

**École polytechnique de Louvain**

# **Dose prediction for protontherapy using neural networks**

Author: **Margerie HUET DASTARAC**  
Supervisors: **John LEE, Edmond STERPIN**  
Readers: **Ana M. BARRAGAN, Christophe DE VLEESCHOUWER**  
Academic year 2018–2019  
Master [120] in Electrical Engineering

## **Abstract**

Proton therapy is an emerging modality of radiotherapy, which is one of the main cancer treatments. Compared to conventional radiotherapy, it has the promising potential sparing more healthy tissues while ensuring that the tumor gets the prescribed dose.

Both modalities require a treatment plan, which is nowadays done in a semi-automatic manner, but still requires several hours of manual works of physicians (for the contouring and prescription) and by the dosimetrist (to generate the plan). Ongoing research aims to speed up the process and tends toward a complete automatic workflow to allow the generation of a plan along the treatment.

This Master thesis contributes to this aim by proposing a model for dose prediction for proton therapy using a UNet architecture. The dosimetric quality of the model was evaluated by comparing predictions with dose delivered in plans generated by a dosimetrist. The method proposed in this study is accurate (less than 4% of mean error between the clinical metrics of the prediction and ground truth, expressed in percent of the highest prescription dose), fast and achievable through dose mimicking.

## **Acknowledgements**

I want to sincerely thank all the members of the MIRO lab. I felt welcomed from my very first days. I really enjoyed my time in the lab and could not imagine being in a better place for this Master thesis.

I want to particularly thank my supervisors, John Lee and Edmond Sterpin. I knew I could turn to them when facing a problem of comprehension or to seek for advice.

I could never thank Ana enough for all the time she took to answer my questions and give me great recommendations.

A big thank you to Elena for her patience and help for the dose mimicking.

And finally, I want to thank Sara for all the data she generated and to be always here to give me with details on the generated plans.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Context</b>	<b>7</b>
2.1	Treatment planning process . . . . .	9
2.1.1	Imaging the patient . . . . .	9
2.1.2	Contouring and dose prescription . . . . .	10
2.1.3	Plan generation . . . . .	12
2.1.4	Evaluation of Plans . . . . .	12
2.2	Robustness . . . . .	13
<b>3</b>	<b>Literature review</b>	<b>14</b>
3.1	Current methods for automatic planning . . . . .	14
3.2	Dose prediction in proton therapy . . . . .	15
3.3	Neural networks used for dose prediction in conventional radiotherapy	15
3.3.1	UNet . . . . .	16
3.3.2	Other CNN architectures . . . . .	17
3.3.3	2-D versus 3-D models . . . . .	18
3.3.4	Generative Adversarial Network . . . . .	18
<b>4</b>	<b>Methods</b>	<b>19</b>
4.1	Network Architecture . . . . .	19
4.1.1	Cost function . . . . .	20
4.1.2	Activation function . . . . .	21
4.1.3	Pooling . . . . .	21
4.1.4	Regularization . . . . .	22
4.1.5	Optimizer . . . . .	22
4.2	Data . . . . .	24
4.3	Training . . . . .	24
4.4	Patching system . . . . .	25
4.5	Adding information about the traversed tissues . . . . .	26
4.5.1	<i>No CT</i> : without Computed Tomography scan . . . . .	27

4.5.2	<i>CT</i> : with Computed Tomography scan . . . . .	27
4.5.3	<i>WEPL</i> : Water Equivalent Path Length image . . . . .	28
4.5.4	<i>Hybrid</i> : <i>WEPL</i> + dose prescription . . . . .	30
4.5.5	<i>CT masked</i> : <i>CT</i> cropped on beams trajectory . . . . .	30
4.5.6	<i>CT masked and Hybrid</i> : <i>CT</i> masked + <i>Hybrid</i> . . . . .	31
4.6	Dose Mimicking . . . . .	32
<b>5</b>	<b>Results</b>	<b>33</b>
5.1	Model Validation . . . . .	33
5.2	Results considering clinical metrics . . . . .	35
5.2.1	Target volume . . . . .	35
5.2.2	Organs at risk . . . . .	36
5.3	Comparison dose prediction and ground truth . . . . .	37
5.4	Computation time . . . . .	39
5.5	Dose Mimicking . . . . .	40
5.5.1	Conventional plan compared to prediction . . . . .	40
5.5.2	Robust plan compared to prediction . . . . .	41
5.5.3	Robust plan compared to conventional plan . . . . .	43
5.5.4	Position of spots . . . . .	44
<b>6</b>	<b>Discussion</b>	<b>45</b>
6.1	Results interpretation . . . . .	45
6.1.1	Results on CTV . . . . .	45
6.1.2	Results on the OARs . . . . .	46
6.1.3	Dose Mimicking . . . . .	47
6.1.4	Computation time . . . . .	48
6.2	Limitations . . . . .	49
6.2.1	Generalization . . . . .	49
6.2.2	Limitation on the mask . . . . .	49
6.2.3	Accuracy . . . . .	49
6.3	Perspectives . . . . .	50
6.3.1	Uncertainty estimation . . . . .	50
6.3.2	GAN . . . . .	50
<b>7</b>	<b>Conclusion</b>	<b>51</b>

# 1 | Introduction

Radiotherapy is one of the most employed cancer treatment modalities, standing among surgery and systemic techniques such as chemotherapy. It is involved in about 50% of treatments nowadays. Radiation therapy consists in irradiating the tumor, generally using an external beam. Photons or charged particles transfer energy to the tissues, resulting in dose deposition (measured in Gray [ $1\text{Gy} = 1\text{ Joule/kg}$ ]). DNA is the part of cells that is damaged when undergoing irradiation. When DNA gets damaged, recovery mechanisms try to reconstruct the lost information. If the damage is too important, reparation can not occur and the cell dies. Fortunately, reparation mechanisms of cancerous cells are less efficient than those of healthy cells. Therefore, instead of delivering the total dose at once, the dose is divided into smaller fractions and delivered every day for 6 weeks. The survival of patients depends on the ability of controlling the tumor while sparing the healthy tissues. A high dose in healthy tissues induces toxicity, which hampers the recovery of healthy cells and may increase the risk of secondary radiation-induced tumors.

Several radiation modalities have been developed. Most centers employ a X-ray beam (technique known as conventional radiotherapy). This technology is well controlled but presents some disadvantages. After a short buildup, the deposited dose decreases linearly along its trajectory, see Figure 1.1, which means that the maximum dose will be deposited at shallow depth in the patient. Conventional radiotherapy is therefore not optimal to treat tumors deep in the patient or close to critical structures. An alternative technology that is emerging is the protontherapy. It consists of delivering proton beams instead X-ray beams. Protons offer the advantageous depth dose profile as it has a finite range and deposits its maximal dose at the end of its trajectory forming a peak of dose called Bragg peak, as shown in Figure 1.1.

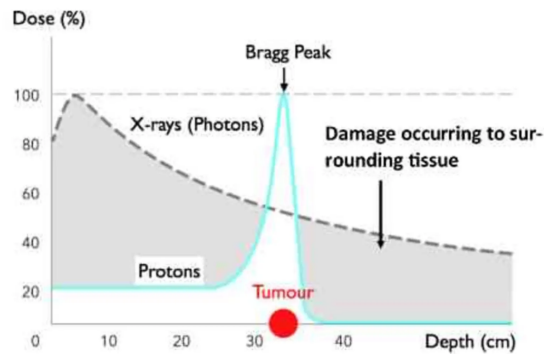


Figure 1.1: Difference between the depth-dose profile of photons and protons. Thanks to the Bragg peak, protons deposit almost no dose beyond the tumor. Credits: [1].

Whatever the chosen modality, a treatment plan needs to be established. It results from a complex design process involving multiple medical specialists and dedicated softwares. The workflow, shown in Figure 1.2, is as follows: the patient gets scanned using a Computed Tomography (CT) scanner, then the physician draws the contours of the tumor. Subsequently, the dosimetrist generates the plan by using an optimization software to determine the beam characteristics required to deliver the final dose distribution (e.g. aperture shapes for each beam angle and energy for each aperture for conventional radiotherapy; or energy of the beams for each pencil beam spot for proton therapy).

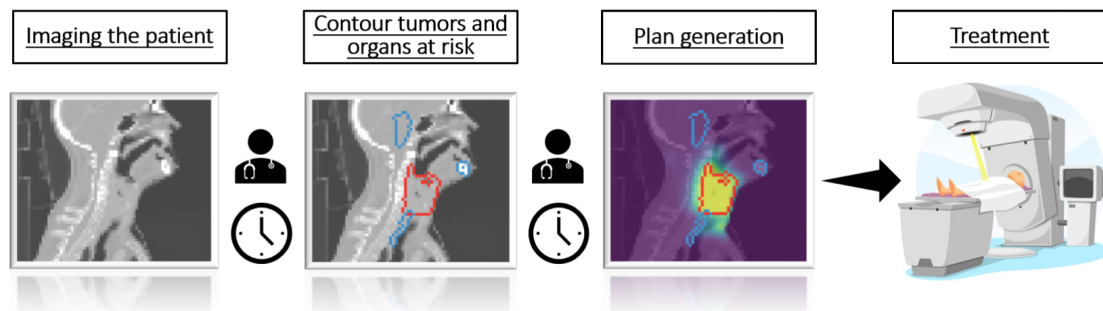


Figure 1.2: Current process to generate a treatment plan. The contouring of tumors (in red) and organs at risk (in blue) is done manually by a physician and takes time. The generation of the plan from these contours is tuned manually by a dosimetrist and takes time as well.

This whole process is labor-intensive, time-consuming and costly. For that reason, the treatment plan is generated once at the beginning and re-planning during the treatment only takes place if a physician finds that the deviation between planning and current situation is too large. Ongoing research is aiming at implementing a fully automatic process, as shown in Figure 1.3.

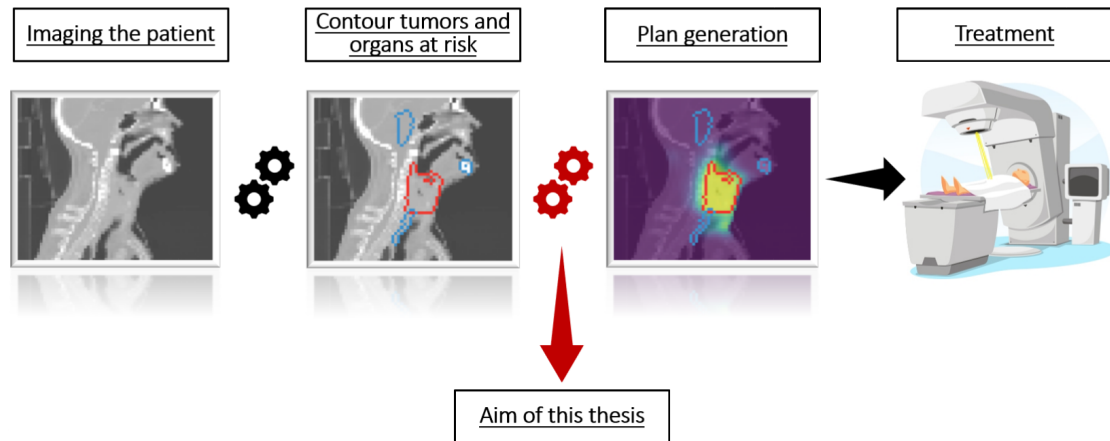


Figure 1.3: Fully automatic process to generate treatment plans. Both the contouring of volumes and plan generation would be done automatically. This thesis proposes a model to generate the plan from the contours.

That automatic workflow would allow one to generate a plan faster and therefore be able to re-plan more frequently. Having a plan specific to each session is the way to deliver the most accurate treatment. However, a typical generation of treatment plan takes several hours, which is too long to adapt the treatment on the fly. The longest parts are the ones involving human expertise: contouring and machine parameter adjusting to achieve the dose prescription. This thesis focuses on the second part for the case of head and neck. For now, no dose prediction model using neural networks have been proposed for protontherapy. The contribution of this thesis is to implement such a model and assess it.

The model predicting the dose distribution in the patient must be fast, accurate and output a dose achievable by the delivery system.

- Fast because if it is a burden that slows down the workflow, it would reduce the patient comfort, who already have to undergo treatment delivery every day for 30 days. The treatment cost would rise too as the time increase would lead to less patient treated in the day.
- Accurate so that it respects the clinical goals in terms of the dose delivered to the tumor and the sparing of organs at risk.
- Achievable as the prediction is not important in itself, it's the result from the system that tries to reproduce this prediction that will be delivered to the patient that is important. If the prediction is an ideal distribution but that can not be achieved, then the model is not valid either.

In this thesis, a deep learning model is proposed to meet these requirements. It has been trained on images of patients suffering from head and neck bilateral cancer. The document is divided as follows: Chapter 2 exposes the context and planning process in protontherapy. Chapter 3 consists of a review of literature. The following two chapters present the method and the results. Chapter 6 presents a discussion of the results and then a conclusion will close this study.

## 2 | Context

Proton radiation therapy is a form of external radiation therapy that uses protons produced by particle accelerators (such as a cyclotron or a synchrotron). Proton beams have physical advantages over photons: they deposit most of the dose at the end of their trajectory. Their range can be adjusted by modifying their energy in order to reach the target. High energy photons will deliver a maximum dose within few centimeters of the skin surface and continue to irradiate beyond the target. For targets deeper than three centimeters, each photon beam will deliver more dose to tissues around the target than in the target [2]. For this reason, photon therapy is generally delivered with a rotating gantry to shoot from multiple beam directions.

Protons behave differently, they enter the body, deliver a small, almost constant dose until their range where the dose deposition peaks, as shown in Figure 1.1. Beyond the Bragg peak, almost no dose is deposited. Thanks to their finite range, proton therapy completely preserves tissues out of the beam paths. As the Bragg peak of a monoenergetic beam is too narrow to cover the whole target volume, several beams of different energies are used to cover it with uniform dose (called SOBP, spread out Bragg Peak), Figure 2.1.

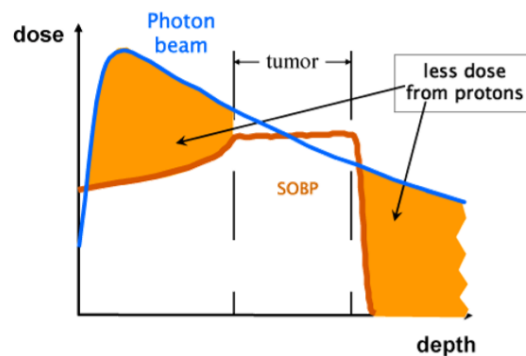


Figure 2.1: For protons, the SOBP curve shows how using multienergetic beam allows to deliver a uniform dose in the depth direction. Credits: [2].

Radiobiological effect (RBE) of protons does not bring a tremendous advantage compared to photons. Taking photons as the reference, the RBE of protons is 1.1, so it is really its physical characteristic that spurs interest [3].

Proton treatments currently fall into two categories: scattering and scanning. In scattering techniques, the cross-plane dose distribution is shaped by a customized aperture, while the distal conformity is controlled by a customized compensator. The aperture is generally made of lead and compensator of plastic material of variable thickness to modulate the range of protons. Using a compensator ensures that the highest beam energy applied to the patient will conform around the distal edge of the target along the beam trajectory, see Figure 2.2. However the proximal conformity cannot be controlled.

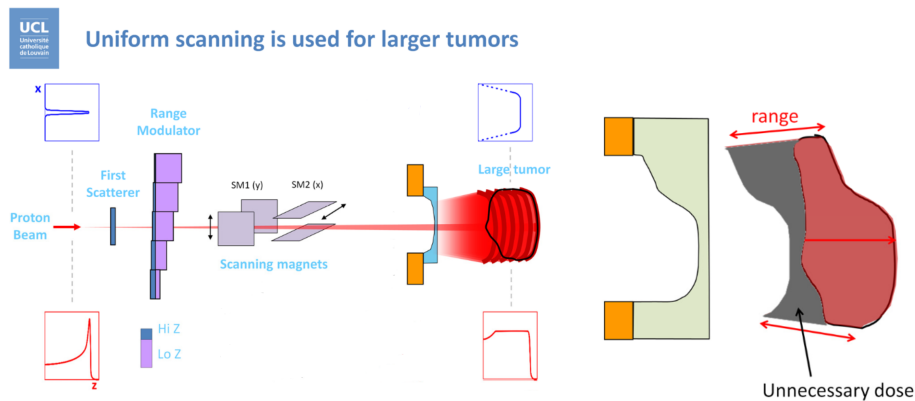


Figure 2.2: Scattering techniques to deliver protontherapy treatment on the left and illustration of the aperture and compensator on the right. The unnecessary dose deposition on the proximal area is depicted. Credits: [4].

The state-of-the-art and most flexible technique to deliver proton therapy treatment is the pencil beam scanning (PBS). It uses an electronically guided scanning system and magnets to place spots of dose and cover all over the tumor volume, as shown in Figure 2.3. As the beam position and depth are able to be tuned, both distal and proximal conformity can be controlled. Additionally, PBS does not require patient specific apertures or compensators.

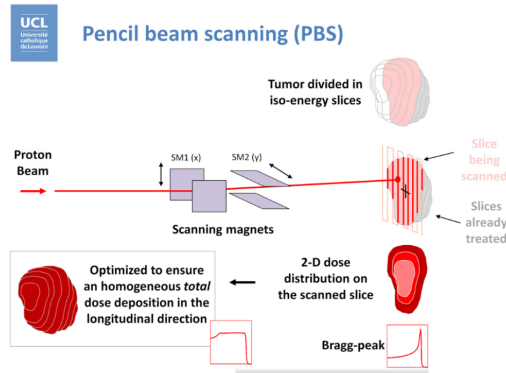


Figure 2.3: Pencil Beam Scanning technique to deliver protontherapy treatment. Credits: [4].

## 2.1 Treatment planning process

Conventional treatment planning takes place in successive steps. It starts from the scan of the patient, generally acquired from Computed Tomography (CT). The physician contours the tumor and the organs at risk to spare and sets dose prescription that needs to be achieved in the tumor. Then the dosimetrist finds the machine parameters to obtain a dose distribution that satisfies the prescription while sparing the organs at risk. These steps are detailed in this section.

### 2.1.1 Imaging the patient

Computed tomography produce cross-sectional images of the patient, as illustrated in Figure 2.4. The image is produced by measuring the attenuation of high energy photons when passing through the patient. A three dimensional image can be assembled from a large number of 2D images, each one being a slice. The intensity of each voxel is proportional to the attenuation coefficient of the tissues inside it and therefore to their electron density. The unit to quantify the intensity is the Hounsfield unit (HU) and ranges from about -1000 HU to 3000 HU. Water at standard temperature and pressure conditions corresponds to 0 HU, and air to -1000 HU [5]. For a voxel with attenuation coefficient  $\mu$ , the corresponding HU is given by:

$$HU = \frac{\mu - \mu_{\text{water}}}{\mu_{\text{water}} - \mu_{\text{air}}}$$

$\mu_{\text{water}}$  and  $\mu_{\text{air}}$  are respectively the attenuation coefficient of water and air.

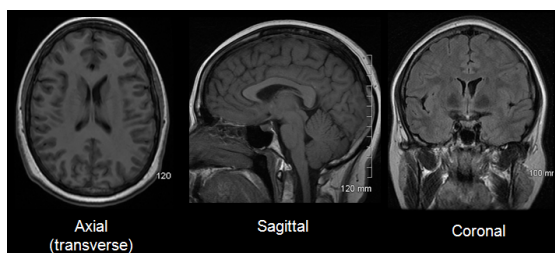


Figure 2.4: Axial, Sagittal and Coronal sections of CT scan of a head. Credits: [6].

The accuracy of the patient positioning and CT calibration for this first scan is really important as all inaccuracies will turn into systematic errors on the treatment delivery. For this reason, devices such as contention masks and laser points are used to position the patient, Figure 2.5. On-board imaging devices are also used in the treatment room to help achieve precise patient positioning.

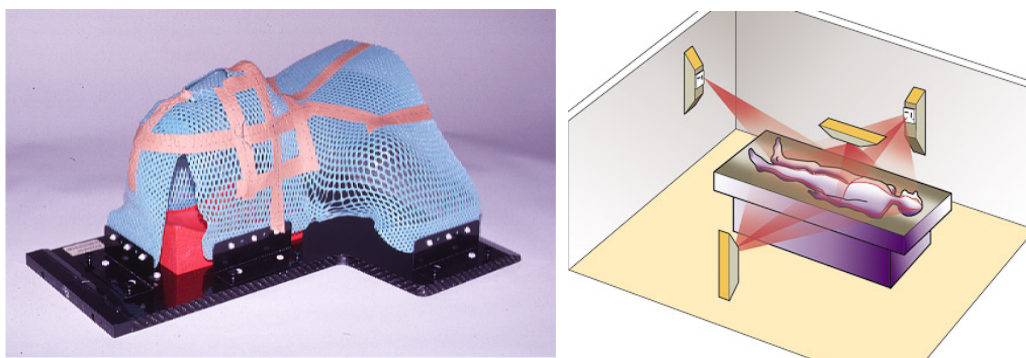


Figure 2.5: Left: patient specific contention mask. Credits: [5]. Right: laser system to position the patient. Credits: [7].

### 2.1.2 Contouring and dose prescription

Once the scan is acquired, a physician needs to contour the GTV, CTV and PTV, which are the target volumes, shown in the left-hand image of Figure 2.6.

- The GTV is the Gross Target Volume, or the macroscopic volume. It is the visible part of the tumor.
- The CTV is the Clinical Target Volume. It extends the GTV to encompass the microscopic extensions of cancerous cells. This step requires the expertise of the physician as the spread of cancerous cells depends on the location of the tumors and anatomical barriers.

- The PTV is the Planning Target Volume. It is a geometrical margin applied to the CTV to take into account various errors such as the one due to the motion or positioning of the patient.

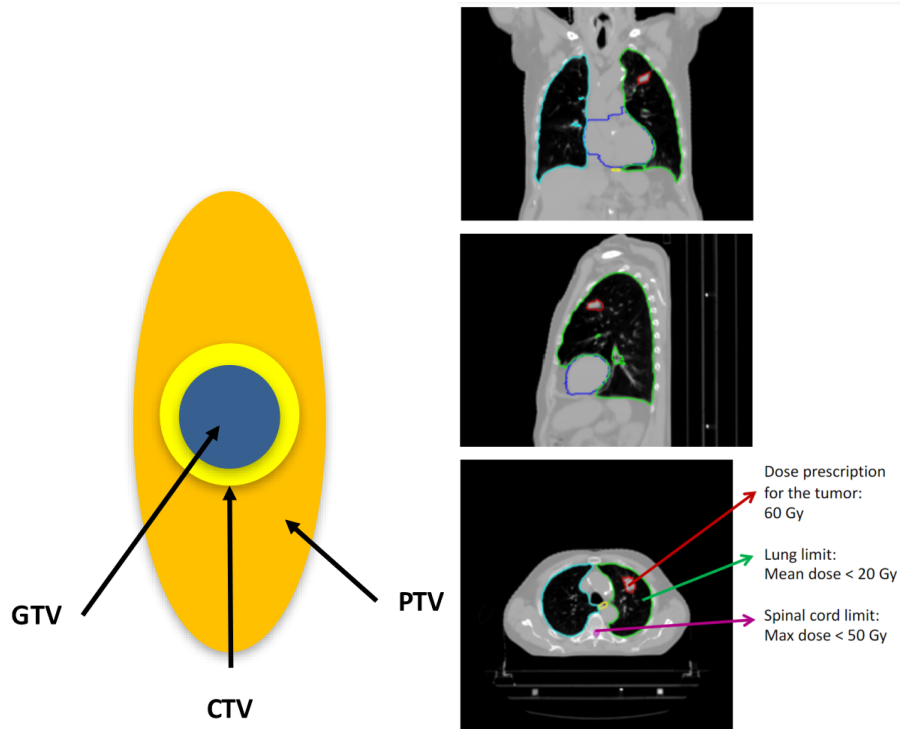


Figure 2.6: Left: Sketch of different contours. Right: Contouring and prescription of tumors and organs in the lung. Credits: [4].

With these contours, the dose prescription is specified, as shown for lungs on the right of Figure 2.6. For instance, in our head and neck case, a prescription of 50Gy in the CTV and a boost of 70Gy were really common. The different organs at risk, such as the spinal cord, the parotids, the eyes and others, need also to be defined in order to minimize the dose deposited in these volumes. This process of contouring takes several hours, and research is ongoing to automate it [8].

### 2.1.3 Plan generation

Minimum and maximum dose constraints are placed on organs. These lead to an optimization problem where typically the objective function is the sum of individual dose objective for each organ. The dosimetrist then optimizes the machine parameters to obtain a dose map that would respect as much as possible the constraints. It takes time as the tuning is an iterative process which may require the addition of constraints to guide the optimization. In the MIRO lab, the software RayStation was used. This software has been developed by RaySearch Labs in Sweden.

In an automatic treatment plan generation workflow, instead of tuning the parameters manually to generate a dose prediction and then the plan, the dose would be predicted directly from the patient scan and contours. Then the full treatment plan would be generated by inverse optimization.

### 2.1.4 Evaluation of Plans

Dose Volume Histograms (DVH) are used to evaluate the dose distribution in the target volumes as well as in the organs at risks (OARs). A DVH is a graphical representation of the dose in a given volume. The x-axis corresponds to the dose and the y-axis gives the percentage of volume.

DVH allows the determination of useful metrics such as DX: the dose delivered in X% of the volume of interest. For instance, metrics such as D95 and D5 are used for the CTV. D2 and  $D_{\text{mean}}$  can be computed for the OARs. D95 is used to evaluate the coverage of the target. D5 is for the maximum dose delivered to 5% of the CTV. D2 is used to evaluate the highest dose delivered in the 2% coldest volume of the organ (used for OARs).  $D_{\text{mean}}$  indicates the mean dose received by the region of interest.

Figure 5.5 is an example of cumulative DVH curves:

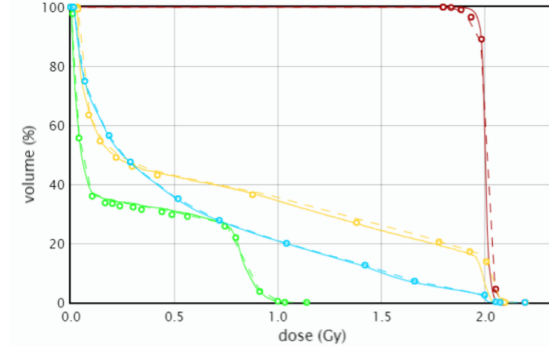


Figure 2.7: Dose-Volume histogram for a fraction of the treatment aiming 2Gy in the tumor, showed in red, the other curves correspond to the dose delivered to the OARs. Credits: [4].

## 2.2 Robustness

Robustness of a dose is achieved through the integration of uncertainties into the treatment plan optimization process. In conventional radiotherapy, robustness is achieved by expanding the CTV with a geometric margin. This margin is determined by considering uncertainties that arise during the treatment delivery process. Source of uncertainties are diverse such as the patient setup on the planning scan and on every treatment session, intrafraction patient motion (such as breathing), tumor contouring etc. The underlying assumption is known as "static dose cloud approximation". It states that the dose might shift due to the uncertainties but that the distribution will not change.

However, applying a simple geometric margin is inadequate for protontherapy [9]. This is due to the fact that the dose distribution varies significantly when the proton beam goes through different tissues. Therefore, robust plan in protontherapy are optimized on the CTV. From the nominal case, scenarios are defined and the plan is then optimized to respect the clinical prescription in all scenarios considered.

## 3 | Literature review

Dose prediction is a problem that has been investigated by the research community. Section 3.1 introduces what methods are currently used for automatic planning. Then more specifically published articles are presented of methods for proton therapy in Section 3.2. No deep learning models have been proposed yet in proton therapy. However, articles have proposed such models in conventional radiotherapy. They are detailed in Section 3.3.

Multiple architectures have been proposed. Most are based on convolutional neural networks, such as variations of the UNet architecture, but also with generative adversarial networks that are addressed in Section 3.3.4.

### 3.1 Current methods for automatic planning

To reduce the manual component of treatment planning, different methods have been developed. They can be divided into two categories: objective-based planning and knowledge-based planning.

Objective-based planning (OBP) relies on a set of objectives such as a list of DVH metrics to be reached for every region of interest. This method has been implemented among others by Pinnacle (Philips Radiation Oncology, Fitchburg, WI) [10].

Knowledge-based planning (KBP) relies on the fact that radiotherapy plans are generally quite similar for patients with similar geometries [11]. A KBP model uses a library of patients treated with the same modality, and tries to match the new patient to this library by using nonrigid registration. Once the patient is matched, the model computes a DVH metric.

These two methods present two limitations. First, they lack spatial information about the dose. Second, both OBP and KBP strategies still require significant

human intervention to define certain parameters needed to create the model, such as the target and OARs optimization goals for OBP [12] or handcrafted features that serve to match the actual patient to those in the library of patients for KBP [13],[14].

## 3.2 Dose prediction in proton therapy

To our present knowledge, no deep learning models have been developed to perform dose prediction in proton therapy yet. Two recent articles [15] and [16] approach the problem from another angle.

In [15], the authors try to predict the shape of two isodoses,  $V_{20\text{Gy}}$  and  $V_{5\text{Gy}}$  of the ipsilateral lung for breast cancer treatments. From their model, they generate expansions of the PTV for different margins ( $\text{PTV} + m$ ) and compare with the isodose volumes in the patient. From this they establish a relationship between the margin and the dose. Their method is used as a guide toward clinical goals but without taking into account the patient's specific CT image. Plus, it does not produce a 3D distribution.

The second one, [16], presents a method to quickly assess the benefit of treating a particular patient with proton therapy rather than conventional radiotherapy. The target audience are centers that do not have proton capabilities (neither the experience nor the treatment planning system) so that they can still make a well-informed referral decision. They are an alternative to a proposed workflow in Germany where a conventional radiotherapy plan and a protontherapy plan are requested before making the decision. This solution is both costly and time-consuming. In the paper, the authors developed a model that uses a library of patients treated with proton therapy, and tries to match the new patient to this library by using nonrigid registration. Once the patient is matched, the model computes a metric to determine whether an OAR would benefit from using proton therapy. Again, this method lacks spatial information about the dose.

## 3.3 Neural networks used for dose prediction in conventional radiotherapy

Neural network is a structure composed of several units, called neurons, that perform operations. These units are combined into layers and the links between neurons of subsequent layers are given weights, see Figure 3.1.

Training of the network indicates the model how to change its internal parameters

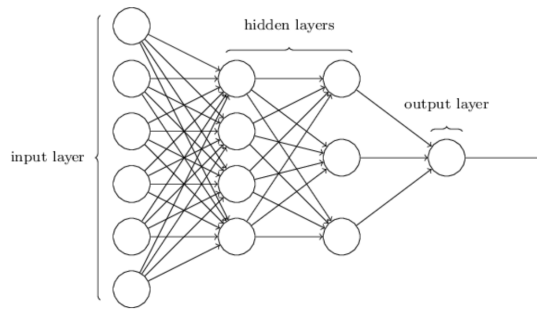


Figure 3.1: Simple illustration of basic neural network.

(such as the weights) in order to discover intricate structure in large data sets. Different type of neural network exist, depending on the architecture and operation performs by the units. Convolutional neural networks have brought about breakthroughs in processing images, video, speech and audio.

Neural networks have been implemented to perform dose prediction in conventional radiotherapy. Among these architectures, UNet has a prominent influence. The advantage of using deep learning is to avoid the reliance on hand-made low dimensional features, like in knowledge based planning techniques.

### 3.3.1 UNet

UNet is an architecture of convolutional neural network (CNN) that has been initially designed for image segmentation by Olaf Ronnenberger et al. in 2015 [17], as shown in Figure 3.2.

UNet is characterized by an encoding part where at each level the resolution of the image is decreased, by a pooling operation or a stride in the convolution operation. Then follows a decoding part where the size of the image is increased again. The advantage is that it allows to get information of the features (using convolutions) at different resolution of the image. Skip connections can also be introduced to reduce the problem of vanishing gradient by injecting features from the encoder and concatenate them with output of upsampling layers. This architecture has proved to be able to perform well even on small training datasets, which is a common problem in medical application.

This architecture has been used by Nguyen et al. 2017 [18], then improved in 2019 [19], by Kearney, Chan et al. 2018 [20] with some variations on the pooling operations, and also by Barragan-Montero et al. in [21].

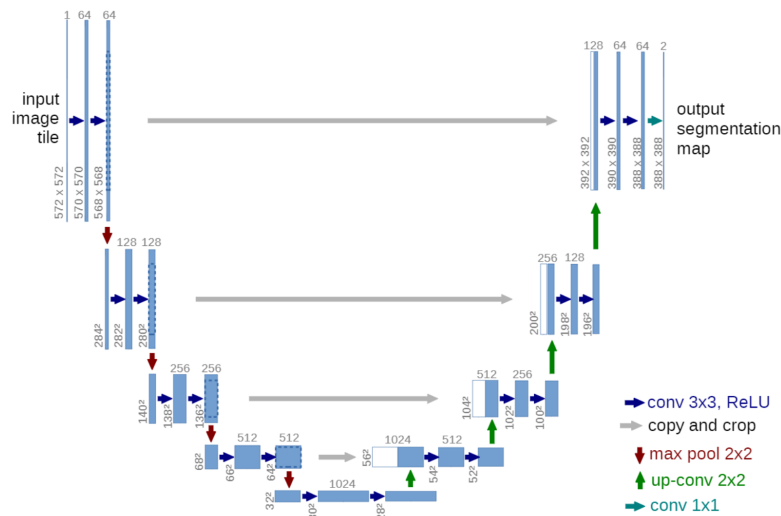


Figure 3.2: Illustration of the UNet from the article of O. Ronnenberger et al. [17]

### 3.3.2 Other CNN architectures

Chen et al. propose in [22] to train two ResNets. One is trained with the prescription for the CTV and organs at risk as input and the other with the input modified to add the beam configuration and then the two are combined. Similarly to the UNet network, they use batch normalization to normalize the features and skip connections to counter the vanishing gradient problem. However in their discussion, they acknowledge that they did not investigate different network architectures. Therefore there was no proof of the superiority of this network compared with UNet. Plus, training two networks is computationally heavy.

Another article where ResNet was also used has been proposed by Fan et al. early 2019 [23]. The architecture is actually composed of two ResNets that are contracting the information. One for the PTV and one for the OARs whose outputs are concatenated and fed to an *anti-ResNet*, which is expanding. The overall seems similar to UNet, with residual connections.

More recently, another mix between UNet and ResNet has been published from Liu et al. 2019 [24].

The similarities between these articles and the UNet show a convergence toward the same type of network.

### 3.3.3 2-D versus 3-D models

Several models perform a prediction slice by slice such as in Nguyen et al. 2017 [18] which predicts the dose for prostate with IMRT or in Chen et al. 2018 [22] for nasopharyngeal cancer. Predicting slice by slice is computationally less intensive than using 3D matrices as the dimension of the input is greatly reduced. However, craniocaudal information is lacking and causes errors at the superior and inferior slices of the tumor and organs contoured.

In Fan et al. 2019 [23], the authors try to mitigate the problem by adding three slices above and below the central slice that is to be predicted.

### 3.3.4 Generative Adversarial Network

An interesting and different approach is to use a generative adversarial network (GAN) to predict the dose [25]. A GAN involves the training of two networks: a generator that performs a task, such as dose prediction and a discriminator that determines whether its input is the real dose coming from the dosimetrist or a generated one. The generator in [25] is a UNet too. Whereas the discriminator is a CNN that take the prediction and outputs a scalar value. During its training, the discriminator receives a ground truth dose and a generated one, then updates its parameters and so on. The GAN method yields very good result. The authors compared it with the UNet of Nguyen et al. [18] that was predicting slice by slice the dose.

## 4 | Methods

### 4.1 Network Architecture

The chosen architecture is similar to the hierarchically densified U-Net (HD-UNET) shown in Figure 4.1, presented in [19].

UNet has the characteristic to be able to learn on very small dataset, due to its symmetry of the encoding and decoding layers. Since the number of patients in this study is relatively small, this characteristic of UNet is an advantage. Furthermore, the three dense operations that have been implemented in the HD-UNET (dense convolution, dense downsampling and dense upsampling shown in the bottom illustration in Figure 4.1) offer a form of residual connections. Indeed, in the dense convolution operation (A on Figure 4.1), the layers are concatenated with the output of the convolution. This architecture is therefore the one reconciling different approaches seen previously.

It is a 4-layer UNet which outputs the 3D dose prediction. The input channels contain the dose prescription on the target volumes, the binary mask of organs at risk (1 inside the volume and 0 outside) and an additional channel adding tissue density information, this will be detailed in Section 4.5. To allow the network to learn efficiently and not being confused by the different resolution of images, coming from different centers and CT scanners, all images were rescaled to voxels of 5\*5\*5 mm.

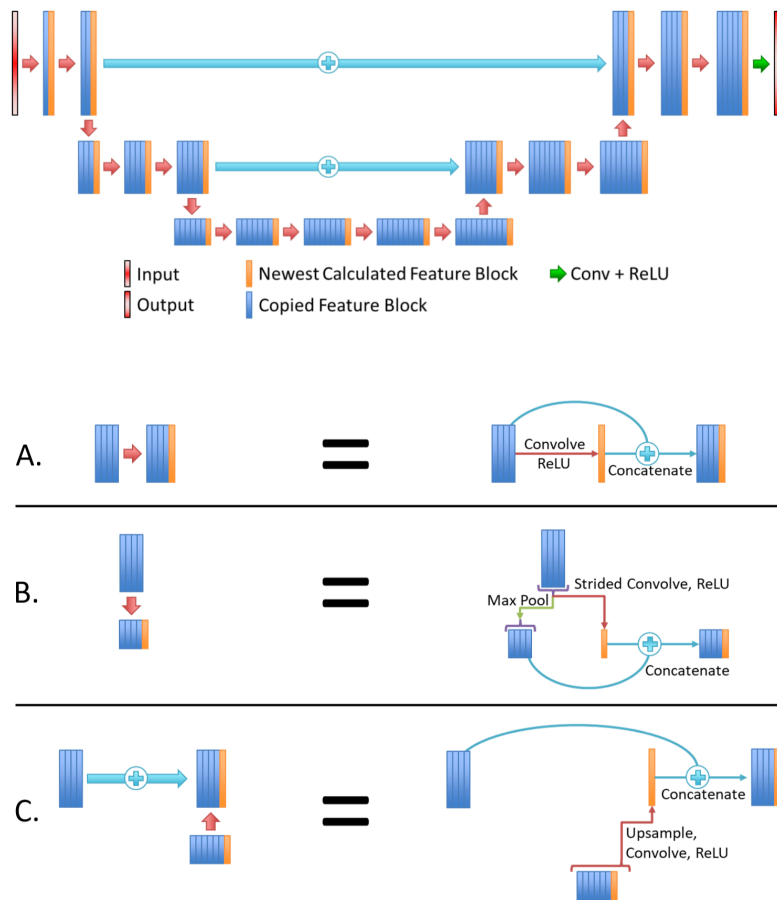


Figure 4.1: Top: General architecture of the H-D U-Net. Bottom: A is the dense convolution operation, represented by a red arrow on the architecture schematics on the top. B is the dense downsampling operation, represented by an orange facing down arrow. C is the dense upsampling operation where layers from the skip connection are concatenated with the convolved lower level layers. Credits: [19]

### 4.1.1 Cost function

The dose prediction problem requires a regression model. Indeed, the aim is to predict the value of each voxel in the dose distribution. The main cost functions used for regression problems are the mean squared error, or L2 norm and the mean absolute error, or L1 norm.

The **mean square error**, MSE, is generally the one used by default when training deep neural networks, but its main disadvantage is that outliers have a large influence. Indeed, the larger the error, the more importance is granted in the metric.

In this case, typical outliers would be for instance CTs that present artifacts.

The **mean absolute error**, MAE, is the sum of absolute differences between the groundtruth and predicted voxel values. MAE presents another major drawback, the slope is constant throughout the x axis, which means that the gradient will be large even for small loss values. This increases the risk to miss the optimum and therefore reduces the precision at the end of training. This problem is mitigated by using a dynamic learning rate that would reduce the step size when getting closer to the minima (such as using an Adam optimizer detailed in Section 4.1.5).

If the outliers can be removed, or their presence must be highlighted, then MSE is recommended. However, if outliers are due to noisy data and should be ignored, the MAE is a good approach. These two loss functions will be tried out to find the best fitting for this application.

### 4.1.2 Activation function

Typical activation units are the rectified linear units (ReLU). They are easy to optimize as they are similar to linear units. The difference with linear units is that ReLU have half of their domain outputs zero. This makes the derivative through a rectified linear unit remain large and consistent whenever the unit is active.

There are some generalisations of rectified linear units that have been proposed. The drawback of classical ReLU is that half of its domain is set to zero. Therefore the model is unable to learn via gradient-based methods from examples with zero activation. Goodfellow et al. in [26] state that most of these generalizations perform comparably to rectified linear units and occasionally improve the performance. Therefore, the classical ReLU activation function was kept.

### 4.1.3 Pooling

Pooling in a UNet allows to reduce the dimensionality of the data passed to a lower level. Therefore, UNet learns features on multiple resolution. Max pooling and Average pooling, shown in Figure 4.2, are two popular methods used for CNN. Max pooling applies a filter to the image and extracts the highest value pixel over the filter. Average pooling outputs the average value of all pixels considered in the filter. The choice of the pooling function depends on the application. Average pooling brings excellent results in classifying the Caltech 101 dataset [5]. Max pooling has been successfully applied to the training of a deep CNN for the imageNet competition [7]. And as pointed out in [27], the two methods have drawbacks, and a stochastic combination mitigates this downside. Max pooling only considers the

maximum element and ignores the others in the pooling region. This can lead to the vanishing of distinguishing features when most of the elements the pooling region have a high magnitude.

Whereas Average pooling will have an effect of smoothing as all the low magnitude elements will be considered. Therefore, the contrast of features in the pooling region will be reduced.

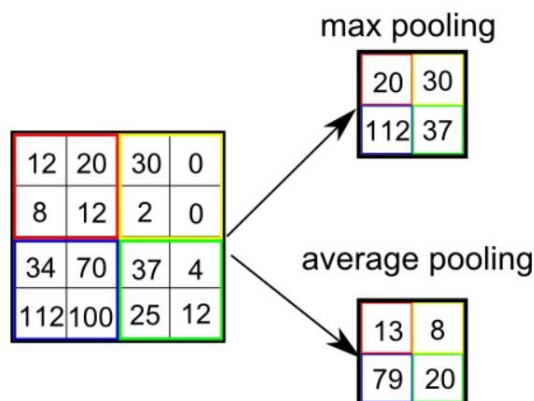


Figure 4.2: Illustration of Max Pooling and Average Pooling operation. Credits: [28]

Both pooling methods will be tested out in the cross-validation.

#### 4.1.4 Regularization

The dropout layer [29], is the regularization method used in Nguyen’s model [19]. Dropout provides a computationally inexpensive but powerful method of regularizing models. Dropout is similar to bagging of several large neural networks. Bagging involves running multiple models and evaluating them on each test sample. It seems impractical to implement it when each model is a large neural network since training takes several hours. Dropout palliates this issue by training sub-networks that are formed by deactivating non-output units from an initial network. This is performed by setting the weights that multiply the output of the unit targeted to zero.

#### 4.1.5 Optimizer

Optimization algorithms that use the entire training set are called batch methods. The parameters are updated after processing all the training examples. Batch

methods are an accurate way to optimize but computing the cost and gradient for the entire training set can be very slow and sometimes intractable on a single machine when the dataset is too big to fit in main memory.

Stochastic gradient descent (SGD) addresses this issue by following the negative gradient of the objective after seeing only a single or a few training examples. A problem remains: choosing a proper learning rate can be difficult [30].

Several optimizer functions with adaptive learning rate have been proposed to solve this issue. Among them stands Adam algorithm (Kingma and Ba, 2014), the name stands for *adaptive moments*. Adam is a method with dynamic learning rate, which means that the step size will be adapted according to the gradient value. Momentum is computed to take past gradients into account to smooth out the steps of gradient descent. The momentum is incorporated directly as an estimates of the first-order moment (with exponential weighting) of the gradient.

The Adam algorithm is as shown in Figure 4.3, coming from the deep learning book of Goodfellow [26]

---

**Algorithm 8.7** The Adam algorithm

---

**Require:** Step size  $\epsilon$  (Suggested default: 0.001)  
**Require:** Exponential decay rates for moment estimates,  $\rho_1$  and  $\rho_2$  in  $[0, 1)$ . (Suggested defaults: 0.9 and 0.999 respectively)  
**Require:** Small constant  $\delta$  used for numerical stabilization (Suggested default:  $10^{-8}$ )  
**Require:** Initial parameters  $\theta$   
Initialize 1st and 2nd moment variables  $\mathbf{s} = \mathbf{0}$ ,  $\mathbf{r} = \mathbf{0}$   
Initialize time step  $t = 0$   
**while** stopping criterion not met **do**  
  Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding targets  $\mathbf{y}^{(i)}$ .  
  Compute gradient:  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$   
   $t \leftarrow t + 1$   
  Update biased first moment estimate:  $\mathbf{s} \leftarrow \rho_1 \mathbf{s} + (1 - \rho_1) \mathbf{g}$   
  Update biased second moment estimate:  $\mathbf{r} \leftarrow \rho_2 \mathbf{r} + (1 - \rho_2) \mathbf{g} \odot \mathbf{g}$   
  Correct bias in first moment:  $\hat{\mathbf{s}} \leftarrow \frac{\mathbf{s}}{1 - \rho_1^t}$   
  Correct bias in second moment:  $\hat{\mathbf{r}} \leftarrow \frac{\mathbf{r}}{1 - \rho_2^t}$   
  Compute update:  $\Delta \theta = -\epsilon \frac{\hat{\mathbf{s}}}{\sqrt{\hat{\mathbf{r}} + \delta}}$  (operations applied element-wise)  
  Apply update:  $\theta \leftarrow \theta + \Delta \theta$   
**end while**

---

Figure 4.3: Adam algorithm, pseudo-code

Adam algorithm was used in the architecture.

## 4.2 Data

The dataset used is composed of 64 patients with head and neck cancer. The proton plans have been generated by a dosimetrist from the CT scans and the contours drawn by physicians. All these patients have bilateral target volumes and were planned with the same beam configuration. Since the dataset is relatively small, selecting only bilateral cancers and same configuration adds some consistency and allows the model to learn useful information. In order to protect the privacy of patients, all the data have been anonymized.

The advantageous symmetry of the head and neck anatomy makes it possible to double the dataset by flipping along the craniocaudal-anterior posterior plan. Many articles presented in Section 3.2 that use neural networks to perform dose prediction in conventional radiotherapy, were using other techniques of data augmentation. Techniques such as non-rigid deformations, noise, and size dilation work well for conventional radiotherapy. However, as protons range depend on the tissues they go through, a non-rigid deformation of the organs will not translate into the same deformation of the dose grid.

In the dataset used for this study, there were always two prescriptions for the target: prescription of 50Gy or 54.25Gy for the CTV and a boost of 70Gy prescribed on a smaller volume.

Thirteen OARs were included as input of the model. These organs were the brainstem, spinal cord, esophagus, larynx, mandible, optic cavity, right and left parotid glands, pharynx, right and left submandibular glands and another mask for the whole body.

The data was exported from RayStation in DICOM (digital imaging and communications in medicine) [31]. DICOM is an international standard introduced to ease the exchange and processing of medical images. Before the image data itself, a header contains several useful information such as the position of the patient, orientation, physical resolution of voxels etc. The implementation of the network being written in Python, a pre-processing step was introduced to transform DICOM images into NumPy arrays.

## 4.3 Training

Four Nvidia GPUs were available in the lab to perform the training and testing of the model. A GeForce GTX Titan with 6GB of memory and three Titan Xp with

12GB of memory each. GPUs were initially designed for video games rendering, task that needs millions of operations performed in parallel. GPUs are used to parallelize simple independent operations such as matrix multiplication to convert 3D coordinates to 2D coordinates. Similarly, neurons in the same layer are independent from each other, therefore they benefit from GPUs' parallelism. GPUs offer also the great advantage over CPUs to provide a larger memory bandwidth [26].

The data was split into subsets of 51 patients for the training and 13 for testing. The training set had been randomly subdivided into a validation set of 11 patients and do a 5 fold cross-validation to determine which method performs better.

## 4.4 Patching system

Ideally, the network should be fed with the full 3D matrices and trained directly. Unfortunately, the input composed of the CT, the CTV and all the OAR masks necessitates extensive computational resources. Therefore, a subset of the matrices, called patches, of size  $96 \times 96 \times 32$  pixels are fed to the network. They are sampled according to the probability that they contain part of the CTV in order to avoid a large number of patches on the sides. All the different patch selections are then merged to produce the full 3D dose distribution. The merging takes place by multiplying each patch prediction with a gaussian filter element-wise. The gaussian filter has a standard deviation of 2.5 and a mean of 0. Then the patches are summed up and weighted according to the gaussian contribution received.

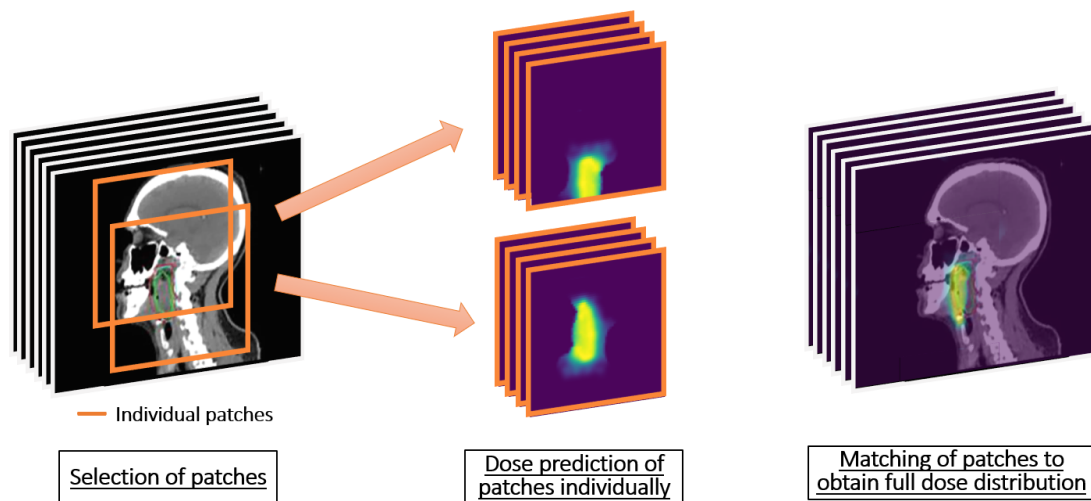


Figure 4.4: Sub-matrices, called patches to feed the network.

## 4.5 Adding information about the traversed tissues

The difference between predicting dose distribution in conventional radiotherapy and in proton therapy is that protons have a finite range that depends on the stopping power ratio of tissues through which they travel. Indeed, the protons have an initial energy and they interact with matter. They mostly lose their energy by transferring it to matter through inelastic collision causing dose deposition. Protons are also mildly scattered through elastic collisions and sometimes absorbed by nucleus reaction.

In protontherapy, tissue heterogeneity will therefore change the location of the Bragg peak in the patient compared to homogeneous tissues. In conventional radiotherapy, heterogeneity of tissues modifies the attenuation of photons. Therefore the deposited dose will be but increased or reduced but the effect is not as drastic as proton therapy where most energy is deposited in one peak as shown in Figure 1.1.

The stopping power is defined as the average energy lost per unit of distance along the track of the particle. In proton therapy, the quantity used is the stopping power ratio.

$$SPR = \frac{SP_{medium}}{SP_{water}}$$

The reason is that the devices are calibrated with ranges in water. The range in a medium is therefore obtained by:

$$R_m = R_w \frac{S_w}{S_m} = R_w \frac{\rho_w \left(\frac{S}{\rho}\right)_w}{\rho_m \left(\frac{S}{\rho}\right)_m} = \frac{R_w}{SPR}$$

Hence it is important to give the information about the stopping power of tissues to the network.

For this reason, six experiments with different inputs have been conducted: *No CT*, *CT*, *WEPL*, *Hybrid*, *CT masked*, *CT masked + Hybrid*. They are described in the following sections.

### 4.5.1 *No CT*: without Computed Tomography scan

No information is added about the tissues electronic density in this experiment. The only input, shown in Figure 4.5, are therefore the CTV contours with prescription and boost and the binary masks of the OARs.

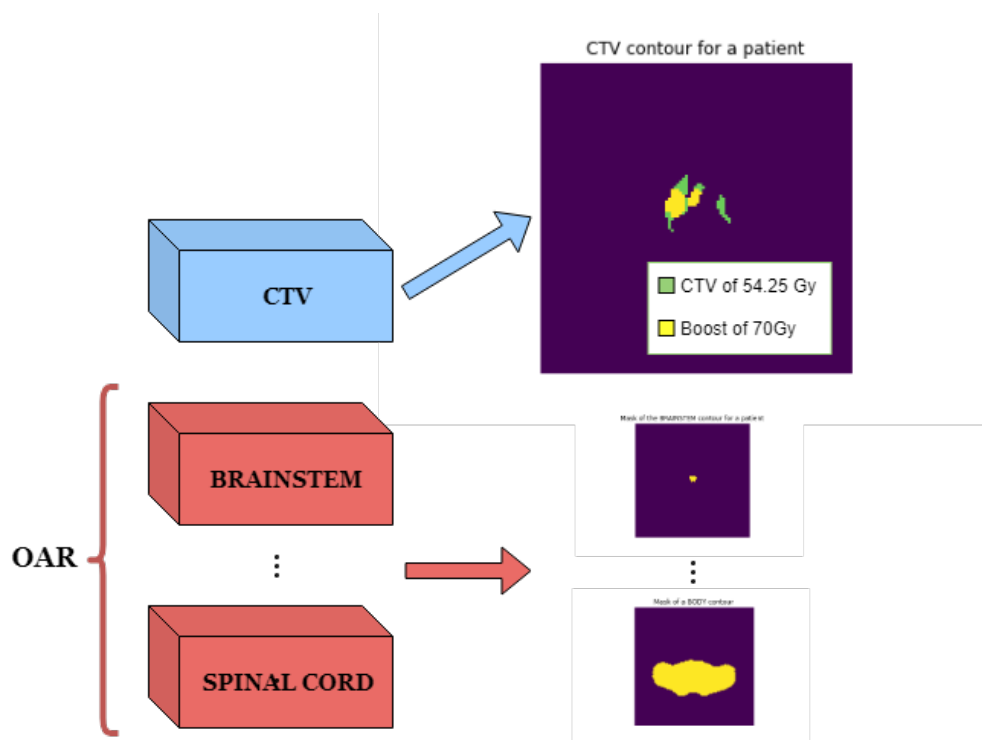


Figure 4.5: As input 14 channels with 3D matrices containing the CTVs and the 13 OAR masks

### 4.5.2 *CT*: with Computed Tomography scan

In this second experiment, a channel has been added containing the CT (Computed Tomography) image of the patch considered. The new input is shown in Figure 4.6. Stopping power of tissues can be determined from Hounsfield units in the CT scan by well known calibration methods.

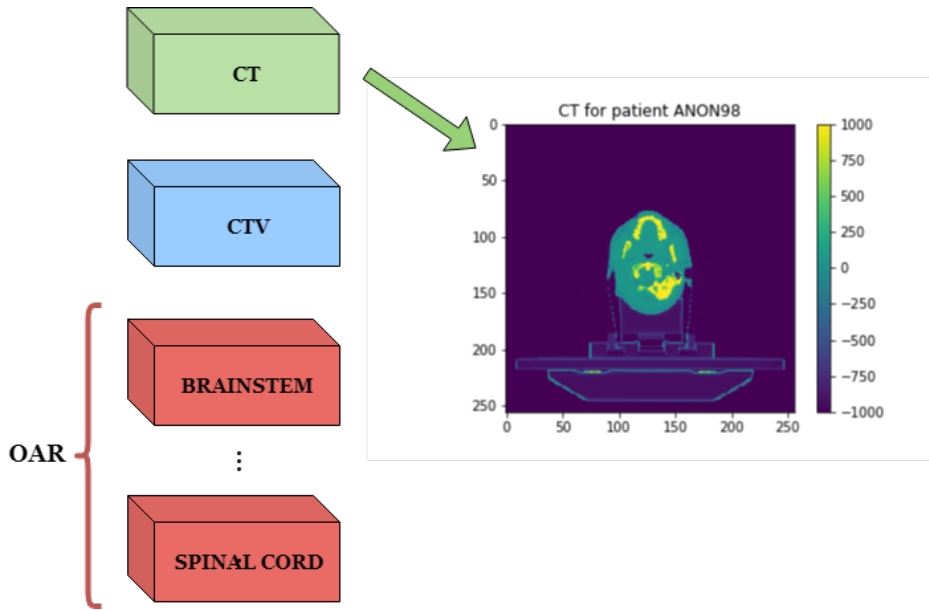


Figure 4.6: As input: 15 channels with 3D matrices containing the CTVs, the 13 OAR masks and the CT image

### 4.5.3 *WEPL*: Water Equivalent Path Length image

In this experiment, instead of adding the CT image, a conversion into WEPL is added.

The range of the proton to reach a given voxel of the CT image is obtained by measuring the distance between this voxel and the border of the image along the beam trajectory. To know the energy required to reach that voxel, the treatment planning system needs to convert this physical range into a WEPL. First, the CT scan is converted in water-equivalent quantities, using the HU to stopping power ratio calibration curve. The WEPL corresponding to a given voxel is computed with a ray-tracing algorithm. To do so, the distance traveled in each crossed voxel is weighted by the voxel stopping power ratio and accumulated.

As shown in Figure 4.7, each beam goes through a range shifter made of Lexan in the RayStation plans, to shift the Bragg peak closer to the surface of the patient. It should be taken into account in the WEPL computation.

Gantry angles from which the beam is delivered are integrated to compute the WEPL, as well as the couch angle. The gantry angles are 60, 120, 240 and 300 degrees. The couch angle is of 10 degrees clockwise for the two first gantry angles

and 10 degrees counterclockwise for the two others. The couch is rotated to limit the irradiation of the shoulders. The setup is illustrated on Figure 4.7:

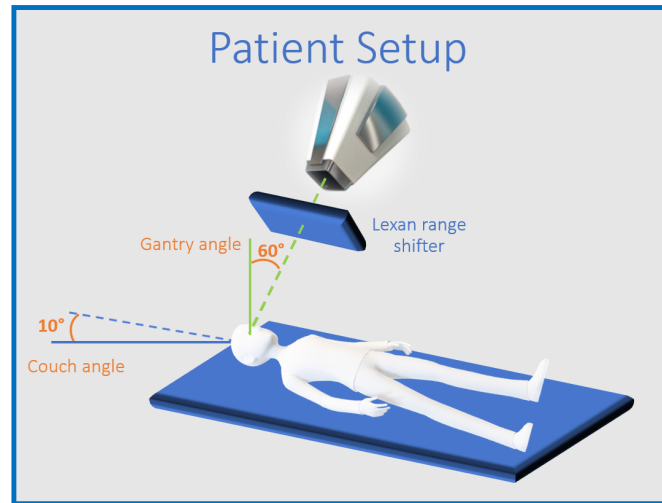


Figure 4.7: Patient setup

The inputs of this experiment are the WEPL on the beam trajectories, the CTV and the OARs, as shown in Figure 4.8. The WEPL image has been computed using the MIROpt software. It is a treatment planning system implemented in the lab by Ana Barragan.

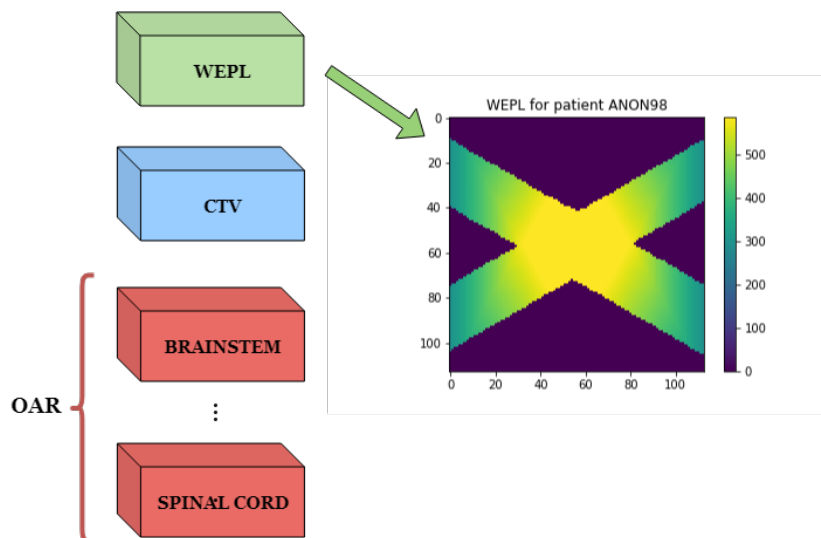


Figure 4.8: As input: 15 channels with 3D matrices containing the CTVs, the 13 OAR masks and the CT image converted into a WEPL image

#### 4.5.4 *Hybrid: WEPL + dose prescription*

To help improve the results, another experiment was conducted. The WEPL image is modified by replacing the values of voxels inside the CTV by the CTV prescription, as shown in Figure 4.9. The motivation of this method is to force the network to seek for an uniform dose over the CTV.

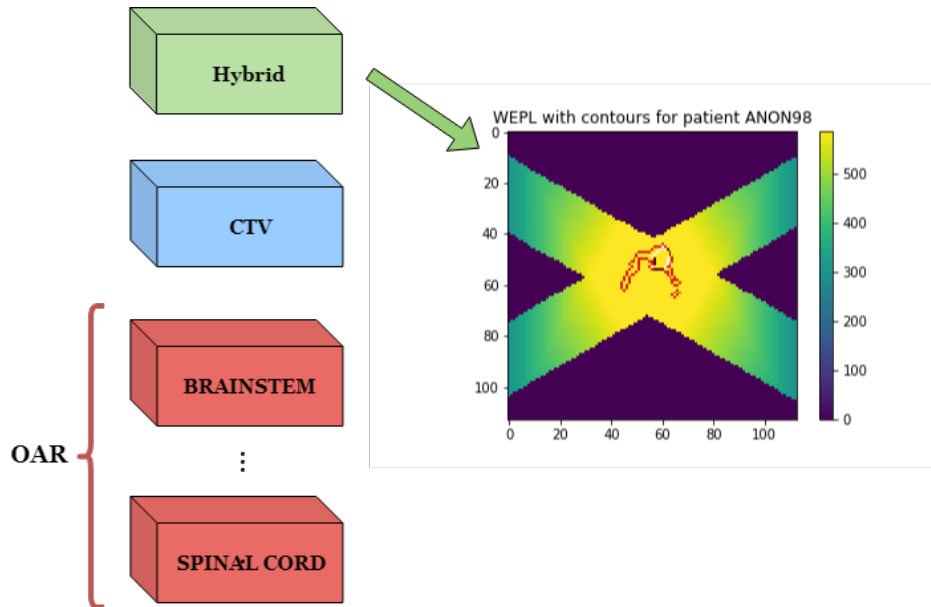


Figure 4.9: As input: 15 channels with 3D matrices containing the CTVs, the 13 OAR masks and the WEPL image where the value of voxels inside the CTV have been replaced by the CTV prescription

#### 4.5.5 *CT masked: CT cropped on beams trajectory*

Since only a set of four beams were used, it is possible to remove unnecessary information by adding the prior knowledge that there is only dose deposition on the beam trajectory. This has been done by masking the CT image with the beam mask, shown in top image of Figure 4.10. The four beams aim at the isocenter, generally the center of the CTV. The beam mask itself has been computed by rotating the CTV of the couch and gantry angle, projecting the shape and then rotate again. The beam mask is shown in the bottom image of Figure 4.10.

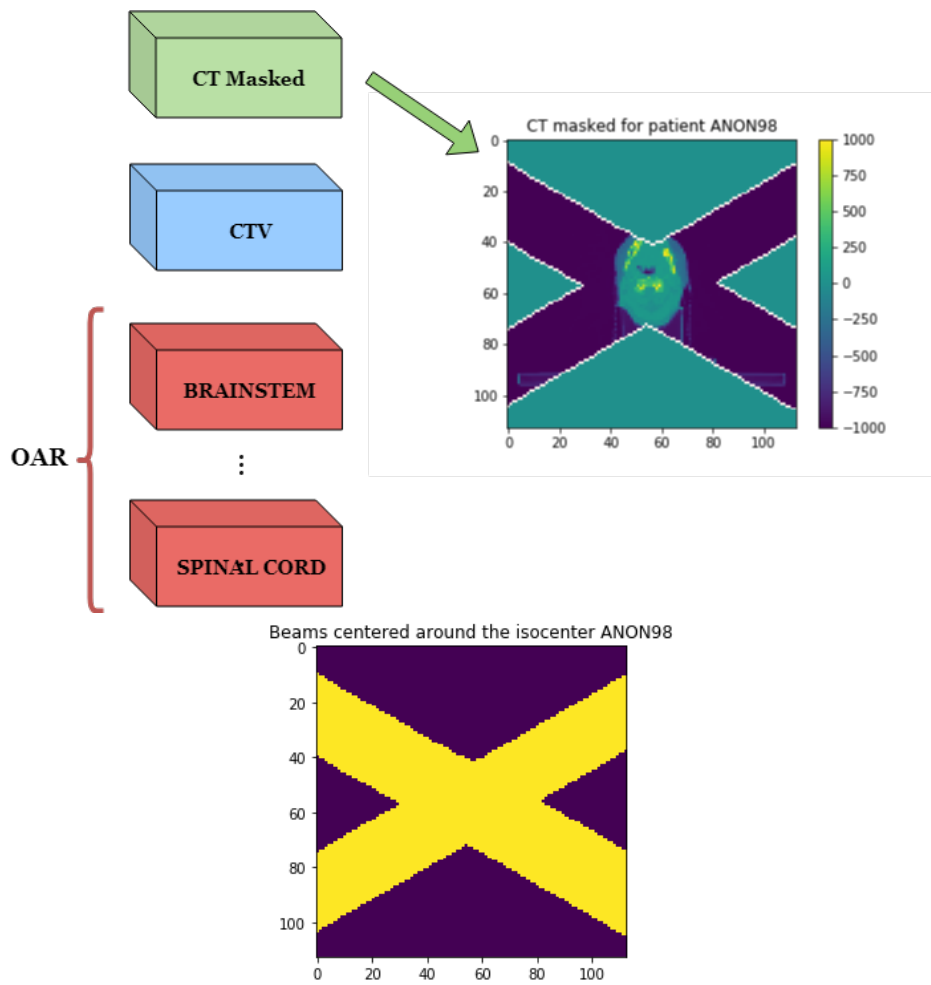


Figure 4.10: As input: 15 channels with 3D matrices containing the CTVs, the 13 OAR masks and the CT image masked to leave only the part where the beam go through the body.

#### 4.5.6 *CT masked and Hybrid: CT masked + Hybrid*

This last experiment was conducted to observe whether the network benefits for receiving two channels: *CT masked* and *Hybrid* in addition to the CTV and the OARs. The motivation for this method is to observe whether the network learns better when receiving the information on the stopping power from these two sources. The input is shown in Figure 4.11.

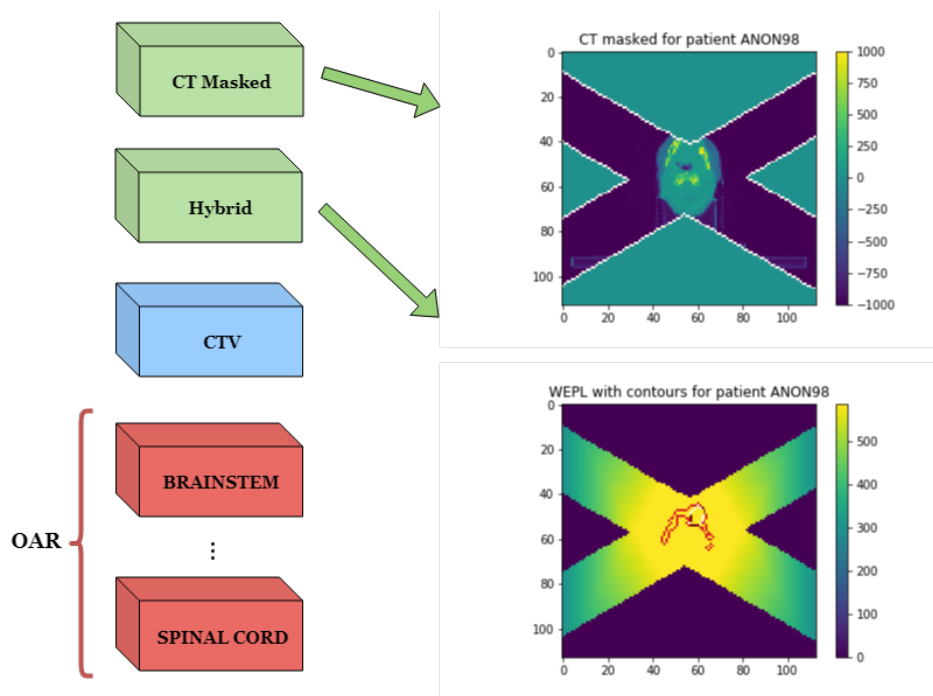


Figure 4.11: As input: 16 channels with 3D matrices containing the CTVs, the 13 OAR masks and the CT image masked and the Hybrid (WEPL where the CTV region has been replaced by dose prescription)

## 4.6 Dose Mimicking

Dose mimicking consists of optimizing automatically the beam parameters to obtain the closest dose distribution possible to the prediction. Dose mimicking is the step required to obtain go from the dose prediction to the treatment plan without the long hours of manual tuning.

To do so, a script was run on RayStation to extract the isodoses from the dose prediction. RayStation then optimizes automatically the weights of the beams and places the pencil beams spots in order to reach these isodoses.

# 5 | Results

## 5.1 Model Validation

The validation of the model was performed using a 5-fold cross validation. From the 64 patients, 13 were kept for testing. For each fold, the remaining 51 are split into a training set of 40 patients and a validation set of 11 patients randomly. The generalisation error is the mean error over all the folds and all patients.

The generalisation error  $E_{\text{gen}}$  is defined as:

$$E_{\text{gen}_{m_i}} = \mu_j \text{MSE}_{m_i,j} \quad \forall j = 0, \dots, 4$$

Where  $m_i$  corresponds to a model among:  $\{NoCT, CT, CT \text{ masked}, WEPL, Hybrid, CT \text{ masked} + Hybrid\}$ .  $m_{i,j}$  represents the corresponding model with the best validation loss in the  $j^{th}$  fold of the cross validation.  $E_{\text{gen}}$  is the mean of the MSE between the prediction and ground truth for each fold with different inputs on the testset. As a first test, the architecture in [19] was tested with the different inputs:

Method	$E_{\text{gen}} \mu \pm \sigma$
No CT	$7.92 \pm 4.38$
CT	$6.88 \pm 3.49$
WEPL	$5.67 \pm 0.60$
CT masked	$2.13 \pm 0.52$
Hybrid	$5.43 \pm 1.07$
CT masked and Hybrid	$19.26 \pm 0.73$

Table 5.1: Mean and standard deviation of the generalization error is displayed.

Table 5.1, shows that the best results are reached with the *CT masked*. Therefore, this method was selected to perform further cross validation. The aim is to investigate whether the architecture modifications considered in Section 4.1 improve the generalization error.

Table 5.2 shows the generalization error when an Average pooling is used instead of the classical Max pooling. In the third column, a trial with MAE as loss function and Max pooling was also tested.

Method	$E_{\text{gen}} \mu \pm \sigma$
CT masked MaxPooling	$2.13 \pm 0.52$
CT masked AveragePooling	$2.49 \pm 0.67$
CT masked MAE	$3.31 \pm 0.80$

Table 5.2: Mean and standard deviation of the generalization error is displayed for the different variation of UNet and the CT masked input

According to the machine learning framework, the best method is the *CT masked* with the max pooling and the MSE loss. MSE is used as loss function since it has the advantage of being differentiable and therefore allows gradient propagation. However, the MSE metric is only a surrogate of the clinical metrics that are interesting for the application such as D95 and D5 for the CTV and D2,  $D_{\text{mean}}$  for the OARs.

Let us define another generalisation error:

$$E_{\text{gen}^*_{m_i}} = \mu_i(\Delta D95_{\text{ctv}_{m_{i,j}}} + \Delta D2_{\text{brainstem}_{m_{i,j}}} + \dots + \Delta D_{\text{mean}_{\text{esophagus}_{m_{i,j}}}}) \quad \forall j = 0, \dots, 4$$

$E_{\text{gen}^*}$  is the mean of the different DVH metrics relevant to ensure that the target is properly covered and that the dose in the organs at risk is limited.

Table 5.3 presents the mean of difference of DVH metrics between the ground truth and the prediction.

Method	$E_{\text{gen}^*} \mu \pm \sigma$
No CT	$4.25 \pm 4.44$
CT	$4.43 \pm 5.07$
WEPL	$4.25 \pm 5.75$
CT masked	$3.85 \pm 3.52$
Hybrid	$4.43 \pm 4.87$
CT masked and Hybrid	$4.22 \pm 4.82$

Table 5.3: Mean and standard deviation of the new generalization error is displayed.  $E_{\text{gen}^*}$  is computed as the mean of the difference of DVH metric between the prediction and ground truth for each fold with different inputs on the testset

## 5.2 Results considering clinical metrics

The results over test set have been displayed using box plot to ease the visualization. Upper and lower boundaries of each box represent the 75th and 25th percentiles respectively, and the horizontal line in the box depicts the median. Whiskers extend to 1.5 times the interquartile range. Dots outside boxes represent outliers. The mean is represented by a x symbol. They will be analysed in Chapter 6, Discussion.

### 5.2.1 Target volume

In Figure 5.1, the three graphics represent the difference of D95 metric for the two dose prescriptions (CTV and CTV boost) and the difference of D5 for the boost. The differences are expressed in percentage of the highest prescribed dose (70Gy). D95 shows the mean dose delivered to 95 % of the CTV volume. As there were two dose prescriptions in the CTV, the two metrics have been separated. D5 shows the eventual hotspot that would be present in the boost region of the CTV.

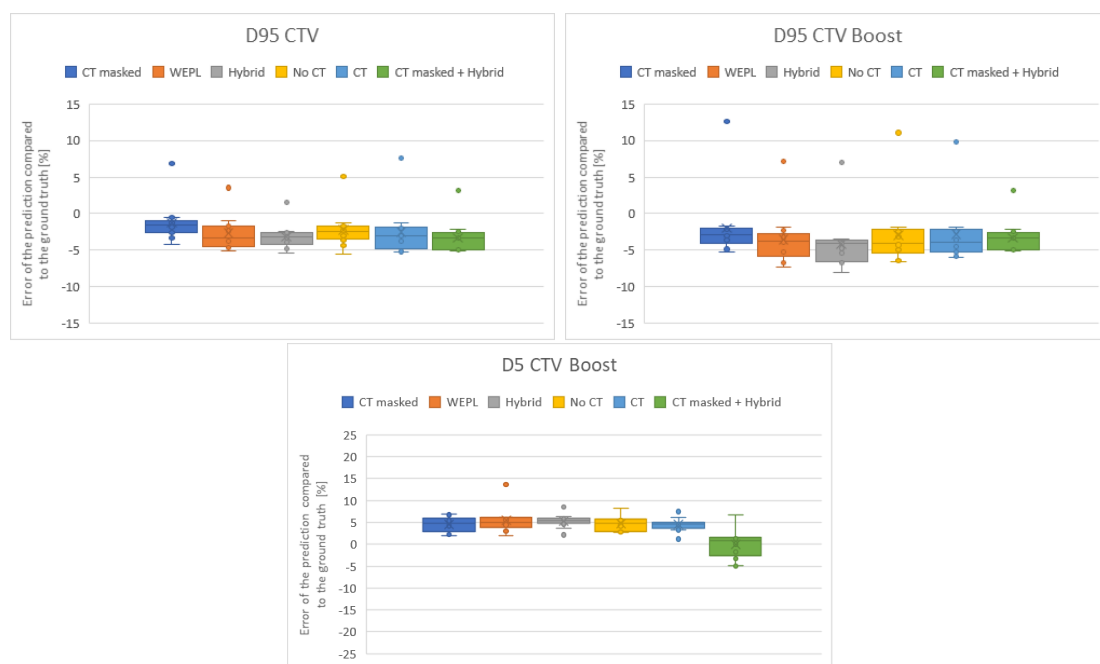


Figure 5.1: Comparison of different method performance for the CTV expressed in percent of the highest dose prescription (70Gy).

## 5.2.2 Organs at risk

Figures 5.2 and 5.3 present either the difference with the ground truth of D2 or of  $D_{\text{mean}}$  metric for the different organs at risk. D2 is interesting for serial organs such as the spinal cord, brainstem and mandible. For these organs, receiving a high dose at one point is critical. For others,  $D_{\text{mean}}$  is the metric used. Similarly to Figure 5.1, the six methods are represented on each graph.

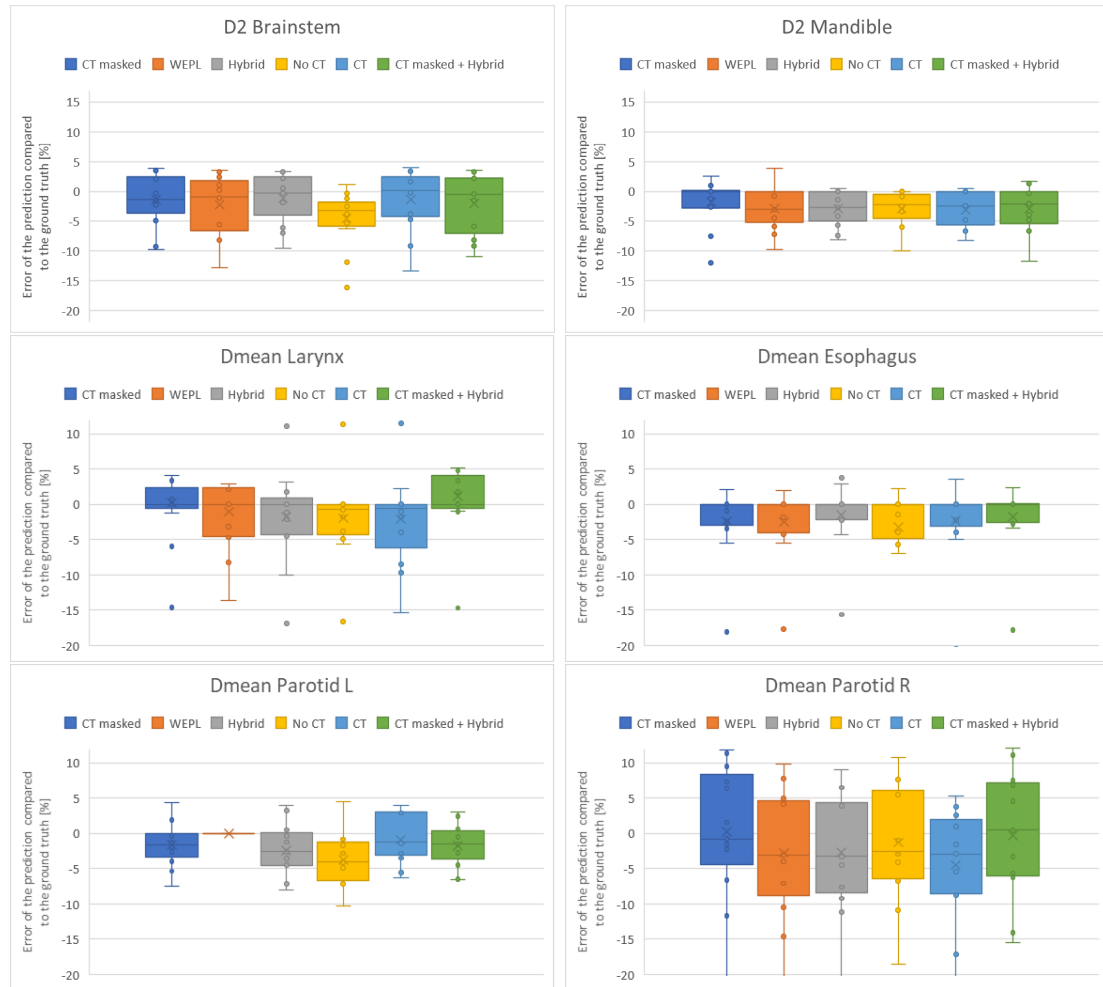


Figure 5.2: Top: Difference between D2 of predicted and ground truth doses for each methods over the brainstem and mandible. Center: difference between  $D_{\text{mean}}$  of predicted and ground truth for larynx and esophagus. Bottom: Difference between  $D_{\text{mean}}$  of predicted and ground truth for each methods over the left and right parotid, glands required to ensure salivation process.

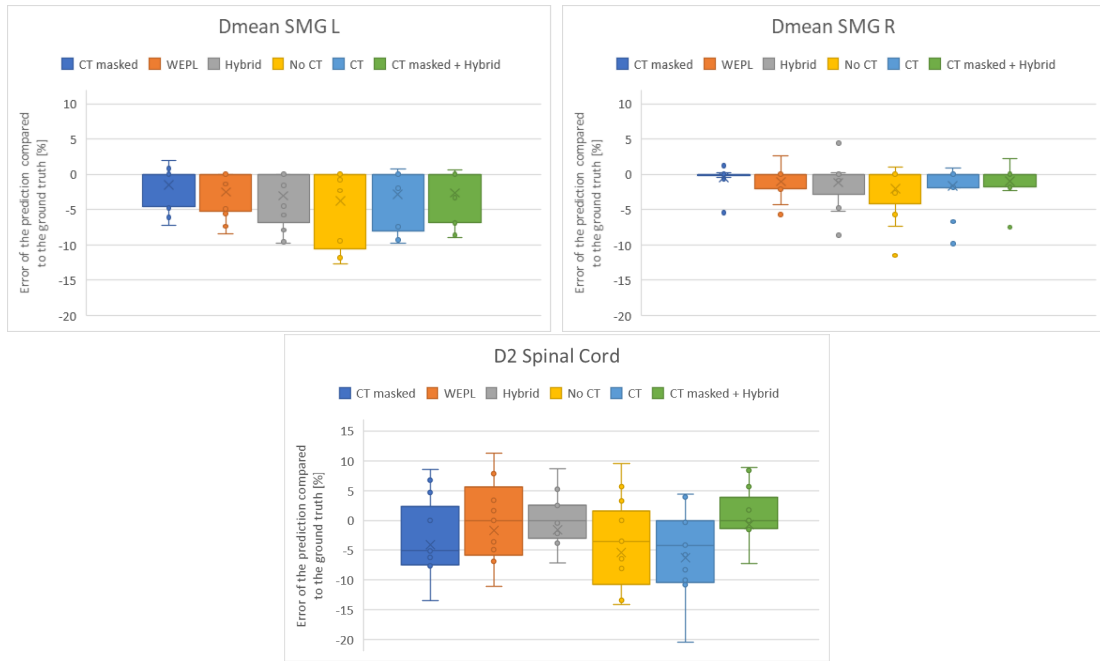


Figure 5.3: Top: Difference between  $D_{\text{mean}}$  of predicted and ground truth doses for each methods over the left and right submandibular glands (SMG). Bottom: Difference between D2 of predicted and ground truth doses for spinal cord.

### 5.3 Comparison dose prediction and ground truth

Figure 5.4 displays the predicted dose and ground truth for different patients. The left column represents the dose distribution of the ground truth. The center column represent the dose prediction outputted by the model and the right column displays the difference between the two. The CTV contours are visible as well.

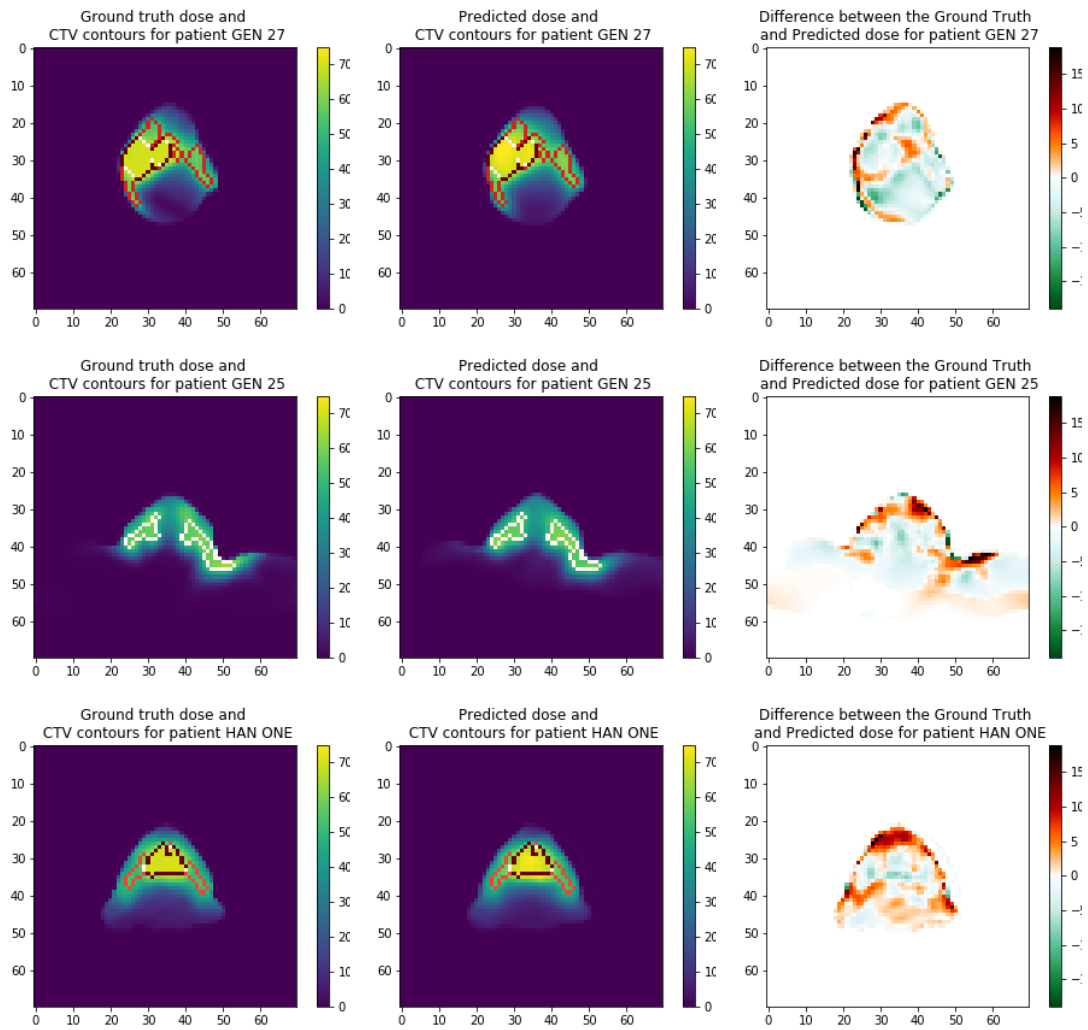


Figure 5.4: Dose maps to compare between the prediction and the ground truth exported from RayStation. The first column is the ground truth, the second is the dose prediction and the third represents the difference between the two. The contours represent the CTV.

Figure 5.5 represents two Dose-Volume histograms. The dashed line is the prediction while the solid line is the ground truth. In light and dark blue, the curves correspond to the two CTV dose prescription. All the other curves are the organs at risk that were contoured for the corresponding patient.

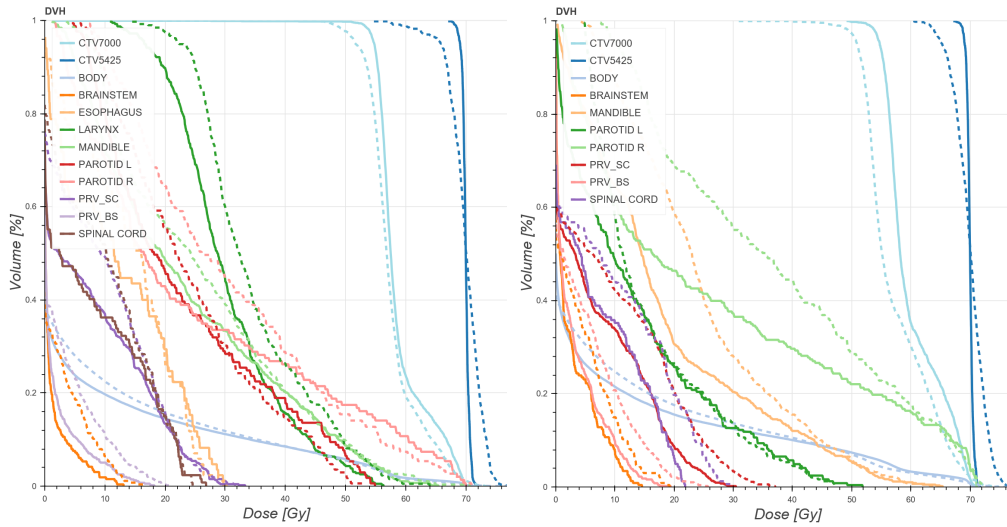


Figure 5.5: Example of DVH obtained for two of the best patients in the test set. The dashed line is the prediction while the solid line is the ground truth.

## 5.4 Computation time

Measuring the required computational time is essential because the purpose of this method is to be integrated in a day to day treatment workflow.

For methods not requiring the WEPL image generation, the computation time required to perform the dose prediction depends on the size of the CT image. For patients with CT size of  $512 \times 512 \times 218$  pixels, then resized to  $164 \times 164 \times 87$  pixels, the prediction time goes up to 6 seconds, whereas patients with CT size of  $256 \times 256 \times 146$  pixels, resized to  $100 \times 100 \times 41$  pixels, lies between 1 and 2 seconds.

For methods requiring the WEPL image generation, an additional 20 seconds are required to generate the WEPL image from the CT scan. The other disadvantage is that it requires opening MIROpt in Matlab to later process the image in Python for the prediction.

About the training, it took 1 hour to train the model with 80 data (the 40 patients that have been flipped as explained in Section 4.2) over 100 epochs.

## 5.5 Dose Mimicking

Dose mimicking has been performed from the dose prediction of a patient. Two plans have been generated: first a conventional plan and then a robust one. This section displays screenshots taken from the RayStation interface.

### 5.5.1 Conventional plan compared to prediction

Figure 5.6 displays screenshots from RayStation software. It shows the comparison between the dose prediction and the plan that RayStation optimized. The two CTV contours are displayed: in pink CTV 70Gy and in light blue CTV 54.25Gy. Spinal cord and parotids are displayed as well in light green and dark blue. Figure 5.7 shows the comparison on the DVH curves.

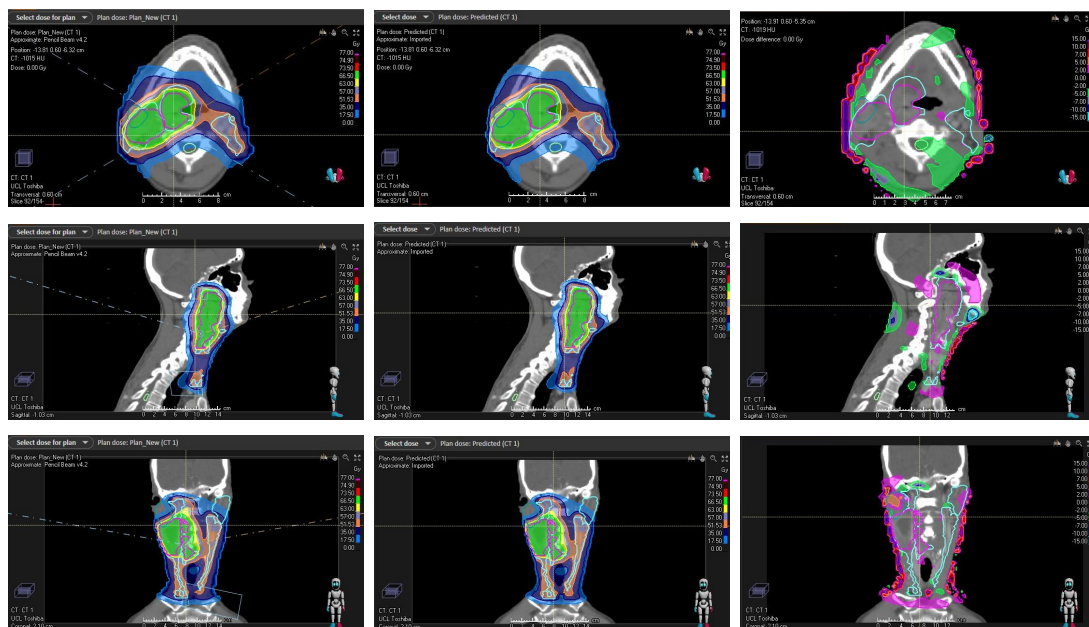


Figure 5.6: Comparison of the dose prediction from UNet and the conventional plan generated from dose mimicking. Left column shows the prediction, center column shows the new plan and right column displays the difference between the two.



Figure 5.7: DVH curves of the prediction in dotted line and of the conventional plan generated in solid line. The pink curve represents the CTV boost of 70Gy, the light blue is the CTV prescription of 54.25Gy and the others are the spinal cord and parotids.

### 5.5.2 Robust plan compared to prediction

A downside of the previous conventional plan generated is that it is not robust. The initial doses generated by the dosimetrist were robust, therefore there is a loss of quality in the plan. However, it is possible to perform the dose mimicking in a robust manner. Two isodoses are computed to be robust to variations: the isodose at 51.50Gy which encompasses the whole volume of CTV 54.25 and another isodise at 66.50Gy which encompasses the CTV 70. Figures 5.8 and 5.9 show the comparison between the prediction and the robust plan.

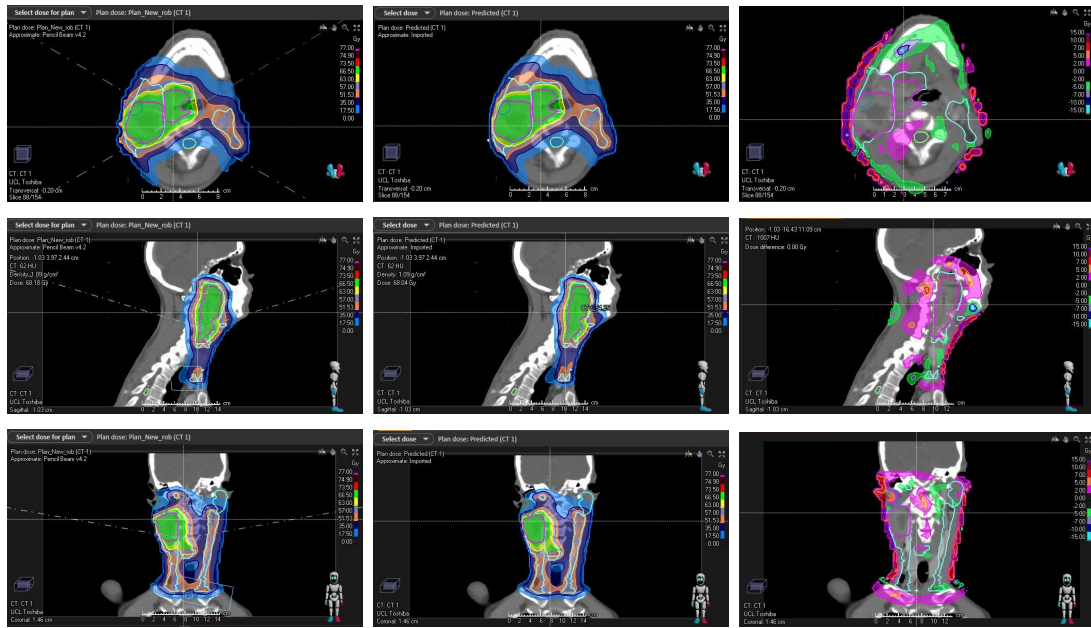


Figure 5.8: Comparison of the dose prediction from UNet and the robust plan generated from dose mimicking. Left column shows the prediction, center column shows the dose from the conventional plan generated in RayStation software and right column displays the difference between the two. All these images are captured from the RayStation interface.

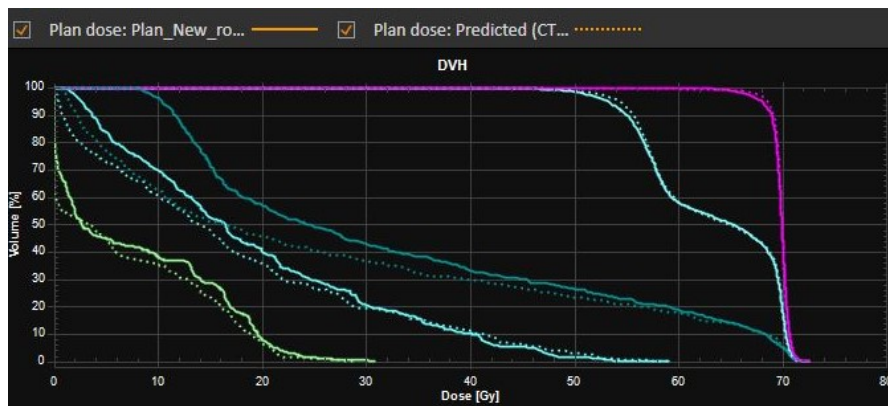


Figure 5.9: DVH curves of the prediction are the dotted lines and of the ones of the generated plan are the solid lines. The pink curve represents the CTV boost of 70Gy, the light blue is the CTV prescription of 54.25Gy and the others are the spinal cord and parotids.

### 5.5.3 Robust plan compared to conventional plan

Figure 5.10 presents a comparison of the dose distribution in the robust plan and conventional plan.

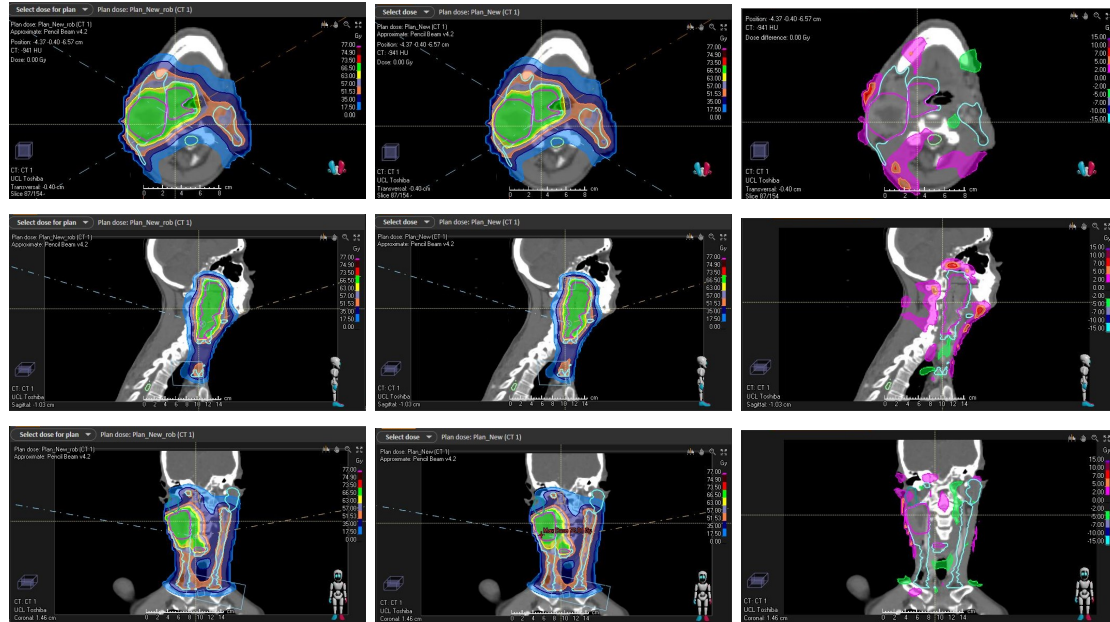


Figure 5.10: Comparison of the robust plan and conventional plan generated from dose mimicking. Left column shows the robust plan, center column shows the dose from the conventional plan and right column displays the difference between the two.

Some DVH metrics are displayed in Table 5.4 and 5.5.

	D95 CTV 54.25	D95 CTV 70
Ground truth	54.75 Gy	69.09 Gy
Prediction	54.34 Gy	68.58 Gy
Conventional plan	53.97 Gy	68.39 Gy
Robust plan	53.23 Gy	68.02 Gy

Table 5.4: Comparison of D95 on the two CTV contours, for the dose prediction, the dose from the conventional plan and from the robust plan.

$D_{\text{mean}}$	Spinal Cord	Parotid L	Parotid R
Ground truth	6.95 Gy	13.32 Gy	26.58 Gy
Prediction	6.90 Gy	17.16 Gy	26.61 Gy
Conventional plan	6.68 Gy	18.38 Gy	30.09 Gy
Robust plan	7.75 Gy	18.98 Gy	32.76 Gy

Table 5.5: Comparison of  $D_{\text{mean}}$  on three OARs: spinal cord and parotids, for the dose prediction, the dose from the conventional plan and from the robust plan and the ground truth.

### 5.5.4 Position of spots

Figure 5.11 shows the dose of the robust plan with the spot positions and the isodoses.

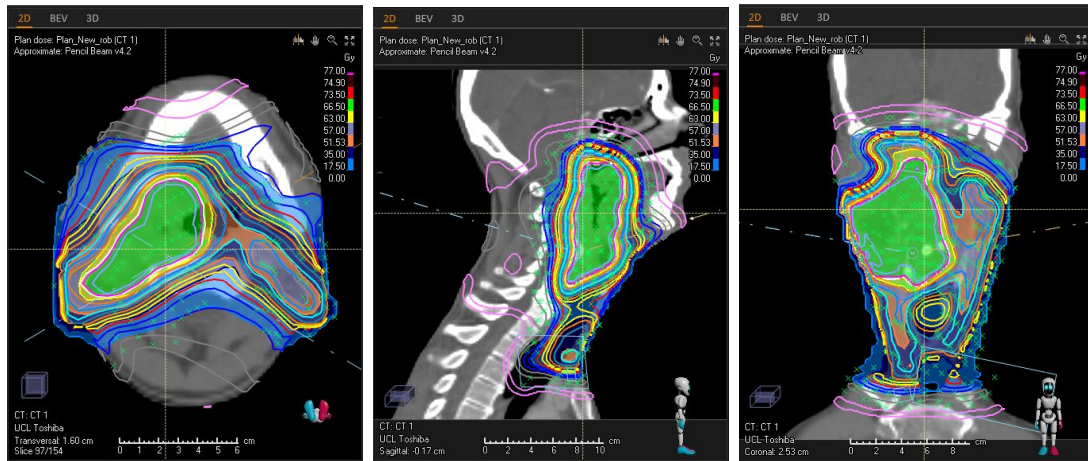


Figure 5.11: Spot positions can be seen by the small green crosses

## 6 | Discussion

### 6.1 Results interpretation

The first cross-validation conducted with the mean of MSE of each fold as generalization error ranked the *CT masked* method first (mean error of 2.13% of the highest prescription dose), followed by the Hybrid and WEPL methods (mean error of 5.43% and 5.67% respectively). The least performing one, according to MSE is the *CT masked + Hybrid* method with a mean error of 19.26%.

As mentioned in Section 5.1, the MSE is only a surrogate of the clinical metrics that physician are interested in when assessing the dose distribution. With the  $E_{\text{gen}^*}$  metric that has been defined in Section 5.1, the *CT masked* method is the best one (error value of 3.85%). However, they do not agree on the second best, which is with this metric, the *CT masked + hybrid* (error of 4.22%).

This difference can be explained by the fact that the first generalisation error is the voxel-wise MSE all over the grid that is considered. Whereas the  $E_{\text{gen}^*}$  metric only considers the difference in the region of interest (target volumes and organs at risk).

The MSE cross-validation, as well as the sum of DVH metrics cross-validation indicate that the *CT masked* method yields the best dose prediction. The second best is the Hybrid for the MSE metric and *CT masked + Hybrid* method for the  $E_{\text{gen}^*}$  metric.

#### 6.1.1 Results on CTV

As shown in the two top graphs of Figure 5.1, D95 of the CTVs present a mean error below 4% of the highest prescribed dose and an interquartile range of less than 4%. When comparing the *CT* and *CT masked* methods for the D95, it can still be observed that the interquartile range is reduced when introducing the mask shown in Section 4.6.5, which justifies the choice of adding it.

The WEPL and Hybrid methods present a mean error of -3.5% for the D95 of the CTV and -4% for the boost, which are higher than *CT masked* that presents a mean error of respectively -1.5% and -2.5%.

On the other hand, when considering solely the D5 of the CTV with a boost prescription of 70Gy, the addition of the two channels *CT masked and Hybrid* yields the smallest mean error.

In terms of difference in D95 on the two prescription volumes, the *CT masked* method presents the lowest mean and standard deviation. This observation confirms the cross validation.

The difference in D5 of the boost prescription volume is smaller for the *CT masked + Hybrid method*.

These two methods stand out from all the others.

### 6.1.2 Results on the OARs

Figure 5.2 and 5.3 present the difference in either D2 or  $D_{\text{mean}}$  metrics, depending on the OAR.

Let's focus mainly on the *CT masked* and the *CT masked + Hybrid* method as there are the most promising methods according to the results on the CTV.

Both present a mean error between 0 and 2.5% of the highest prescription dose for all organs, except *CT masked* that has a mean error of 4% for the spinal cord. *CT masked* presents a smaller mean difference for five contours: larynx, mandible, left and right submandibular glands and right parotid. A difference of more than 1% is found in only one organ: the spinal cord with an error of 4% for *CT masked* against 1% for *CT masked + Hybrid*. Overall, concerning the mean, perform similarly.

Interquantile ranges however are smaller or equal for all organs except the spinal cord with *CT masked* method. This observation indicates that the best performing one is the *CT masked*. To palliate the problem of the spinal cord, putting an emphasis on the channel containing spinal cord by increasing the weight could be investigated.

Overall, the best performing technique seems to be the *CT masked* method. It could benefit from tuning on the weight of the spinal cord channel to really outperform the *CT masked + Hybrid* on every front.

### 6.1.3 Dose Mimicking

Dose mimicking is the key step to assess the feasibility of the automatic planning method based on a neural network presented in this thesis. Indeed, as mentioned in the introduction, the dose prediction in itself is not enough. The clinician needs the machine parameters that lead to the delivery of the dose. Dose mimicking is the optimization of parameters by the treatment planning system to fit the isodoses of the dose prediction.

The resulting dose from dose mimicking in Figure 5.6, two observations stand out. First, the differential image is really needed to assess the differences as the two first columns of Figure 5.6 are very similar. There is no difference between the achievable dose and the prediction in the CTV. The dose in the rest of the body presents a small difference of 2Gy in the pink region such as by the jaw in the sagittal view. A reduction of the dose by 5Gy can be observed in green on the axial view. Similar features can be observed in Figure 5.7 where the contours with high dose are really conform but the low doses show some slight differences. The dose distribution in the body is therefore very similar in the contours specified.

Second observation is that the largest difference is located around the body with red-purple dots, representing a difference of 5 to 15 Gy. This can be explained by the position of the pencil beam spots. Indeed, in the conventional workflow, the treatment planning system knows that it has to put all the spots in the CTV. However, as now the objective is to reproduce the isodoses, there are spots disseminated in a larger region. This phenomena is shown in Figure 5.11. The same problem justifies the difference of the low doses of Figure 5.7. A solution would be to find a way to specify where the spots should be located.

In Sections 5.5.2, it was possible to obtain a robust plan from RayStation. Figure 5.8 shows similar characteristics as the non robust plan: very similar dose distribution in the CTVs, some variation ranging from -5Gy and 2Gy in the body and the errors due to the misplaced spots on the outside. In Tables 5.4 and 5.5, D95 of the two CTVs from the robust dose are slightly less (53.23Gy and 68.02Gy) than the non robust one (53.97Gy and 68.39Gy). However, being able to have the same robustness characteristic as the initial dose generated by the dosimetrist is a big advantage.

In Figure 5.10, slices were selected in order to see where lies the difference: in the back of the neck area.

The dose mimicking confirms the motivation of this thesis. Indeed, the treatment planning system managed to retrieve a dose close to the prediction without any manual tuning.

In addition, it was possible to generate robust dose from the prediction. The robust dose is slightly less good but it allows to keep the robustness in various scenarios. One lacking feature that would improve further the reproduction in the low dose areas is to be able to specify that pencil beam spots must lie in the CTV area.

#### 6.1.4 Computation time

There is a clear difference between the computational time required for methods requiring the WEPL image : *WEPL*, *Hybrid* and *CT masked + Hybrid* methods; and the others: *No CT*, *CT* and *CT masked*.

The time require ranges from 21 to 26 seconds for methods with the WEPL and 1 to 6 seconds for the others. The reason why computing the WEPL is so slow is that the Matlab functions used are not very time efficient. This time would certainly be reduced if dedicated functions are created to compute WEPL in a fast way. However, for now, this aspect is a strong downside to the *WEPL* methods, and in particular to the *CT masked + Hybrid* that performed well in the dose prediction.

In Section 5.4, a difference of prediction time has been observed for patients with different CT image sizes. The difference for CT scans of size (100, 100, 41) and (164, 164, 87) is due to the higher number of patches. For the first one, 8 patches were computed whereas for the latter 27 patches where considered.

The longest step is the merging of the different patches which requires to loop on the whole grid. Indeed, individual patch prediction requires 0.15 seconds to be predicted. This aspect should definitely be optimized in future improvements.

The best methods in term of time efficiency are the ones not requiring the WEPL computation. This is another advantage of the *CT masked* over the *CT masked + Hybrid* method. Improvement can be done over the merging of the different patches to keep a time required of less than to second even for larger CT images.

## 6.2 Limitations

The results enumerated in the previous section are very promising. However, this model has still some limitations and aspects that can be improved.

### 6.2.1 Generalization

The data used in this study was limited to a small number of patients, only 50 for the training. In articles presented in Section 3, training sets were either larger: 195 in [23], 106 in [20], 80 in [19] and in [22] and they could use extensive data augmentation by deforming the anatomy and the dose accordingly. Furthermore, the patients were selected so that only bilateral cases and they were planned using the same beam configuration. This helped the network to perform dose prediction on such a small dataset. It would be interesting to see whether the network is able to perform that well when a larger dataset where unilateral cancer cases are added.

### 6.2.2 Limitation on the mask

The restriction of the information about the CT is on the geometric projection of the shape of the CTV in the beam eye view. It has been observe that applying this mask improves the prediction of the network. However, there are two limitations. First, it ignores the scattering of the protons. Indeed, protons undergo scattering when they are involved in elastic collision. Second, when using robust optimization, the optimizer will put some dose around the target in order to cover all the uncertainties. Adding a margin on the mask to give information about tissue density in these additional areas could be considered to have a more accurate prediction of the dose distribution.

### 6.2.3 Accuracy

To allow the network to learn efficiently and not being confused by the different resolution of images, coming from different centers and CT scanners, all images where rescaled to voxels of  $5*5*5$  mm. The grid is therefore relatively coarse. It would be interesting to see if reducing the size of the voxels to  $2.5*2.5*2.5$ mm is still feasible considering the increased computational time and resources required for the prediction.

## 6.3 Perspectives

In addition to the improvements that can be brought to the model presented in this thesis (such as the spots position for the dose mimicking or optimizing computation time), this section presents possible perspectives that could be investigated.

### 6.3.1 Uncertainty estimation

An important aspect that is lacking in our model is the estimation of the uncertainty on the output. Whether the inputs are correct or not, whether the type of cancer is what it had been trained on or not, the network will output a dose prediction. This is an issue for the day-to-day delivery. A check of the respect of clinical goals is a good indicator whether the results are wrong to allow the tracking down of the issue.

However, a very interesting approach would be to retrieve, from the network itself, an estimation of the uncertainty on the accuracy of the output. Theory about Bayesian neural network have been developed [32]. The Bayesian neural network decomposes uncertainty into model uncertainty, model misspecification, and inherent noise. The major difference is that the parameters are not fixed weights but distributions. It would require rethinking the architecture, but it seems to be a really promising approach. In this article from 2018 [33], the authors managed to estimate the uncertainty on the segmentation output by their network.

### 6.3.2 GAN

This thesis has demonstrated that UNet is able to efficiently predict dose distribution for proton has long that the *CT masked* is provided in input. It would be interesting to investigate if using a GAN with our network as generator would improve the prediction. The article [25] did a comparison with a UNet prediction but in 2D only.

## 7 | Conclusion

The goal of this thesis was to propose a neural network model to perform 3D dose prediction of a proton therapy treatment for head and neck. Such a model is necessary to implement an automatic proton therapy treatment workflow and to allow the generation of a treatment plan. The model developed in this study fulfilled all the requirements established in the introduction: it is accurate (less than 4% of mean error between the clinical DVH metrics, expressed in percent of the highest prescription dose), fast and achievable through dose mimicking.

From the different experiments conducted, UNet has shown to learn more useful information when the input includes the *CT masked* information. The input is therefore composed of :

- The 3D CT image that has been masked to leave only the values on the beam trajectory, the rest is put to zero.
- A channel with a 3D matrix of the prescription on Clinical Target Volume.
- 13 channels with binary masks (1 inside the volume and 0 outside) of organs at risk.

This model outputs a prediction of the 3D dose distribution. When comparing with the ground truth over a testset of patients, the mean difference between D95 in CTV is less than 2.5% of the highest prescribed dose. For all organs at risk but spinal cord, the mean difference of  $D_{\text{mean}}$  or D2 of less than 1% of the highest prescribed dose. For spinal cord, a mean error of 4% has been found.

Additionally, the model takes only 1 to 6 seconds to predict the full dose distribution. This confirms that this method can be used in an automatic workflow without slowing down the process.

More than just the prediction of dose distribution, dose mimicking has been implemented to obtain a robust treatment plan. The dose of this plan is close to the prediction: a difference of 1Gy between D95 of CTV 54.25Gy and a difference

of 0.5Gy between D95 of CTV 70Gy. The difference on the organs at risk is larger 1.8Gy for left parotid and 4Gy for right parotid. Dose mimicking can still be improved for the low dose areas by finding a way to specify the that pencil beam spots must lie in the CTV area.

The two following perspectives would be interesting to investigate:

- The field of Bayesian networks could lead to the output dose distribution, accompanied by an estimation of the uncertainty over the accuracy of the dose prediction. Such an estimation would be very informative as the stakes are high in the radio-oncology field.
- Another approach is to use the generative adversarial network framework of learning. An article cited in the literature review published results that outperformed UNet in conventional radiotherapy and predicting 2D slice by 2D slice. It would be interesting to see if the results are improved also in the case of 3D dose prediction in proton therapy.

# Bibliography

- [1] AVO | The Potential of Proton Therapy.
- [2] Amichetti Maurizio Nancy J Tarbell Barbara Rombi, Shannon M MacDonald and Torunn I Yock. Proton radiotherapy for childhood tumors: an overview of early clinical results. *Journal of Nuclear Medicine Radiation Therapy*, 4(4):1–10, 2013.
- [3] Harald Paganetti, Andrzej Niemierko, Marek Ancukiewicz, Leo E Gerweck, Michael Goitein, Jay S Loeffler, and Herman D Suit. Relative biological effectiveness (rbe) values for proton beam therapy. *International Journal of Radiation Oncology\*Biology\*Physics*, 53(2):407 – 421, 2002.
- [4] G. Janssens E. Sterpin, J. Lee. Engineering challenges in protontherapy, gbio2070 course. 2017.
- [5] J. Lee A. Bol. Medical imaging, gbio2050 course. 2017.
- [6] Christopher Johnson M.D. PICU Author. Ct scan of children with minor head trauma. 10 2018.
- [7] DIACOR. Centralite patient alignment lasers. 3 2016.
- [8] Stanislav Nikolov, Sam Blackwell, Ruheena Mendes, Jeffrey De Fauw, Clemens Meyer, Cian Hughes, Harry Askham, Bernardino Romera-Paredes, Alan Karthikesalingam, Carlton Chu, Dawn Carnell, Cheng Boon, Derek D’Souza, Syed Ali Moinuddin, Kevin Sullivan, DeepMind Radiographer Consortium, Hugh Montgomery, Geraint Rees, Ricky Sharma, and Olaf Ronneberger. Deep learning to achieve clinically applicable segmentation of head and neck anatomy for radiotherapy, 09 2018.
- [9] Peter Park, Xiaorong Zhu, Andrew K Lee, Narayan Sahoo, Adam Melancon, Lifei Zhang, and Lei Dong. A beam-specific planning target volume (ptv) design for proton therapy to account for setup and range uncertainties. *International journal of radiation oncology, biology, physics*, 82:e329–36, 06 2011.

- [10] Ilma Khaferllari, Eugene Wong, Karl Bzdusek, Michael Lock, and Jeff Z. Chen. Automated imrt planning with regional optimization using planning scripts. *Journal of applied clinical medical physics / American College of Medical Physics*, 14:4052, 01 2013.
- [11] Michael B Sharpe and Kevin Moore. Point/counterpoint: Within the next ten years treatment planning will become fully automated without the need for human intervention. *Medical physics*, 41:120601, 12 2014.
- [12] Yidong Yang, Eric C Ford, Binbin Wu, Michael Pinkawa, Baukelien Triest, Patrick Campbell, Danny Y Song, and Todd McNutt. An overlap-volume-histogram based method for rectal dose prediction and automated treatment planning in the external beam prostate radiotherapy following hydrogel injection. *Medical physics*, 40:011709, 01 2013.
- [13] Kelly C. Younge, Robin B. Marsh, Dawn Owen, Huaizhi Geng, Ying Xiao, Daniel Spratt, Joseph Foy, Krithika Suresh, Q Jackie Wu, Fang-Fang Yin, Samuel Ryu, and Martha M. Matuszak. Improving quality and consistency in nrg oncology rtog 0631 for spine radiosurgery via knowledge-based planning. *International Journal of Radiation Oncology\*Biology\*Physics*, 100, 01 2018.
- [14] Satomi Shiraishi and Kevin Moore. Knowledge-based prediction of three-dimensional dose distributions for external beam radiotherapy. *Medical Physics*, 43:378–387, 01 2016.
- [15] Chuan Zeng, Kevin Sine, and Dennis Mah. Contour-based lung dose prediction for breast proton therapy. *Journal of Applied Clinical Medical Physics*, 19(6):53–59, 2018.
- [16] David C. Hall, Alexei V. Trofimov, Brian A. Winey, Norbert J. Liebsch, and Harald Paganetti. Predicting patient-specific dosimetric benefits of proton therapy for skull-base tumors using a geometric knowledge-based method. *International Journal of Radiation Oncology\*Biology\*Physics*, 97(5):1087 – 1094, 2017.
- [17] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. pages 234–241, 2015.
- [18] Dan Nguyen, Troy Long, Xun Jia, Weiguo Lu, Xuejun Gu, Zohaib Iqbal, and Shucui Jiang. Dose prediction with u-net: A feasibility study for predicting dose distributions from contours using deep learning on prostate imrt patients. 09 2017.

- [19] Dan Nguyen, Xun Jia, David Sher, Mu-Han Lin, Zohaib Iqbal, Hui Liu, and Shucui Jiang. Three-dimensional radiotherapy dose prediction on head and neck cancer patients with a hierarchically densely connected u-net deep learning architecture. *Physics in Medicine and Biology*, 64, 01 2019.
- [20] Vasant Kearney, Jason W Chan, Samuel Haaf, Martina Descovich, and Timothy D Solberg. DoseNet: a volumetric dose prediction algorithm using 3d fully-convolutional neural networks. *Physics in Medicine & Biology*, 63(23):235022, dec 2018.
- [21] Ana M. Barragan-Montero, Dan Nguyen, Weiguo Lu, Mu-Han Lin, Xavie Geets, Edmond Sterpin, and Steve Jiang. Three-dimensional dose prediction for lung imrt patients with deep neural networks: Robust learning from heterogeneous beam configurations. 12 2018.
- [22] Xinyuan Chen, Kuo Men, Yexiong Li, Yi Lu, and Jianrong Dai. A feasibility study on an automated method to generate patient-specific dose distributions for radiotherapy using deep learning. *Medical Physics*, 46, 10 2018.
- [23] Jiawei Fan, Jiazhou Wang, Zhi Chen, Chaosu Hu, Zhen Zhang, and Weigang Hu. Automatic treatment planning based on three-dimensional dose distribution predicted from deep learning technique. *Medical Physics*, 46(1):370–381, January 2019.
- [24] Zhiqiang Liu, Jiawei Fan, Minghui Li, Hui Yan, Zhihui Hu, Peng Huang, Yuan Tian, Junjie Miao, and Jianrong Dai. A deep-learning method for prediction of three-dimensional dose distribution of helical tomotherapy. *Medical Physics*, 46, 03 2019.
- [25] Rafid Mahmood, Aaron Babier, Andrea Mcniven, Adam Diamant, and Timothy C. Y. Chan. Automated treatment planning in radiation therapy using generative adversarial networks. 07 2018.
- [26] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [27] Dingjun Yu, Hanli Wang, Peiqiu Chen, and Zhihua Wei. Mixed pooling for convolutional neural networks. pages 364–375, 10 2014.
- [28] Jelo Salomon and Bianca Schoen Phelan. Lung cancer detection using deep learning. 04 2018.
- [29] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks

from overfitting. *Journal of Machine Learning Research*, 15:1929–1958, 06 2014.

- [30] Sebastian Ruder. An overview of gradient descent optimization algorithms. *CoRR*, abs/1609.04747, 2016.
- [31] Digital imaging and communications in medicine (dicom). *NEMA Publications PS 3.1-PS 3.12. The National Electrical Manufacturers Association*. Rosslyn, VA, 41, 1992,1993,1994,1995.
- [32] Radford M. Neal. *Bayesian Learning for Neural Networks*. Lecture Notes in Statistics. Springer-Verlag, New York, 1996.
- [33] Terrance DeVries and Graham W. Taylor. Leveraging uncertainty estimates for predicting segmentation quality. 07 2018.

**UNIVERSITÉ CATHOLIQUE DE LOUVAIN**  
École polytechnique de Louvain

Rue Archimède, 1 bte L6.11.01, 1348 Louvain-la-Neuve, Belgique | [www.uclouvain.be/epl](http://www.uclouvain.be/epl)