background (or inversely, background is wrongly classified as rectum). This makes sense since  
308 there are few differences in contrast between the rectum, anal canal, and colon. The impact of such  
309 errors on the prostate and rectum, as well as the required contour quality for clinical use in adaptive  
310 radiotherapy, is such that additional quality assessment with a contours review process is needed.  
311 This should be done by radiation oncologists and will be the subject of future research.

312

313 Our DL approach also outperforms or achieves the same performance as patient specific  
314 models for the bladder. Those models rely on PCA to extract principal modes of deformation from  
315 landmarks placed on the bladder's contour and across several contoured images for each patient  
316 being considered. The drawback for clinical use of such approaches is that (i) a different model  
317 is required for every patient and organ and (ii) several images per patient are needed to build the model.

318

319 Concerning alternative DL methods, the current work slightly outperforms our initial conference  
320 paper, Brion et al. [24], on bladder segmentation with 3D U-net. This is probably due to the larger  
321 training database and/or the multi-class formulation used in this work, since three organs were  
322 segmented instead of one. Only 41 of the patients used in our conference paper were kept in this  
323 study. This is because the remaining patients had either had their prostates removed or lacked fully

324 annotated scans. New patients were also added. The two datasets are thus different. However,  
325 Schreier et al.'s work is the closest to this study. Hence, we did a more thorough comparison with  
326 their findings. They obtained a higher DSC than we did for all the organs considered in this study.  
327 This might be explained by the fact that they used more samples in their training set (300 CT and  
328 300 pseudo CBCT scans compared with 74 CT and 42 CBCT scans). However, it is hard to determine  
329 whether this is the only explanation for their better results. Indeed, in Figure 4, we see that the  
330 DSC increases more slowly as the number of training samples increases. Interestingly, they ran the  
331 patch-wise 3D U-net proposed by Hänsch et al. on their test set and got DSCs of 0.927, 0.860, and 0.816  
332 for the bladder, rectum, and prostate, respectively. Those results are higher than the results obtained  
333 on the bladder (DSC = 0.88) and rectum (DSC = 0.71) by Hänsch et al. So, their test set might be of a  
334 higher quality than ours, which can be a limitation on their approach in clinical practice, where low  
335 quality images are common. Another shortcoming is that they report their results on a dataset that  
336 includes both CBCT and CT scans (10%). It is therefore unclear how well their method would perform  
337 on a dataset containing only CBCT scans (such as ours). As a final remark, their proposed generation  
338 of pseudo CBCT scans from clinically contoured CT scans is a powerful tool for solving the problem of  
339 CBCT annotations. However, such knowledge of artificial data generation might not be present in all  
340 hospitals. To summarize this comparison, we consider the two publications to be complementary, with  
341 our strengths being the size of our test set, detailed comparison with registration approaches, and  
342 detailed study of the impact of additional CT scans in the training database.

343

## 344 5. Conclusions

345 In this work, a 3D U-net DL model was trained on CBCT and CT scans in order to segment the  
346 bladder, rectum, and prostate on CBCT scans. The proposed approach significantly outperformed all  
347 the DIR-based segmentation methods applied on our dataset in terms of DSC, JI, and SMBD for the  
348 bladder and rectum. The conclusions are more mitigated concerning the prostate, where the DL-based  
349 segmentation did not significantly outperform alternative approaches. A Bland-Altman analysis on the  
350 manual and predicted organs volumes revealed a low bias on the predicted volumes for all organs but  
351 higher precision (i.e., a larger spread of the volumes) for the prostate than for the other organs. Also,  
352 the study shows that the cross-domain data augmentation consisting in adding CT to the CBCT scans  
353 in the training set significantly improved the segmentation results. A further step will be to highlight  
354 these improvements by showing the better tumor coverage and reduction in the doses delivered to  
355 organs at risk that it allows.

356 **Author Contributions:** Conceptualization: J.L., E.B., P.D., C.D.V., J.A.L. and B.M.; methodology: J.L., E.B., P.D.,  
357 C.D.V., J.A.L. and B.M.; software: J.L., E.B. and P.D.; validation: J.L. and E.B.; formal analysis, J.L., E.B. and  
358 P.D.; investigation: J.L., E.B. and P.D.; resources: J.L. and E.B.; data curation: J.L. and E.B.; writing—original draft  
359 preparation: J.L., E.B. and P.D.; writing—review and editing: J.L., E.B., P.D., C.D.V., J.A.L. and B.M.; visualization:  
360 J.L., E.B. and P.D.; supervision: C.D.V., J.A.L. and B.M.; project administration: J.L., E.B. and B.M.; funding  
361 acquisition: C.D.V. and B.M.

362 **Funding:** Jean Léger is a Research Fellow of the Belgian national scientific research foundation *Fonds de la Recherche*  
363 *Scientifique - FNRS*. Elliott Brion is supported by the Walloon Region under grant RW-DGO6-Biowin-Bidmed. Paul  
364 Desbordes is a member of the Prother-wal project funded by the Walloon Region (Belgium). John A. Lee and  
365 Christophe De Vleeschouwer are Senior Research Associates with the Belgian F.R.S.-FNRS.

366 **Acknowledgments:** We thank CHU-UCL-Namur (Dr. Vincent Remouchamps) and CHU-Charleroi (Dr. Nicolas  
367 Meert) for providing the data. The authors want to thank Sara Teruel Rivas for her technical support in the data  
368 acquisition and annotation and to Dr. Geneviève Van Ooteghem for the verification of the contours. Finally, we  
369 thank Gabrielle Leyden for editing the final revision of this paper.

370 **Conflicts of Interest:** The authors declare no conflict of interest.

## 371 Abbreviations

372 The following abbreviations are used in this manuscript:

373

CBCT	Cone beam computed tomography
CT	Computed tomography
CTV	Clinical target volume
DIR	Deformable image registration
DL	Deep learning
DSC	Dice similarity coefficient
DVF	Deformation vector field
374 EBRT	External beam radiation therapy
FCN	Fully convolutional neural network
GPU	Graphical processing unit
JI	Jaccard index
LoA	Limit of agreement
OAR	Organ at risk
ROI	Region of interest
SMBD	Symmetric mean boundary distance

## 375 References

- 376 1. Brousmiche, S.; Orban de Xivry, J.; Macq, B.; Seco, J. SU-E-J-125: Classification of CBCT Noises in Terms of  
377 Their Contribution to Proton Range Uncertainty. *Medical Physics* **2014**, *41*, 184–184.
- 378 2. Peng, C.; Ahunbay, E.; Chen, G.; Anderson, S.; Lawton, C.; Li, X.A. Characterizing interfraction variations  
379 and their dosimetric effects in prostate cancer radiotherapy. *International Journal of Radiation Oncology\**  
380 *Biology\* Physics* **2011**, *79*, 909–914.
- 381 3. Ghilezan, M.; Yan, D.; Martinez, A. Adaptive Radiation Therapy for Prostate Cancer. *Seminars in Radiation*  
382 *Oncology* **2010**, *20*, 130–137. doi:10.1016/j.semradonc.2009.11.007.
- 383 4. Pos, F.; Remeijer, P. Adaptive Management of Bladder Cancer Radiotherapy. *Seminars in Radiation Oncology*  
384 **2010**, *20*, 116–120. doi:10.1016/j.semradonc.2009.11.005.
- 385 5. Wang, Y.; Efstathiou, J.A.; Sharp, G.C.; Lu, H.M.; Ciernik, I.F.; Trofimov, A.V. Evaluation of the dosimetric  
386 impact of interfractional anatomical variations on prostate proton therapy using daily in-room CT images.  
387 *Medical physics* **2011**, *38*, 4623–33. doi:10.1118/1.3604152.
- 388 6. Moteabbed, M.; Trofimov, A.; Sharp, G.C.; Wang, Y.; Zietman, A.L.; Efstathiou, J.A.; Lu, H.M. Proton  
389 therapy of prostate cancer by anterior-oblique beams: Implications of setup and anatomy variations.  
390 *Physics in Medicine and Biology* **2017**, *62*, 1644–1660. doi:10.1088/1361-6560/62/5/1644.
- 391 7. Rigaud, B.; Simon, A.; Castelli, J.; Lafond, C.; Acosta, O.; Haigron, P.; Cazoulat, G.; de Crevoisier, R.  
392 Deformable image registration for radiation therapy: principle, methods, applications and evaluation.  
393 *Acta Oncologica* **2019**, pp. 1–13.
- 394 8. Oh, S.; Kim, S. Deformable image registration in radiation therapy. *Radiation oncology journal* **2017**, *35*, 101.
- 395 9. Motegi, K.; Tachibana, H.; Motegi, A.; Hotta, K.; Baba, H.; Akimoto, T. Usefulness of hybrid deformable  
396 image registration algorithms in prostate radiation therapy. *Journal of applied clinical medical physics* **2019**,  
397 *20*, 229–236.
- 398 10. Takayama, Y.; Kadoya, N.; Yamamoto, T.; Ito, K.; Chiba, M.; Fujiwara, K.; Miyasaka, Y.; Dobashi, S.; Sato,  
399 K.; Takeda, K.; others. Evaluation of the performance of deformable image registration between planning  
400 CT and CBCT images for the pelvic region: comparison between hybrid and intensity-based DIR. *Journal*  
401 *of radiation research* **2017**, *58*, 567–571.
- 402 11. Zambrano, V.; Furtado, H.; Fabri, D.; Lütgendorf-Caucig, C.; Góra, J.; Stock, M.; Mayer, R.; Birkfellner, W.;  
403 Georg, D. Performance validation of deformable image registration in the pelvic region. *Journal of radiation*  
404 *research* **2013**, *54*, i120–i128.
- 405 12. Thor, M.; Petersen, J.B.; Bentzen, L.; Høyer, M.; Muren, L.P. Deformable image registration for contour  
406 propagation from CT to cone-beam CT scans in radiotherapy of prostate cancer. *Acta Oncologica* **2011**,  
407 *50*, 918–925.
- 408 13. Söhn, M.; Birkner, M.; Chi, Y.; Wang, J.; Yan, D.; Berger, B.; Alber, M. Model-independent, multimodality  
409 deformable image registration by local matching of anatomical features and minimization of elastic energy.  
410 *Medical physics* **2008**, *35*, 866–878.
- 411 14. Thirion, J.P. Image matching as a diffusion process: an analogy with Maxwell's demons **1998**.

- 412 15. Woerner, A.J.; Choi, M.; Harkenrider, M.M.; Roeske, J.C.; Surucu, M. Evaluation of deformable image  
413 registration-based contour propagation from planning CT to cone-beam CT. *Technology in cancer research &*  
414 *treatment* **2017**, *16*, 801–810.
- 415 16. König, L.; Derksen, A.; Papenberg, N.; Haas, B. Deformable image registration for adaptive radiotherapy  
416 with guaranteed local rigidity constraints. *Radiation Oncology* **2016**, *11*, 122.
- 417 17. Chai, X.; van Herk, M.; Betgen, A.; Hulshof, M.; Bel, A. Automatic bladder segmentation on CBCT for  
418 multiple plan ART of bladder cancer using a patient-specific bladder model. *Physics in Medicine & Biology*  
419 **2012**, *57*, 3945.
- 420 18. van de Schoot, A.; Schooneveldt, G.; Wognum, S.; Hoogeman, M.; Chai, X.; Stalpers, L.; Rasch, C.; Bel, A.  
421 Generic method for automatic bladder segmentation on cone beam CT using a patient-specific bladder  
422 shape model. *Medical physics* **2014**, *41*.
- 423 19. Kazemifar, S.; Balagopal, A.; Nguyen, D.; McGuire, S.; Hannan, R.; Jiang, S.; Owrangi, A. Segmentation of  
424 the prostate and organs at risk in male pelvic CT images using deep learning. *arXiv preprint arXiv:1802.09587*  
425 **2018**.
- 426 20. Cha, K.H.; Hadjiiski, L.; Samala, R.K.; Chan, H.P.; Caoili, E.M.; Cohan, R.H. Urinary bladder segmentation  
427 in CT urography using deep-learning convolutional neural network and level sets. *Medical physics* **2016**,  
428 *43*, 1882–1896.
- 429 21. Çiçek, Ö.; Abdulkadir, A.; Lienkamp, S.S.; Brox, T.; Ronneberger, O. 3D U-Net: learning dense volumetric  
430 segmentation from sparse annotation. International Conference on Medical Image Computing and  
431 Computer-Assisted Intervention. Springer, 2016, pp. 424–432.
- 432 22. Haensch, A.; Dicken, V.; Gass, T.; Morgas, T.; Klein, J.; Meine, H.; Hahn, H. Deep learning based  
433 segmentation of organs of the female pelvis in CBCT scans for adaptive radiotherapy using CT and CBCT  
434 data. *Int J Comput Assist Radiol Surg* **2018**, *13*, 179–180.
- 435 23. Hänsch, A.; Dicken, V.; Klein, J.; Morgas, T.; Haas, B.; Hahn, H.K. Artifact-driven sampling schemes for  
436 robust female pelvis CBCT segmentation using deep learning. Medical Imaging 2019: Computer-Aided  
437 Diagnosis. International Society for Optics and Photonics, 2019, Vol. 10950, p. 109500T.
- 438 24. Brion, E.; Léger, J.; Javaid, U.; Lee, J.; De Vleeschouwer, C.; Macq, B. Using planning CTs to enhance  
439 CNN-based bladder segmentation on cone beam CT. Medical Imaging 2019: Image-Guided Procedures,  
440 Robotic Interventions, and Modeling. International Society for Optics and Photonics, 2019, Vol. 10951, p.  
441 109511M.
- 442 25. Schreier, J.; Genghi, A.; Laaksonen, H.; Morgas, T.; Haas, B. Clinical evaluation of a full-image deep  
443 segmentation algorithm for the male pelvis on cone-beam CT and CT. *Radiotherapy and Oncology* **2020**,  
444 *145*, 1–6.
- 445 26. Ronneberger, O.; Fischer, P.; Brox, T. U-net: Convolutional networks for biomedical image segmentation.  
446 International Conference on Medical image computing and computer-assisted intervention. Springer, 2015,  
447 pp. 234–241.
- 448 27. Hatton, J.A.; Greer, P.B.; Tang, C.; Wright, P.; Capp, A.; Gupta, S.; Parker, J.; Wratten, C.; Denham, J.W.  
449 Does the planning dose–volume histogram represent treatment doses in image-guided prostate radiation  
450 therapy? Assessment with cone-beam computerised tomography scans. *Radiotherapy and Oncology* **2011**,  
451 *98*, 162–168.
- 452 28. Giavarina, D. Understanding bland altman analysis. *Biochemia medica: Biochemia medica* **2015**, *25*, 141–151.
- 453 29. Weistrand, O.; Svensson, S. The ANACONDA algorithm for deformable image registration in radiotherapy.  
454 *Medical physics* **2015**, *42*, 40–53.
- 455 30. Janssens, G.; Jacques, L.; de Xivry, J.O.; Geets, X.; Macq, B. Diffeomorphic registration of images with  
456 variable contrast enhancement. *Journal of Biomedical Imaging* **2011**, *2011*, 3.
- 457 31. Mattes, D.; Haynor, D.R.; Vesselle, H.; Lewellen, T.K.; Eubank, W. PET-CT image registration in the chest  
458 using free-form deformations. *IEEE transactions on medical imaging* **2003**, *22*, 120–128.
- 459 32. Oktay, O.; Ferrante, E.; Kamnitsas, K.; Heinrich, M.; Bai, W.; Caballero, J.; Cook, S.A.; De Marvao, A.;  
460 Dawes, T.; O’Regan, D.P.; others. Anatomically constrained neural networks (ACNNs): application to  
461 cardiac image enhancement and segmentation. *IEEE transactions on medical imaging* **2017**, *37*, 384–395.
- 462 33. Ravishankar, H.; Venkataramani, R.; Thiruvenkadam, S.; Sudhakar, P.; Vaidya, V. Learning and  
463 incorporating shape models for semantic segmentation. International Conference on Medical Image  
464 Computing and Computer-Assisted Intervention. Springer, 2017, pp. 203–211.



465 34. Kamnitsas, K.; Baumgartner, C.; Ledig, C.; Newcombe, V.; Simpson, J.; Kane, A.; Menon, D.; Nori, A.;  
466 Criminisi, A.; Rueckert, D.; others. Unsupervised domain adaptation in brain lesion segmentation with  
467 adversarial networks. *International Conference on Information Processing in Medical Imaging*. Springer,  
468 2017, pp. 597–609.

469 **Sample Availability:** Access to the dataset is subjected to the authorization of the partner hospitals' ethics  
470 committees. The dataset is not available by default.

471 © 2020 by the authors. Submitted to *Appl. Sci.* for possible open access publication under the terms and conditions  
472 of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).