

ACCEPTED MANUSCRIPT

# Development of robustness evaluation strategies for enabling statistically consistent reporting

To cite this article before publication: Edmond Sterpin *et al* 2020 *Phys. Med. Biol.* in press <https://doi.org/10.1088/1361-6560/abd22f>

## Manuscript version: Accepted Manuscript

Accepted Manuscript is “the version of the article accepted for publication including all changes made as a result of the peer review process, and which may also include the addition to the article by IOP Publishing of a header, an article ID, a cover sheet and/or an ‘Accepted Manuscript’ watermark, but excluding any other editing, typesetting or other changes made by IOP Publishing and/or its licensors”

This Accepted Manuscript is © 2020 Institute of Physics and Engineering in Medicine.

During the embargo period (the 12 month period from the publication of the Version of Record of this article), the Accepted Manuscript is fully protected by copyright and cannot be reused or reposted elsewhere.

As the Version of Record of this article is going to be / has been published on a subscription basis, this Accepted Manuscript is available for reuse under a CC BY-NC-ND 3.0 licence after the 12 month embargo period.

After the embargo period, everyone is permitted to use copy and redistribute this article for non-commercial purposes only, provided that they adhere to all the terms of the licence <https://creativecommons.org/licenses/by-nc-nd/3.0>

Although reasonable endeavours have been taken to obtain all necessary permissions from third parties to include their copyrighted content within this article, their full citation and copyright line may not be present in this Accepted Manuscript version. Before using any content from this article, please refer to the Version of Record on IOPscience once published for full citation and copyright details, as permissions will likely be required. All third party content is fully copyright protected, unless specifically stated otherwise in the figure caption in the Version of Record.

View the [article online](#) for updates and enhancements.

# Development of robustness evaluation strategies for enabling statistically consistent reporting

E. Sterpin<sup>1,2</sup>, Sara T Rivas<sup>2</sup>, F. Van den Heuvel<sup>3,4</sup>, B. George<sup>3</sup>, J.A. Lee<sup>2</sup>, and K. Souris<sup>2</sup>

<sup>1</sup> KU Leuven, Department of Oncology, Laboratory of Experimental Radiotherapy, Leuven, Belgium.

<sup>2</sup> Université catholique de Louvain, Institut de Recherche Expérimentale et Clinique, Center of Molecular Imaging, Radiotherapy and Oncology (MIRO), Brussels, Belgium

<sup>3</sup> CRUK/MRC Oxford Institute for Radiation Oncology, University of Oxford, Oxford, United Kingdom

<sup>4</sup> Dept of Haematology/Oncology, Oxford University Hospitals NHS Foundation Trust, Oxford, United Kingdom

## Abstract

Robustness evaluation of proton therapy treatment plans is essential for ensuring safe treatment delivery. However, available evaluation procedures feature a limited exploration of the actual robustness of the plan and generally do not provide confidence levels. This study compared established and more sophisticated robustness evaluation procedures, with quantified confidence levels.

We have evaluated several robustness evaluation methods for 5 bilateral head-and-neck patients optimized considering spot scanning delivery and with a conventional CTV-to-PTV margin of 4 mm. Method 1) good practice scenario selection (GPSS) (e.g +/- 4 mm setup error 3% range uncertainty); 2) statistically sound scenario selection (SSSS) either only on or both on and inside isoprobability hypersurface encompassing 90% of the possible errors; 3) statistically sound dosimetric selection (SSDS). In the last method, the 90% best plans were selected according to either target coverage quantified by  $D_{95}$  (SSDS\_ $D_{95}$ ) or to an approximation of the final objective function (OF) used during treatment optimization (SSDS\_OF). For all methods, we have considered systematic setup and systematic range errors. A mix of systematic and random setup errors were also simulated for SSDS, but keeping the same conventional margin of 4 mm.

All robustness evaluations have been performed using the fast Monte Carlo dose engine MCsquare. Both SSSS strategies yielded on average very similar results. SSSS and GPSS yield comparable values for target coverage (within 0.4 Gy). The most noticeable differences were found for

1  
2  
3 the CTV between GPSS, on the one hand, and SSDS\_ $D_{95}$  and SSDS\_OF, on the other hand (average  
4 worst-case  $D_{98}$  were 2.5 and 1.6 Gy larger than for GPSS, respectively). Simulating explicitly random  
5 errors in SSDS improved almost all DVH metrics.  
6  
7

8 We have observed that the width of DVH-bands and the confidence levels depend on the  
9 method chosen to sample the scenarios. Statistically sound estimation of the robustness of the plan in  
10 the dosimetric space may provide an improved insight on the actual robustness of the plan for a given  
11 confidence level.  
12  
13  
14

## 15 16 17 18 **I. Introduction** 19

20  
21  
22 External beam radiotherapy aims at delivering sufficient dose to tumor tissue while preserving  
23 surrounding organs. In order to achieve this goal, sophisticated irradiation techniques, such as the use  
24 of protons, have been developed in order to conform doses to target volumes. However, radiotherapy  
25 treatment delivery can be affected by many sources of uncertainties: patient positioning, inter- and  
26 intra-fraction movements, and imperfect conversion of imaging data into physical quantities. In order  
27 to secure target coverage and avoid accidental organ-at-risk irradiation, robust planning methods have  
28 been developed to ensure that delivered doses keep meeting the objectives and constraints despite  
29 uncertainties. In conventional X-ray radiotherapy, this objective is typically achieved by safety margins  
30 with the concept of planning target volume (PTV) and planning risk volume (PRV). In proton therapy,  
31 the typical margin strategies suffer from notorious shortcomings, because of the sensitivity of proton  
32 therapy dose distributions to the uncertainties of the position of the Bragg peak and failure of the  
33 static-dose cloud approximation, which assumes that patient shifts do not change the dose  
34 distributions <sup>1,2</sup>.  
35  
36  
37  
38  
39  
40  
41  
42

43 As a result, many methods of robust planning and robustness evaluation have been proposed  
44 in the literature for proton therapy <sup>1-8</sup>. In the case of the most advanced intensity modulated proton  
45 therapy (IMPT) techniques, robust planning typically consists of a minimax problem that is solved by  
46 optimizing the worst-case scenario among a set of predefined possible scenarios <sup>2</sup>. Generally, this set  
47 of scenarios includes (systematic) positioning errors, image conversion errors, and in some cases the  
48 movement of organs represented by additional image sets which are included as additional scenarios  
49 <sup>9-12</sup>. Irrespective of the considered robust planning method, it remains necessary to evaluate the actual  
50 robustness of the plan. The evaluation can be done more thoroughly than during robust optimization.  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60 Robustness evaluation often includes, for example, random errors, organ deformations, interplay  
effects, etc. This is because robustness evaluation is computationally less demanding than robust  
optimization (for instance, no need to store influence matrices). However, the methods typically

1  
2  
3 reported for robustness evaluation are also based on a relatively simple sampling of error scenarios,  
4 for example some (systematic) positioning errors combined with image conversion errors. Other  
5 authors have also incorporated random errors<sup>2,3</sup>.  
6  
7

8 Most robustness evaluation methods reported in the literature and used in some commercial  
9 planning systems may feature several biases because of pragmatic choices imposed by limited  
10 computing resources and due to a lack of consensus in the involved concepts. A first bias lies in the  
11 direct combination of pre-sampled uncertainties, leading to the selection of very unlikely scenarios,  
12 for example setup errors of +/- 5 mm combined with density errors of +/- 3%, i.e., the simultaneous  
13 selection of two extremes in the probability distributions. This amounts to combining extremes of  
14 marginal probability distributions, while the joint probability distribution should be sampled instead.  
15 Korevaar et al have already pointed that issue and have performed robustness evaluation using a  
16 statistically consistent but limited set of scenarios<sup>13</sup>. A second bias is the lack of consistently calculated  
17 confidence levels, in order to clearly define what is meant by a "worst-case". Indeed, the worst-case  
18 scenario is the least favorable scenario among a pre-defined selection set (otherwise, the most  
19 extreme case can always be envisaged). In the best-known margin calculation recipe, the value of the  
20 final margin depends on a choice of the number of patients for which one wishes to ensure target  
21 coverage (typically 90%)<sup>14</sup>. This confidence level is not always reported in the literature when it comes  
22 to robustness assessments. In addition, a lack of clarity remains on how to calculate this confidence  
23 level. Specifically, should it be calculated in the error space, i.e., as the percentage of possible scenarios  
24 covered by a given robustness test? Or should it be calculated in the dose space, that is, as the  
25 percentage of dose distributions meeting a given clinical endpoint?  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37

38 In this publication, we compare several robustness evaluation methods, with explicitly  
39 calculated confidence levels in either the error space or the dose distribution space.  
40  
41  
42  
43  
44

## 45 **II. Materials and methods**

### 46 **II.A. Definitions and notations**

47  
48  
49  
50 We define the robustness of a treatment plan as the capability of this plan to continue  
51 satisfying clinical objectives and/or constraints despite uncertainties, for a certain confidence level. As  
52 in van Herk et al<sup>14</sup>, treatment errors can be classified in treatment preparation errors (e.g., systematic  
53 errors) and treatment execution errors (e.g., random errors). Like van Herk's formalisation, we  
54 suppose knowledge of the probability density functions (*pdf*) of these errors. In general, these are  
55  
56  
57  
58  
59  
60

assumed to be normal (Gaussian) with standard deviations  $\Sigma$  and  $\sigma$  for systematic and random errors, respectively. For the remaining of this manuscript, we will limit ourselves to the following errors:

1. Setup errors ( $se$ ;  $(x_{se}, y_{se}, z_{se})$ ) characterized by 3D Gaussian *pdfs* for both systematic and random errors, with vector standard deviations  $\Sigma_{se}$  and  $\sigma_{se}$ , respectively
2. Range uncertainties (RU; due for instance to improper image conversion), characterized by a 1D Gaussian *pdf* with  $\Sigma_{RU}$  as standard deviation.

However, these considerations can be generalized to an arbitrary number of types of uncertainties. When appropriate, the generalization of the developed methods will be addressed.

## II.B. Computation of confidence levels

For a given *pdf*, confidence intervals define a range within which a population parameter resides for a given confidence level. In van Herk's margin recipe, a typical confidence level chosen is 90% which leads to the 2.5 factor that multiplies the standard deviation in the well-known formula for 3D-conformal dose distributions:  $M_{PTV} = 2.5\Sigma + 0.7\sigma$ . This means that 90% of the possible systematic errors within the patient population will be covered by the margin recipe. However, such a margin recipe fails notoriously in proton therapy because it is based on the static-dose cloud approximation<sup>15-17</sup>. Moreover, the number of fractions is assumed infinite, which allows a simple model to approximate how random errors blur the dose distributions. This simplification leads to the term  $0.7\sigma$  in the margin recipe, considering a typical 95% of the dose prescription as the minimum dosimetric coverage. A reduced number of fractions requires either a more complex model or the conversion of part of the random error into a systematic component, as acknowledged in van Herk et al<sup>14</sup>. Such approach will also not hold in proton therapy because of the failure of the static dose cloud approximation.

Thus, more complex models and formalisms are needed in proton therapy to assess the robustness in lieu of simplistic margin recipes. First, we need to distinguish occurrences of errors and the combined effect of these occurrences (the sum over each fraction) over the entire course of a treatment, referred here as *treatment scenarios* or, shorter, *scenarios*. In the simplified context mentioned here, a scenario will therefore be characterized by a systematic error sampled from a Gaussian distribution with a vector of standard deviations  $\Sigma_{se}$ , and a sequence of errors for each treatment fraction that are randomly sampled from a Gaussian distribution with a vector of standard deviations  $\sigma_{se}$ . Second, it is very unlikely to provide closed-form analytical expressions to characterize

1  
2  
3 how uncertainties affect the dose distributions. Therefore, it is not feasible to derive simple margin  
4 recipes with satisfactory mathematical grounds.  
5

6 In general, a confidence interval for an estimator of interest consists in giving the narrowest  
7 range of values for that estimator, such that the *pdf* integrates to 0.9 (90%) over that range. In practice,  
8 the *pdf* can be sampled and sorted, after which the suitable bounds can be reported.  
9

### 10 11 12 13 **II.B.1. Confidence levels in the dosimetric space** 14

15  
16 A straightforward way to compute a confidence level for a dosimetric estimator is to generate  
17 dose distributions for many scenarios and compute the probability that a certain rule on this dosimetric  
18 estimator will be realized (for instance,  $D_{95} > xx$  Gy with a probability of  $yy$  (or confidence)). This will  
19 be referred as the computation of a confidence level in the *dosimetric space*. In such approach, we can  
20 provide the percentage of times, i.e., the confidence level, that each objective/constraint defined by  
21 the radiation oncologist will be satisfied. Another possibility would be to provide a bandwidth for a  
22 value of interest and an associated confidence level. For instance, we could provide the range of  $D_{95}$   
23 for the CTV, corresponding to the 90% highest  $D_{95}$  values. This is a relevant metric to estimate the  
24 probability of covering the target as desired. However, this might cause to focus too much on target  
25 coverage. In order to provide a fair balance between target coverage and organs-at-risk exposure,  
26 another possibility would be to select the best 90% *objective function* (OF) values. The value of the  
27 objective function of the accepted plan, with the penalties (/ objective function weights) for each organ  
28 included in the objective function, provides a good estimate of the clinical compromise accepted by  
29 the physicist and the physician at the end of the optimization process. Thus, it provides a quantification  
30 of the clinical quality of the plan. Therefore, the classification of the best simulated dose distributions  
31 according to the value of their associated objective functions seems ideal from a clinical point of view.  
32

33 Because the confidence levels are estimated from random sampling of the errors, they will be  
34 subject to statistical noise. Therefore, enough scenarios must be simulated for estimating confidence  
35 levels with sufficient accuracy. The number of scenarios needed to achieve a given statistical accuracy  
36 on the confidence level can be determined using the method developed in Souris et al, where the  
37 statistical uncertainty on the estimated confidence level considered is computed dynamically during  
38 the robustness evaluation process<sup>18</sup>.  
39

40 Working in the dosimetric space requires the generation of many dose distributions that must  
41 be computed in a practical fashion. The key difficulty resides in the generation of the dose distributions  
42 that must be done in a practical fashion. Fast Monte Carlo dose engines associated with clever  
43 statistical stopping criteria<sup>18</sup> or other methods like polynomial chaos expansion<sup>19</sup> can help for this task.  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## II.B.2. Confidence levels in the error space

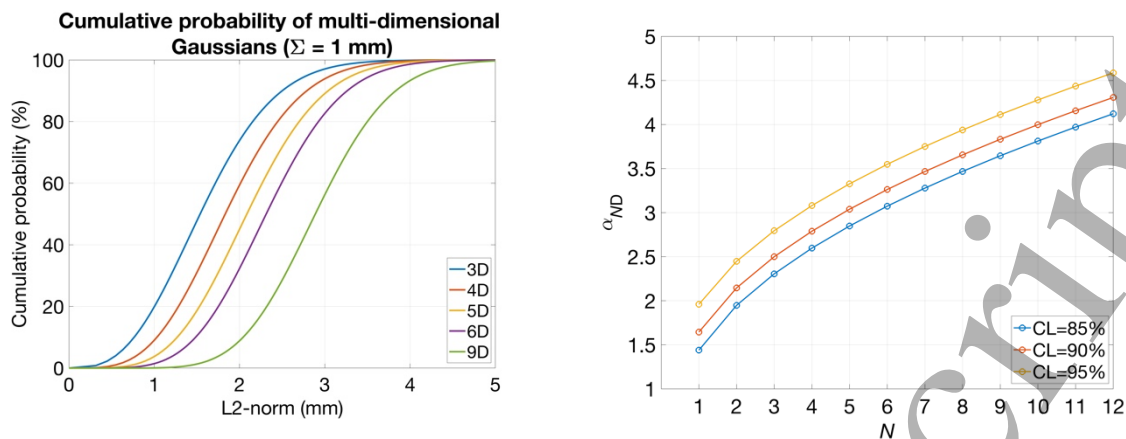
In current practice, robustness evaluation tools are limited to the generation of some occurrences of systematic setup and range errors according to parameters defined by the user. Random errors are typically not simulated. Van Der Voort et al. have suggested to consider random errors using empirical relations that can convert a combination of systematic and random errors into pure systematic errors<sup>20</sup>. Another method has been suggested by the group of PSI, using a relatively small subset of possible errors, a priori limited by an 85% confidence interval line<sup>1,3</sup>. For the reminder of the argument, we will assume that random errors are either neglected or converted to systematic errors as in Van Der Voort et al<sup>20</sup>.

If dose distributions are unknown, computing confidence levels in proton therapy is not as straightforward as in photon therapy. The main reason is that one cannot easily approximate the effect an error may have on the dose distributions. Consequently, each type of error needs to be considered separately. In the context of independent setup errors and range uncertainties, this leads to the sampling of errors in a 4D space with reduced axis ( $x' = \frac{x_{se}}{\Sigma_{setup,x}}$ ,  $y' = \frac{y_{se}}{\Sigma_{setup,y}}$ ,  $z' = \frac{z_{se}}{\Sigma_{setup,z}}$ ,  $RU' = \frac{RU}{\Sigma_{RU}}$ ), where  $\Sigma$  is the standard deviation. In this space, equiprobable errors will be located on the surface of a hypersphere with equation  $x'^2 + y'^2 + z'^2 + RU'^2 = \alpha_{4D}^2$ . The parameter  $\alpha_{4D}$  denotes the (reduced) radius of the hypersphere. The left side of the last equation represents a chi-square distribution with 4 degrees of freedom. The behavior of the cumulative chi-square distribution is illustrated in Figure 1 (a) for different numbers of degrees of freedom.

A confidence level in the *error space* can now be approximately computed. To ensure robustness against 90% of all possible scenarios, we need to select all possible configurations within a hypersphere with radius of approximately 2.8 as seen from Figure 1. If we hypothesize that the worst-case scenarios are located on the surface of the hyper-sphere, then one can assume that this confidence level of 90% will be achieved by only simulating the points distributed over the hyper-sphere. However, this hypothesis is not necessarily true and will be tested in one of the robustness evaluation strategies introduced in section II.D

If range uncertainties are removed, we come back to the 3D case and  $\alpha_{3D}$  equals the well-known 2.5 value. Figure 1 (b) displays how  $\alpha_{ND}$  varies depending on the number of dimensions. It is a direct translation of the value of the L2 norm in Figure 1 (a) at 90% cumulative probability.

One thing important to note here is that the selection of the scenarios will strongly depend on the dimensionality of the problem. More extreme scenarios will have to be selected for a higher number of dimensions and a fixed confidence level (because of the corresponding increase of the radius  $\alpha$  of the hyper-sphere).



**Figure 1 (a)** Examples of cumulative probabilities for isotropic multi-dimensional independent normal (Gaussian) distributions with 1 mm standard deviation; **(b)** values of the reduced radius  $\alpha_{ND}$  of the isoproability hyper-sphere with respect to the number of dimensions for 85, 90, and 95% confidence levels (CL). For 3 dimensions and 90% confidence level (CL),  $\alpha_{3D}$  equals the typical 2.5 value found in margin recipe of van Herk et al<sup>14</sup>.

### II.C. Patient test cases

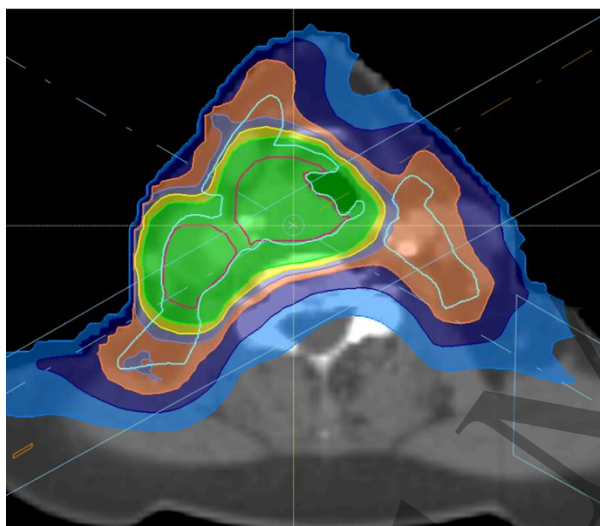
Five bilateral head-and-neck patients were considered for illustrating the notions described above. Some tumor characteristics are detailed in the appendix (Table S 1). The patients were treated by conventional radiotherapy. Hence, the proton treatment plans were optimized for the purpose of this study. The target was the PTV, obtained by expanding the CTV by a 4 mm isotropic margin. The treatment plans included two prescriptions, 70 Gy and 54 Gy on tumor and elective volumes, respectively. The proton treatment plan was composed of 4 scanned beam incidences ((350,60):(350,120):(10,240):(10,300) in degrees for couch and gantry angles, respectively). Treatment plans were optimised to ensure adequate coverage of the PTV, without robustness parameters (treatment plans were not robustly optimized). The minimum requirements were  $D_{98} > 90\%$  of prescription dose,  $D_{95} > 95\%$  of prescription dose,  $D_5 < 105\%$  of prescription dose. However, when possible to respect OAR constraints, we tried to achieve at least 95% of prescribed dose for  $D_{98}$ . Constraints to OARs were set according to the clinical rules of Cliniques Universitaires Saint-Luc used for conventional photon therapy. The OARs subject to sparing and their associated dose limits are listed in Table S 2 in the supplemental material. When possible, the dose to OARs were further diminished provided that it did not compromise PTV coverage. The treatment plans were optimised using RayStation (from RaySearch, research license 5.99). The achieved dose distributions and the used beam angles are illustrated in axial slices for each patient in Figure 2

The spot positions and weights were exported to a local robustness evaluation tool developed by Souris et al<sup>18</sup>. This robustness evaluation tool is based on a validated Monte Carlo dose engine called

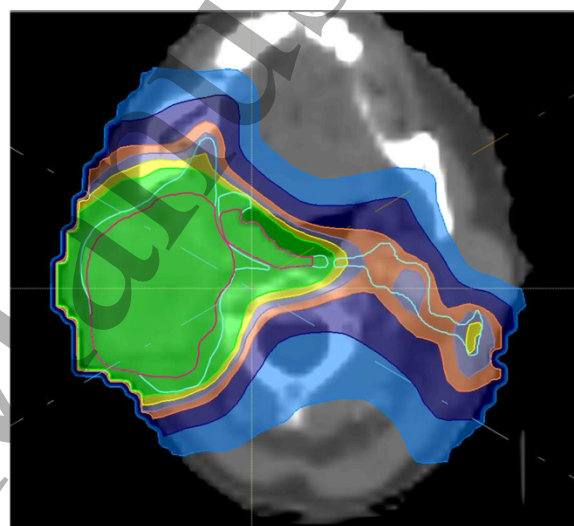


1  
2  
3 MCsquare<sup>21</sup>. For the purpose of the present study, systematic setup errors and image conversion  
4 errors were simulated by shifting the patient and applying a density scaling according to sampled  
5 values of setup errors and image conversion errors.  
6  
7

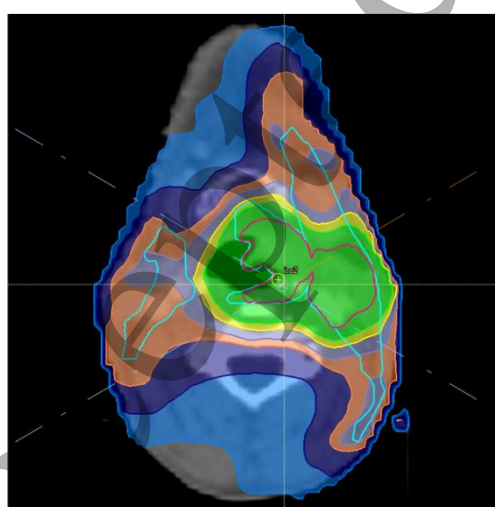
8 The values chosen for the standard deviations were as follows. For the tests without random  
9 errors,  $\Sigma_{\text{setup}} = 1.6 \text{ mm}$ ,  $\sigma_{\text{setup}} = 0 \text{ mm}$ , and  $\Sigma_{\text{RU}} = 1.8\%$ . The values were chosen in order to  
10 represent 4 mm and 3% errors at 90% confidence level in their respective spaces (3D for setup errors  
11 ( $\alpha_{3\text{D}} = 2.5$ ), 1D for range uncertainties ( $\alpha_{4\text{D}} = 1.67$ )). For the tests with random errors, the values  
12 chosen were  $\Sigma_{\text{setup}} = 1.3 \text{ mm}$ ,  $\sigma_{\text{setup}} = 1.0 \text{ mm}$ , and  $\Sigma_{\text{RU}} = 1.8\%$ . Such combination of systematic  
13 and random setup errors leads to a margin of 4 mm using the simplified van Herk formula ( $2.5\Sigma +$   
14  $0.7\sigma$ ). It is also in line with the empirical relationships shown in Figure 3 of van der Voort et al<sup>20</sup>.  
15  
16  
17  
18  
19  
20  
21



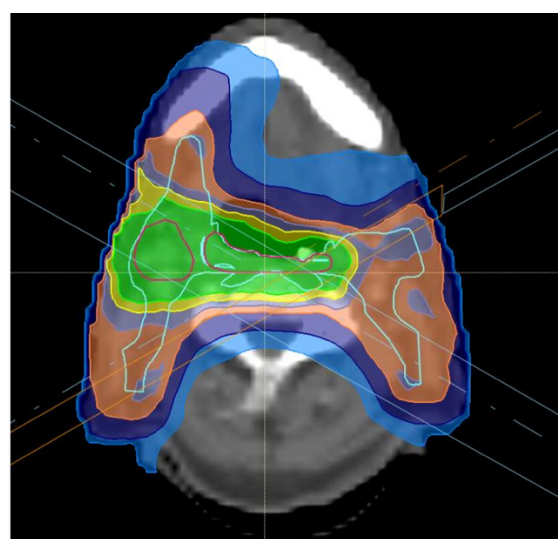
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38 Patient 1  
(a)



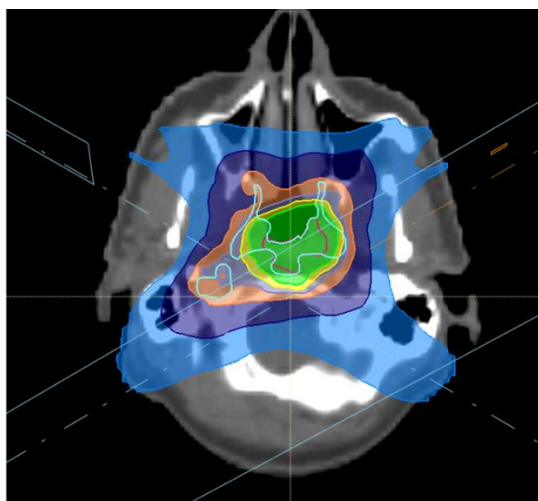
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57 Patient 2  
(b)



58  
59  
60 Patient 3  
(c)



Patient 4  
(d)



Patient 5  
(e)

Figure 2 Axial slices of the 5 patients selected for this study with overlaid dose distributions

## II.D. Robustness evaluation strategies investigated

We summarize here the robustness evaluation strategies investigated. A short overview is also given in Table 1. In all robustness evaluation strategies, the nominal scenario is kept in the simulated set of dose distributions.

### II.D.1. Strategy 1: Good practice scenario selection (GPSS) of flat systematic setup and range errors

In many robust optimization/evaluation approaches, scenarios are selected pragmatically according to good practice rules. In general, the CTV to PTV margin is replaced with a systematic setup error of comparable magnitude and the range uncertainty parameter takes typically three values, +RU, 0 and -RU where RU ranges from 2.5 to 3.5 % in most publications. Random errors are typically ignored or converted into systematic errors, for instance using the approach developed by Van der Voort et al<sup>20</sup>. For this strategy, the setup error and RU parameters equalled 4 mm and 3% (consistent with  $\Sigma_{\text{setup}} = 1.6 \text{ mm}$  and  $\Sigma_{\text{RU}} = 1.8\%$ ). In typical clinical practice, only a few scenarios are sampled in the directions  $x$ ,  $y$ , and  $z$ , i.e. positive and negative extreme values along each axis (no diagonals). By combining with range errors, it amounts to 20 scenarios in total, excluding the nominal scenario. However, it is not possible with such strategy to estimate a confidence level with acceptable accuracy, as the errors in the spatial directions  $x$ ,  $y$ , and  $z$  are sampled too coarsely. Therefore, we have simulated

1  
2  
3 more scenarios by including those on the diagonals between the  $x$ ,  $y$ , and  $z$  axes. In such case, the  
4 setup errors are selected on the 3D-sphere, at 90% confidence level in 3D (using  $\alpha_{3D}=2.5$ ). The total  
5 amount of scenarios then reaches 80 without the nominal scenario.  
6  
7

8 In this configuration, a confidence level can be estimated by integrating the joint probability  
9 density function inside the 4D hyper-cylinder defined by the 3D setup errors (distributed over a sphere)  
10 and the range errors. This was approximated numerically by generating randomly setup and range  
11 errors and counting the ones that are inside the hyper-cylinder. This amounts to 81% of possible errors.  
12 It is important to mention here that this way of computing the confidence level assume continuity of  
13 the errors in the error space and also that the worst errors are located on the edges of the explored  
14 space.  
15  
16  
17  
18  
19

20 For the sake of completeness, we have also simulated the GPSS case with 20 scenarios only.  
21 The results are reported in supplementary materials.  
22  
23  
24

### 25 **II.D.2. Strategy 2: Statistically sound scenario selection (SSSS)**

26  
27  
28 Two configurations were tested in this study. In the first configuration, scenarios were sampled  
29 uniformly *on* the hyper-surface of the 4D hyper-sphere delimited by the equation  $x'^2 + y'^2 + z'^2 +$   
30  $RU'^2 = \alpha_{4D}^2$ , where  $\alpha_{4D} = 2.8$  to ensure a 90% confidence level in the error space (SSSS (ON) Figure  
31 3 (b)). In such case, one may assume that this confidence level is secured in the error space provided  
32 that robustness for scenarios inside the hyper-surface is also warranted. In the second configuration,  
33 scenarios were *also* uniformly sampled *within* the hyper-sphere, in order to better approximate a true  
34 90% confidence level<sup>19</sup> computed in the error space (Figure 3 (c)). In SSSS (IN), we also sample  
35 hypersurfaces within the 90% hyper-surface with a different radius. The number of scenarios per  
36 surface is 80 ( $3^4$  minus the nominal case). In SSSS (IN), we sample 3 hypersurfaces (at (reduced) radii  
37 2.2 and 1.1) hence 240 scenarios. One can note that errors and scenarios lead eventually here to the  
38 same meaning, because only systematic errors are sampled.  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48

### 49 **II.D.3. Strategy 3: Statistically sound dosimetric selection (SSDS)**

50  
51  
52 We consider here a Monte Carlo robustness evaluation tool, i.e. errors are randomly sampled  
53 according to their *pdfs*. It is worth mentioning that the dose engine associated with this tool can be  
54 anything, either Monte Carlo or analytical. A random error sampling approach would be an excellent  
55 candidate for performing robustness evaluation because 1) errors can be sampled without any  
56 statistical bias from their actual *pdfs*; 2) random errors can be simulated naturally; and 3) it enables an  
57  
58  
59  
60

evaluation of the confidence level in the dosimetric space. A weakness of this approach is that the number of treatment scenarios to simulate may be substantial. To ensure its practical viability, dose computation must be performed at a low computational cost. Fast Monte Carlo dose engines may be used for this task, but, in such case, the number of errors and scenarios to simulate must be limited to what is necessary. Therefore, this requires the introduction of a convergence criteria and variance reduction techniques, as described in Souris et al<sup>18</sup>. For the purpose of this study, we have tried to minimize the statistical noise as much as reasonably achievable. In Souris et al<sup>18</sup>, it was shown that 300 scenarios were sufficient to ensure convergence of the DVH error bands. In the present study, we have therefore simulated 1000 scenarios for ensuring low noise levels on the reported values (for instance, the lowest  $D_{95}$  value at 90% confidence level). The number of particles per scenario was about  $10^8$  to ensure a statistical uncertainty in the target below 2% (one standard deviation). Simulations were performed on a 2x Intel(R) Xeon(R) Gold 6248 CPU.

The SSDS method allows flexibility in the way the scenarios are selected. We implemented two scenario selection methods. In the  $D_{95}$  method, the 90% best scenarios according to target coverage for the high dose CTV (quantified here by the  $D_{95}$ ) were selected for reporting. In the OF method, the 90% best scenarios according to the value of the objective function were selected for reporting. The value of the objective function was computed as a weighted sum of all clinical objectives used in the TPS for the treatment plan optimization. Four objective types, namely minimum dose, maximum dose, mean dose, and DVH objectives, were implemented in the objective function using quadratic terms as described in Oelfke et al<sup>22</sup> (see Table S 2 in the supplementary material).

For each scenario selection method in the dosimetric space, two tests were performed. In the SSDS (S) strategy, only systematic errors were considered. In the SSDS (R) strategy, both systematic and random errors were considered. The standard-deviations selected for both examples are provided in section II.C.

Robustness evaluation strategy	Description	$\Sigma_{\text{setup}}$ (mm)	$\sigma_{\text{setup}}$ (mm)	$\Sigma_{\text{RU}}$ (%)
GPSS	Good practice scenario selection in the error space: selection of setup errors onto 90% 3D sphere, and a positive and a negative range value	1.6	0.0	3.0*
SSSS (ON)	Statistically sound selection in the error space onto 90% isoprobability line of the 4D hypersphere	1.6	0.0	1.8
SSSS (IN)	Statistically sound selection in the error space onto and inside 90% isoprobability surface of the 4D hypersphere	1.6	0.0	1.8
SSDS $D_{95}$ (S)	Statistically sound selection in the dosimetric space for the 90% best CTV $D_{95}$	1.6	0.0	1.8

SSDS_OF (S)	Statistically sound selection in the dosimetric space for the 90% best objective function values	1.6	0.0	1.8
SSDS_D <sub>95</sub> (R)	Statistically sound selection in the dosimetric space for the 90% best CTV D <sub>95</sub>	1.3	1.0	1.8
SSDS_OF (R)	Statistically sound selection in the dosimetric space for the 90% best objective function values	1.3	1.0	1.8

**Table 1. Summary of the robustness evaluation strategies studied and their associated robustness parameters**  
 \*For GPSS, only extreme values of the distributions are considered for range errors, not the standard deviations.

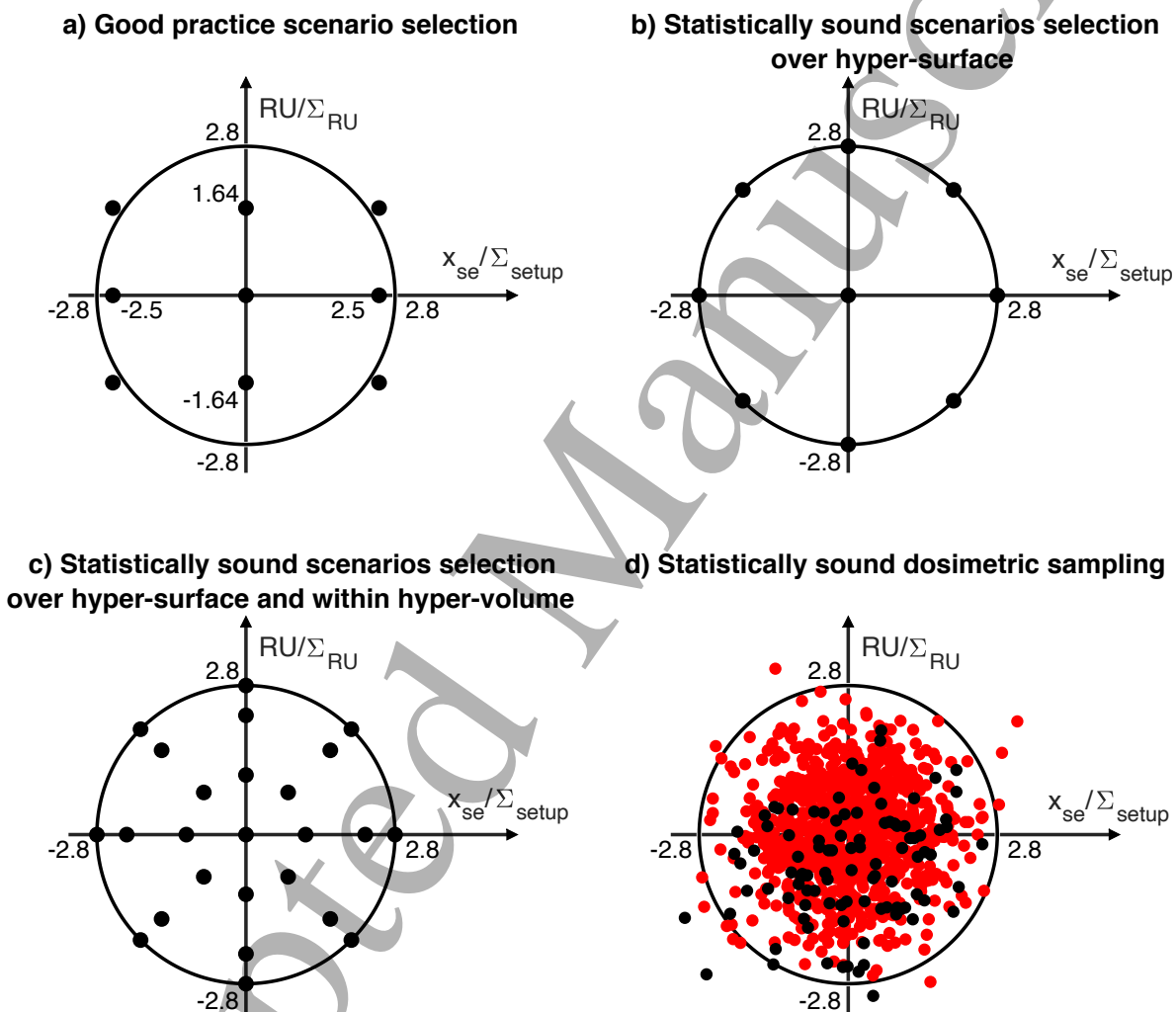


Figure 3 Illustration of the scenario selection methods in a 2D slice (a,b,c), or in a projection (d), of the 4D iso-probability hypersphere. Each dot represents a simulated scenario. The hypersphere contains 90% of all possible scenarios. Figure (a) corresponds to good practice scenario selection (GPSS). Figure (b) corresponds to statistically sound selection of the scenarios (SSSS) at the surface of the isoprobability hyper-sphere including 90% of possible scenarios (SSSS (ON)). Figure (c) corresponds to SSSS at the surface and within the isoprobability hyper-sphere including 90% of possible scenarios (SSSS (IN)). Figure (d) corresponds to statistically sound dosimetric sampling with selection on best 90% D<sub>95</sub> (SSDS\_D<sub>95</sub>) with a projection of all

scenarios onto the plane defined by the  $x$  and  $range$  dimensions. Red dots correspond to selected dose distributions; black dots to discarded dose distributions.

### III. Results

#### III.A. Results for the nominal plans

The results obtained for PTV coverage, quantified by the metrics  $D_{98}$ ,  $D_{95}$  and  $D_5$ , in the nominal plan using MCsquare are provided in Table 2. This computation was necessary to ensure that the dose distributions in the nominal configuration computed by MCsquare met target coverage criteria.

PTV metric (Gy)	Patient results for high dose PTV 70Gy				
	P1	P2	P3	P4	P5
$D_{98}$	67.6	64.9	67.2	67.9	66.7
$D_{95}$	68.2	67.0	67.9	68.4	67.5
$D_5$	71.6	73.4	71.7	72.2	71.6

Table 2 Metric assessing PTV dose coverage for the nominal plan used in this study. The dose was computed with MCsquare in the nominal case. Target coverage objectives were at least  $D_{95} > 95\%$  and  $D_{98} > 90\%$  of prescribed dose (70 Gy), thus 66.5 Gy and 63 Gy, respectively. Overdosage were limited by the constraint  $D_{95} < 105\%$  (thus 73.5 Gy). When possible, we tried to achieve  $D_{98} > 95\%$  of prescribed dose.

#### III.B. Comparison of the robustness evaluation methods

Robustness evaluation has been performed for the strategies described in Table 1. Table 3 provides the differences between worst-case and the nominal DVH metrics averaged over all patients, for each robustness evaluation strategy. Table 4 displays the same data as Table 3, this time with respect to the results yielded by the GPSS method (instead of the nominal plans in Table 3). Individual DVH metrics are illustrated for patient 3 in Figure 4, and detailed for the same patient in Table 5. The same results for the other patients are available in appendix (Tables S 3-6 and Figures S 1-4).

For each strategy, the time needed to compute one scenario was about 150 seconds.

##### III.B.1. Considering systematic errors only

As shown in Table 3 and in the individual results (Table 5, Figure 4, Tables S 3-6 and Figures S 1-4), SSSS (ON) and SSSS (IN) provide very similar DVH metrics. Therefore, they will not be distinguished anymore to present the results. For the high dose CTV, GPSS and SSSS yield similar results, with an average of worst-case  $D_{98}$ ,  $D_{95}$ , and  $D_5$  within 0.5 Gy. Results for individual patients are also similar for the high dose CTV, with differences within 0.7 Gy (Table 5 and Tables S 3-6). For the low dose CTV, GPSS and SSSS yield slightly divergent results, with average differences within 0.4 Gy and 0.3 Gy for  $D_{98}$  and  $D_{95}$ , respectively. The maximum variability occurred for patient 5, with SSSS yielding a  $D_{98}$  and a  $D_{95}$  0.9 Gy larger (Table S 6).

When comparing GPSS to the SSDS methods for target coverage (Table 4), differences are more substantial. Considering systematic errors only (S), and  $D_{98}$  of the high dose CTV, the worst-case scenario is on average 2.8 Gy and 2.0 Gy larger for  $SSDS_{D_{95}}$  and  $SSDS_{OF}$  compared to GPSS, respectively. For  $D_{95}$ , it is 2.6 and 1.7 Gy, respectively. Comparing GPSS and  $SSDS_{D_{95}}$ , the differences reported are maximum 5.0 Gy and 4.4 Gy higher for  $D_{98}$  and  $D_{95}$ , respectively (patient 2, Table S 4). For the low dose CTV, maximum average differences within 0.4 Gy are observed between both SSDS evaluation methods and GPSS. SSSS and  $SSDS_{OF}$  yield on average very similar results for the low dose target (Table 4).

These results are confirmed visually in Figure 4, where it can be noticed that DVH-bands for the high dose CTV (red) are broader for GPSS and SSSS, than for both SSDS strategies.

For organs-at-risk, the average differences reported are within 1.6 Gy for all metrics between all methods (Table 3). It is difficult to distinguish clear trends looking at individual patient results (Table 5 and Tables S 3-6). However, one can notice that GPSS often reports the lowest values for OARs.  $SSDS_{OF}$  yields in general similar or lower values than SSSS. Sometimes,  $SSDS_{D_{95}}$  (S) yields substantially larger values than other evaluation methods. For instance, for patient 2,  $D_{mean}$  of the left parotid is more than 1.5 Gy larger for  $SSDS_{D_{95}}$  than all other methods (Table S 4).

### III.B.2. Considering systematic and random errors

Simulating explicitly random errors during robustness evaluation yields similar or improved DVH metrics with respect to their counterparts with systematic error only. One can notice in Table 3 an average improved coverage of the low dose CTV up to 0.7 Gy for  $D_{98}$  ( $SSDS_{OF}$ ). For OARs, similar observations can be made, with an improvement of all OAR DVH metrics when random errors are simulated explicitly (i.e. not translated to their approximatively equivalent systematic errors). For instance, the mean to the left and right parotids improved on average in a range from 0.5 Gy to 0.9 Gy.

Strategy	# of scen	CL (%)	$D_{98}$ CTV 70Gy (Gy)	$D_{95}$ CTV 70Gy (Gy)	$D_5$ CTV 70Gy (Gy)	$D_{98}$ CTV 54Gy (Gy)	$D_{95}$ CTV 54Gy (Gy)	$D_{mean}$ left prtd (Gy)	$D_{mean}$ right prtd (Gy)	$D_{mean}$ oral cavity (Gy)	$D_2$ spinal cord (Gy)	$D_2$ brain stem (Gy)
GPSS	80	81	-4.9	-3.9	2.0	-3.9	-2.9	6.4	5.8	3.8	6.7	5.3
SSSS (ON)	80	90	-4.4	-3.6	1.8	-3.5	-2.6	6.1	5.7	3.5	7.0	5.6
SSSS (IN)	240	90	-4.4	-3.6	1.8	-3.5	-2.6	6.1	5.7	3.5	7.0	5.6
SSDS_ $D_{95}$ (S)	1000	90	-2.1	-1.3	1.9	-3.9	-3.1	6.7	6.2	4.1	7.0	5.5
SSDS_OF (S)	1000	90	-2.9	-2.2	1.6	-3.5	-2.6	5.3	5.6	3.3	7.0	5.5
SSDS_ $D_{95}$ (R)	1000	90	-2.0	-1.1	1.9	-3.4	-2.6	6.1	5.6	3.6	5.8	4.7
SSDS_OF (R)	1000	90	-2.4	-1.8	1.3	-2.8	-2.2	4.8	4.7	2.8	6.2	4.9

**Table 3 Dose differences between the worst-case and the nominal DVH metrics for the target and organs-at-risk, averaged over the 5 patients (# of scen = number of scenarios; CL = confidence level). The meaning of each robustness evaluation strategy is detailed in Table 1. The abbreviation “prtd” stands for “parotid”.**

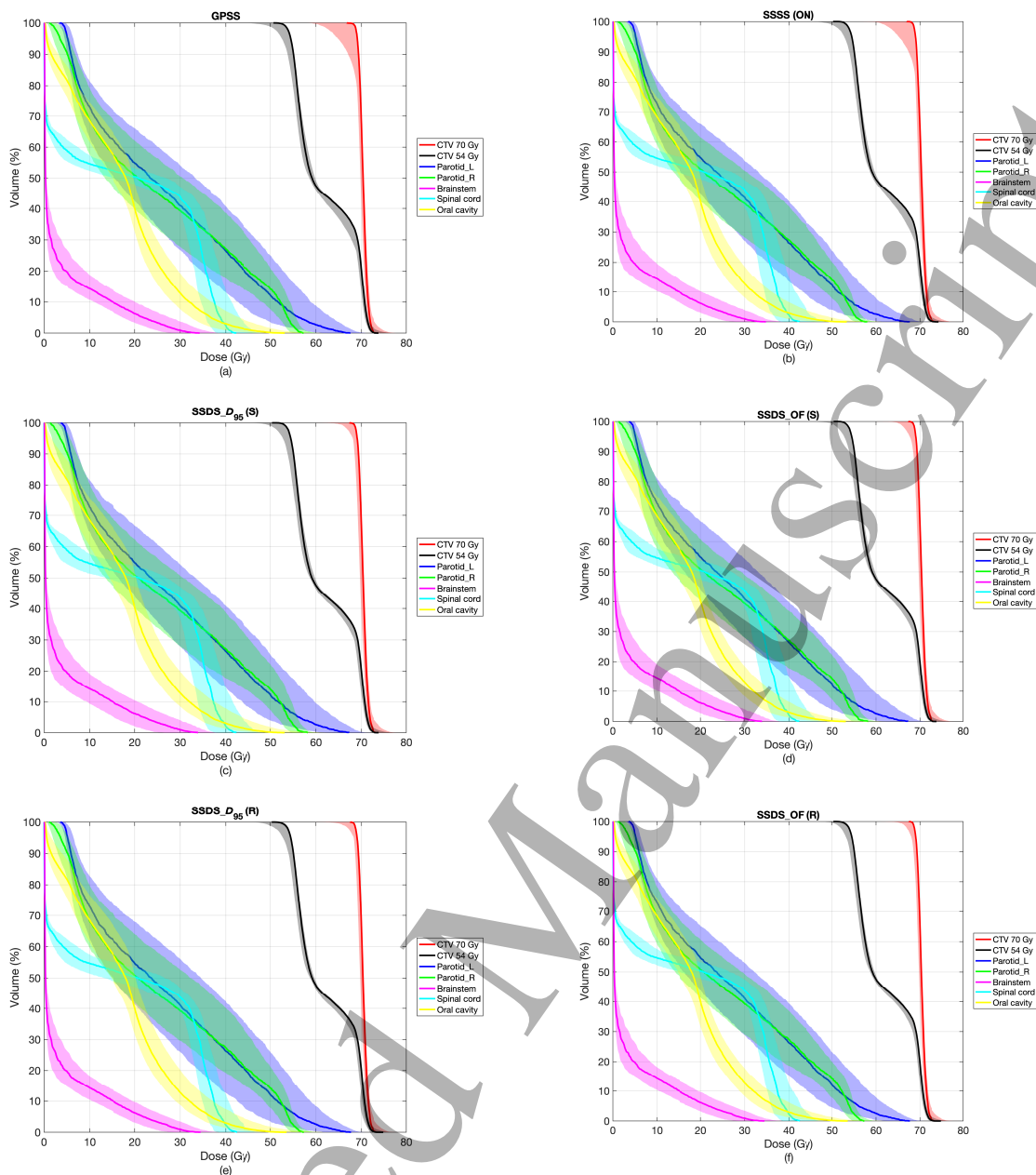
Strategy	# of scen	CL (%)	$D_{98}$ CTV 70Gy (Gy)	$D_{95}$ CTV 70Gy (Gy)	$D_5$ CTV 70Gy (Gy)	$D_{98}$ CTV 54Gy (Gy)	$D_{95}$ CTV 54Gy (Gy)	$D_{mean}$ left prtd (Gy)	$D_{mean}$ right prtd (Gy)	$D_{mean}$ oral cavity (Gy)	$D_2$ spinal cord (Gy)	$D_2$ brain stem (Gy)
SSSS (ON)	80	90	0.5	0.3	-0.2	0.4	0.3	-0.3	-0.1	-0.3	0.3	0.3
SSSS (IN)	240	90	0.5	0.3	-0.2	0.4	0.3	-0.3	-0.1	-0.3	0.3	0.3
SSDS_ $D_{95}$ (S)	1000	90	2.8	2.6	-0.1	0.0	-0.2	0.3	0.4	0.3	0.3	0.2
SSDS_OF (S)	1000	90	2.0	1.7	-0.4	0.4	0.3	-1.1	-0.2	-0.5	0.3	0.2
SSDS_ $D_{95}$ (R)	1000	90	2.9	2.8	-0.1	0.5	0.3	-0.3	-0.2	-0.2	-0.9	-0.6
SSDS_OF (R)	1000	90	2.5	2.1	-0.7	1.1	0.7	-1.6	-1.1	-1.0	-0.5	-0.4

**Table 4 Average absolute differences of the DVH metrics for the 5 patients with respect to GPSS taken as a reference ('# of scen' = number of scenarios; 'CL' = confidence level). The meaning of each robustness evaluation strategy is detailed in Table 1. The abbreviation “prtd” stands for “parotid”.**



Robustness evaluation strategy	# of scen	CL (%)	$D_{98}$ CTV 70Gy (Gy)	$D_{95}$ CTV 70Gy (Gy)	$D_5$ CTV 70Gy (Gy)	$D_{98}$ CTV 54Gy (Gy)	$D_{95}$ CTV 54Gy (Gy)	$D_{mean}$ left prtd (Gy)	$D_{mean}$ right prtd (Gy)	$D_{mean}$ oral cavity (Gy)	$D_2$ spinal cord (Gy)	$D_2$ brain stem (Gy)
<b>Worst-case</b>												
GPSS	80	81	62.7	64.8	72.4	50.9	52.7	32.4	30.1	19.1	45.0	32.6
SSSS (ON)	80	90	63.1	65.2	72.6	51.1	52.8	32.4	29.8	19.7	45.3	33.1
SSSS (IN)	240	90	63.1	65.2	72.6	51.1	52.8	32.4	29.8	19.7	45.3	33.1
SSDS_ $D_{95}$ (S)	1000	90	66.8	68.1	72.6	51.2	52.7	32.9	30.1	20.1	45.4	33.4
SSDS_OF (S)	1000	90	66.3	67.7	72.6	51.2	52.7	32.4	29.9	19.5	45.4	33.4
SSDS_ $D_{95}$ (R)	1000	90	67.0	68.3	72.6	51.0	52.7	33.5	30.7	20.4	44.3	32.4
SSDS_OF (R)	1000	90	66.9	68.0	72.4	51.9	53.1	32.1	29.5	19.2	44.9	32.4
<b>Nominal</b>												
	1	NA	68.7	69	71.9	53.8	54.4	26.2	24.9	17	40.2	27.6

**Table 5 Results of the robustness evaluation for patient 3 ('# of scen' = number of scenarios; 'CL' = confidence level). The worst-case are shown for each robustness evaluation strategy. For comparison purposes, the nominal values are also displayed. The meaning of each robustness evaluation strategy is detailed in Table 1. The abbreviation "prtd" stands for "parotid".**



**Figure 4** Results of the robustness evaluation for patient 3. The meaning of each robustness evaluation strategy (mentioned in the title of every graph) is detailed in Table 1.

#### IV. Discussion

The results show that for the patients investigated, SSSS yields the same results whether scenarios are simulated inside the isoprobability sphere or only on the surface. This is in line with previous findings<sup>4, 5</sup>. It is, however, impossible to strictly exclude that a few scenarios inside the hypersphere could lead to unexpected loss of target coverage or unexpected OAR exposure. For instance, range errors induced by setup errors and explicitly simulated range errors could compensate

1  
2  
3 for some particular points on the surface of the hypersphere but not inside, leading to eventually less  
4 perturbed dose distributions for some extreme errors. But checking the interior of the hypersphere  
5 will inevitably lead to an important increase of the scenarios to simulate (from 80 to 240 in our  
6 examples). Therefore, one may consider for practical purposes to explore only the surface of the  
7 hypersphere (i.e., the most extreme errors).  
8  
9

10  
11 A striking result is that GPSS leads to larger error bands for target coverage, smaller worst-  
12 cases doses overall to OARs, and a smaller confidence level of 81%. In practice, this may lead to the  
13 decision of replanning because of a lack of target coverage, with the inevitable downside of increasing  
14 the dose to OARs, with some that are already slightly underestimated (e.g. 0.3 Gy average difference  
15 for  $D_2$  brainstem between GPSS and SSSS). Those results are intuitively expected. Because the GPSS  
16 strategy only explores 81% of the possible scenarios (assuming robustness against intermediate errors)  
17 AND arbitrarily select extreme scenarios with a very low probability (i.e. outside the 90% hypersphere,  
18 thus inconsistent with generally accepted confidence levels (90%)), this leads to an over-conservative  
19 approach for the target (because of the extreme cases considered) and a possible under-estimation of  
20 the OARs (because of a larger number of unexplored scenarios). An additional source of inconsistency  
21 is the arbitrary selection of scenarios with different probabilities (for instance  $(x_{se}, y_{se}, z_{se}, RU)$  may  
22 equal (4 mm, 0, 0, 0) or (4 mm, 0, 0, 3%) as shown in Figure 3 (a); the first scenario is more likely to  
23 occur). In clinical practice, GPSS is often implemented differently, with a coarser selection of the  
24 scenarios in the directions  $x$ ,  $y$ , and  $z$ . In such case, the computation of a reliable confidence level  
25 becomes very problematic. However, we observe similar results for GPSS either with 20 or 80  
26 scenarios, as it can be seen by comparing Table 3 and Table S 7, which report average results within  
27 0.3 Gy for the targets and 0.8 Gy for the OARs.  
28  
29

30  
31 The SSSS method will lead overall to the most conservative approach, as shown in Table 3 and  
32 Table 4. Because of the effect of dimensionality (Figure 1), SSSS forces the exploration of scenarios that  
33 are typically not considered in clinical practice (for photons and protons), for instance an error up to  
34  $2.8\sigma_{\text{setup}}$ , which is larger than the more familiar  $2.5\sigma_{\text{setup}}$ . The effect of the dimensionality has already  
35 been addressed by Korevaar et al <sup>13</sup>. If more errors are included, for instance baseline shifts and/or  
36 rotations, the errors to explore would be more extreme as shown by Figure 4. This is a key weakness  
37 of the SSSS method. Because we are blind to the effect of the uncertainties on the dose distributions,  
38 the selection can only be performed on or within isoprobability hypersurfaces in order to ensure  
39 statistical consistency. As a consequence, the space to explore will increase with the types of errors to  
40 explore. In practical cases, the dimensionality of the error space is typically 4D, which leads to a mild  
41 increase of the errors to explore (from 2.5 to  $2.8\sigma_{\text{setup}}$ ). But if a robustness evaluation system aims at  
42 improved generalizability of the evaluation, it may need to explore more dimensions (inter-fractional  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 anatomical change, breathing variability, etc), which will inevitably lead to an explosion of the  
4 magnitude of the errors and to extremely conservative treatment plans. In the context of the PTV  
5 margin recipe, this blindness is overcome by assuming a simple hypothesis related to the dose  
6 distributions: the static-dose cloud approximation. This allows a simple sum of the associated random  
7 variables – i.e. quadratic sum of variances in margin recipe – so that the problem remains a 3D  
8 problem. This hypothesis is rightly forbidden in proton therapy, hence the dimensionality problem that  
9 appears here.

15 The SSDS methods aim precisely at overcoming the downsides of GPSS (inconsistency and  
16 arbitrariness) and SSSS (over-conservativeness) discussed above. Because the problem under  
17 consideration is eventually a 3D problem (dose distributions are 3D objects), it is more powerful to  
18 explore the scenarios in the dosimetric space. In such case, all the potential redundancies in the error  
19 space will be captured. Moreover, extreme errors that may have a low impact on the dose distribution  
20 (for instance, a motion parallel to a highly contributing treatment field), can be included in the DVH  
21 bands naturally. This can be observed in Figure 5, where substantial errors, outside the isoproability  
22 hypersphere, could lead to an acceptable dose distribution. Because what is important in the end is  
23 the *confidence level* (i.e., the probability of meeting a criterion or not), statistically unlikely errors can  
24 be included safely provided that the final probability (or confidence level) is correctly computed. This  
25 leads to a more optimistic estimation of target coverage (2.3 and 1.5 Gy higher on average for  $D_{98}$  of  
26 the high dose CTV, for SSDS\_ $D_{95}$  (S) and SSDS\_OF (S) compared to SSSS, respectively). And a mild  
27 increase (for SSDS\_ $D_{95}$  or decrease (for SSDS\_OF) of DVH metrics of OARs within 0.8 Gy (on average  
28 over the 5 patients) compared to SSSS. It is interesting also to mention that such considerations were  
29 already addressed for establishing confidence levels for PTV margin recipes. In van Herk et al 2000<sup>14</sup>,  
30 it is written that ‘the margin for treatment preparation (systematic) errors is chosen as a confidence  
31 interval that is spherically symmetric. However, an infinite number of 90% confidence intervals may  
32 be chosen that are not spherically symmetric. This observation leaves some room for optimization.’ In  
33 a follow-up paper, Witte<sup>23</sup> (2017) showed by Monte Carlo simulations how the margin can be  
34 optimized to reduce OAR dose while maintaining minimum CTV dose.

48 However, a new problem that arises is the adequate selection of the scenarios in the  
49 dosimetric space. In other words, what is the worst dose distribution? How do we define “worst”?  
50 Table 3 and Table 4 show that the reported worst-case will differ significantly depending on the  
51 scenario selection method. If we focus on target coverage and select the 90% best  $D_{95}$  (SSDS\_ $D_{95}$ ), we  
52 obtain the most optimistic result for high dose target coverage, at the expense of generally higher DVH  
53 metrics for OARs. Such approach would be ideal in cases with no compromise with respect to OARs.  
54 We would then achieve the best estimate of target coverage, for a confidence level of 90%. However,  
55 if there are compromises to be made with OARs, then the worst-case dose to OAR will be on the  
56  
57  
58  
59  
60

1  
2  
3 pessimistic side, which may lead to exceed clinical constraints causing the reoptimization of a plan and  
4 eventually a deterioration of target coverage.  
5

6 A solution to the issue of scenario selection based on target coverage only would be to capture  
7 the clinical compromise made at the planning level and display the 90% *best* dose distributions, with  
8 respect to both target coverage and OAR sparing. We propose here to achieve this by computing for  
9 each scenario the *objective function* as accepted by the radiation oncologist and the medical physicist  
10 before robustness evaluation. The objective function provides a quantitative assessment of the quality  
11 of the plan from a clinical point-of-view, since it integrates clinical objectives and constraints, as well  
12 as objective function weights used for optimization that are implicitly approved by the radiation  
13 oncologist. Such approach could also naturally be translated to a model-based dose distribution  
14 assessment, for instance using tumor control and normal tissue complication probabilities.  
15

16 The SSDS\_OF method yields less optimistic numbers for high dose target coverage than  
17 SSDS\_ $D_{95}$ , but those are still significantly larger than GPSS and SSSS (2.0 and 1.5 Gy larger for  $D_{98}$  on  
18 average, respectively). However, the results obtained for OARs are on average comparable to both  
19 GPSS and SSSS. Interestingly, SSDS\_OF also yields results for the low dose target comparable to SSSS.  
20 Therefore, SSDS\_OF seems to better capture the plans that will lead to the best clinical compromises.  
21

22 One potential issue of the SSDS\_OF method is that objective functions vary by nature from one  
23 patient to another depending on the tuning of objective/constraint weights in order to achieve a  
24 clinically acceptable compromise between target coverage and OAR sparing. This may lead to  
25 undesired variability in robustness reporting. However, such feature could also be seen as an  
26 advantage. Two identical robustness evaluation results may lead to different appreciations by a  
27 radiation oncologist depending on individual patient characteristics. For instance, more attention can  
28 be given to a particular organ-at-risk in a given patient. Such patient-specific characteristics are at least  
29 partially entailed implicitly in the objective function. As a consequence, selecting the best dose  
30 distributions according to the value of the objective function will tend to be more faithful to the clinical  
31 compromises made at the treatment optimization level, and therefore reduce variability in patient  
32 reporting from a clinical perspective. Such approach also motivates the radiation oncologist to better  
33 formalize the clinical goals he/she aims to achieve before the robust optimization phase starts. This is  
34 in line with an improved standardization of the treatment planning workflow, which is essential for its  
35 automation.  
36

37 It is important to note that the computation of confidence levels in the dosimetric space has  
38 already been illustrated by Perko et al <sup>19</sup> with the polynomial chaos expansion. Perko et al have also  
39 identified the potential of working in the dosimetric space to estimate the magnitude of the errors to  
40 be included in the robustness evaluation to achieve given statistical criteria, e.g. coverage of the CTV  
41 in a given fraction of the patients. The only requirement to work in the dosimetric space is to have a  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 fast dose engine available in order to generate enough dose distributions to compute statistical  
4 quantities. This is exactly the purpose of the polynomial chaos expansion method that proposes a novel  
5 approach to generate a virtually infinite number of dose distributions after taking the time to generate  
6 a comprehensive dose calculation model (based on about 100 pre-computed scenarios). In our work,  
7 we use a fast Monte Carlo dose engine associated with statistically defined stopping criteria to  
8 generate the required scenarios. Another difference with the study by Perko et al, is that the authors  
9 evaluate the robustness for each volume of interest separately, while we attempt to evaluate methods  
10 to select scenarios globally. The approach of Perko et al could be trivially adapted to our methodology.  
11 An advantage of a global approach is that it naturally takes into account correlations between the DVH  
12 metrics since a set of dose distributions is selected.  
13  
14  
15  
16  
17  
18  
19

20 The explicit simulation of random errors leads to results that are on average more optimistic  
21 than their counterparts with systematic errors only. We remind here that we have always used sets of  
22  $(\Sigma, \sigma)$  that lead to a consistent CTV-to-PTV margin of 4 mm using the simplified formula of van Herk et  
23 al  $(2.5\Sigma + 0.7\sigma)$ . This indicates that this formula might be overconservative for the patients  
24 investigated in this study. More aggressive plans could therefore be achieved using a statistically sound  
25 robustness evaluation method that includes random errors. For instance, SSDS\_OF (R) yields a worst-  
26 case  $D_{98}$  for low dose target that is on average 0.5 Gy higher than SSDS\_OF (S). For the right parotid,  
27 the worst-case  $D_{\text{mean}}$  is 0.9 Gy lower for SSDS\_OF (R) than SSDS\_OF (S). It is interesting to compare  
28 SSDS\_OF (R) with GPSS by analyzing the last line of Table 4. SSDS\_OF (R) estimates a better target  
29 coverage, overall more optimistic organ-at-risk sparing, and all this for a higher confidence level (90%  
30 versus 81%).  
31  
32  
33  
34  
35  
36  
37

38 It is not the purpose of this paper to suggest a procedure for robustness evaluation. First of all,  
39 such procedures will strongly depend on the tumor site considered, the advancement of computing  
40 technology, the number of effects we want to consider, and clinical practice. For instance, the group  
41 at PSI has suggested a robustness evaluation procedure built up across many publications that is well  
42 suited for locations with small systematic errors<sup>5</sup>. The computation of confidence levels was also  
43 included for the effect of fractionation<sup>3</sup>. Other groups have suggested to include variable  
44 radiobiological models in their evaluation<sup>24</sup>. However, most robustness evaluation strategies reported  
45 in the literature select separately setup errors and range errors according to good practice rules,  
46 without considering the computation of confidence levels, neither in the error space nor in the  
47 dosimetric space<sup>8, 16, 25</sup>. As mentioned before, Perko et al do compute appropriately confidence levels  
48 in the dosimetric space using the polynomial chaos expansion method<sup>19</sup>. Finally, we have reported  
49 here worst-case DVH metrics for both target volumes and OARs. One could argue that for parallel-like  
50 OAR, like lungs, DVH metrics averaged over the entire set of dose distributions could be more  
51 meaningful. In such case, the problem is made trivial for our SSDS methods since we can simply average  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 all DVH metrics over all simulated scenarios. SSSS (IN) should also work. However, adaptations will be  
4 required for GPSS and SSSS (ON) since those sample only extreme scenarios, whilst the accurate  
5 computation of average DVH metrics would require also intermediate values.  
6  
7

8 The choice of a robustness evaluation procedure entails also pragmatic considerations such as  
9 the time needed to execute the procedure. The SSDS methods are time consuming because enough  
10 scenarios need to be simulated in order to minimize the impact of the statistical noise on the reported  
11 values. In Souris et al<sup>18</sup>, about 300 scenarios seemed adequate to ensure convergence of the results.  
12 An intrinsic advantage of Monte Carlo simulations is that the computation time does not scale  
13 necessarily with complexity. For instance, random errors can be simulated comprehensively with  
14 minimal impact on computation time. Yet, we report here 153 s computation time per scenario, which  
15 leads to a total computation time of 13 h for 300 scenarios, which is the maximum limit one may  
16 consider in clinical practice (this would correspond to calculations performed overnight). However,  
17 such computation time would only be acceptable for a final check, but not for an iterative approach  
18 where treatment plans are re-optimized several times according to the results of the robustness  
19 evaluation. Therefore, significant improvements are needed to warrant dosimetric selection of  
20 scenarios in the clinical practice. This may be achieved by improving the speed of the Monte Carlo dose  
21 engine, or the introduction of variance reduction techniques for enabling more efficient sampling of  
22 the scenarios, as suggested in Souris et al<sup>18</sup>. The polynomial chaos expansion method can also be used  
23 to reduce somewhat the number of dose computations needed, and hence speed-up the overall  
24 process<sup>19</sup>.  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35

36 The distinction between the error space and the dosimetric space has been made in the  
37 current study for protons only. In general, such distinction is not made in photon therapy because of  
38 the usual hypothesis of shift invariance of the dose distributions. If the hypothesis is true, the issue of  
39 robustness for target coverage can be formulated as a geometric problem, which leads to safety  
40 margin recipes. However, such hypothesis is not necessarily true (for instance, misplaced shoulders in  
41 head-and-neck tumors that cause undesired attenuation). Therefore, photon-based treatment plans  
42 could also benefit from comprehensive robustness evaluation strategies, which would also help for  
43 defining common dose metrics to evaluate proton and photon plans. One can also note that photon-  
44 based plans may still benefit from a comprehensive robustness evaluation in the dosimetric space  
45 under the hypothesis of shift-invariance of the dose distributions, for instance to reveal robustness  
46 improvements due to non-perfect conformity to the target, or to generate DVH-bands using advanced  
47 metrics like the value of the OF.  
48  
49  
50  
51  
52  
53  
54  
55

56 Finally, it is important to mention that the results presented here were achieved using PTV-  
57 based treatment plans, that is, non-robustly optimized. Many papers have shown that robust  
58 optimization is more suitable to ensure adequate plan robustness<sup>10</sup>. Qualitatively, our conclusions  
59  
60

1  
2  
3 should remain valid if we apply our robustness evaluation methods to robust optimized plans, although  
4 this must be confirmed in further studies. Quantitatively, robust optimization is expected to mitigate  
5 the differences observed during the present study between the various robustness strategies.  
6  
7

8  
9 However, complex treatment plans with adjacent target volumes and OARs might lead to  
10 challenging clinical trade-offs, even in the context of robust optimization. In such case, having at one's  
11 disposal a statistically fair and comprehensive evaluation strategy will help to provide the patients with  
12 the best treatment plans, with improved safety. Another limitation of our study resides in the  
13 computation of the objective function in the evaluation phase. We have tried to reproduce the best  
14 we could the objective function used in the RayStation. However, hidden terms or unforeseen  
15 mathematical expressions could be used in the RayStation's objective function and would not be  
16 captured by our computation. It would be interesting to compare our results for SSDS\_OF to those that  
17 would be obtained using the objective function used within the RayStation. Another option would be  
18 to design objective functions exclusively for evaluation.  
19  
20  
21  
22  
23  
24  
25  
26

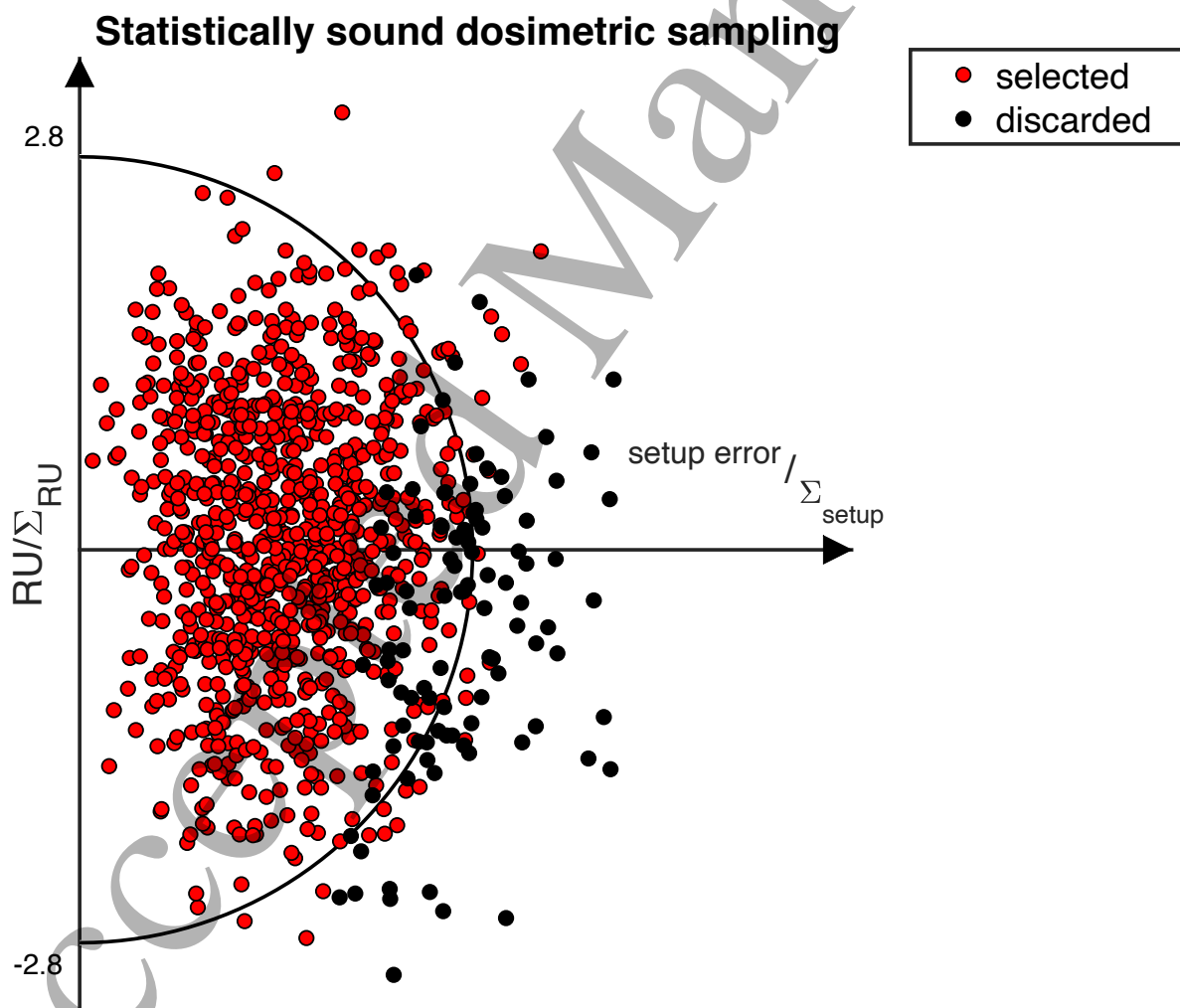




Figure 5 Representation of the scenarios selected by the robustness evaluation  $SSDS_{D_{95}}$ . Each scenario is represented with respect to the simulated range uncertainty and the norm of the sampled setup error ( $se$ )  $\sqrt{x'^2 + y'^2 + z'^2}$  where  $x' = \frac{x_{se}}{\Sigma_{setup}}$ ,  $y' = \frac{y_{se}}{\Sigma_{setup}}$ ,  $z' = \frac{z_{se}}{\Sigma_{setup}}$ . Black dots are excluded from the DVH-band while red dots are included in the DVH band. The red dots represent the best 90% values of the  $D_{95}$ . The large black circle represents all equiprobable configurations in the (4D) error space (90% confidence level in the error space). Diagram not to scale.

## V. Conclusions

Robustness evaluation is a critical step in proton therapy treatment planning. Typically, we aim at evaluating worst-case scenarios within a reasonable set of possible treatment errors. Depending on the outcome of the robustness evaluation, treatment plan optimization may be resumed for enhancing the quality of the plan in terms of target coverage and/or organs-at-risk dose. Therefore, the information delivered by the chosen robustness evaluation strategy must be as accurate and as comprehensive as possible.

We have provided several ways to evaluate statistically the robustness of the plan. An approach based on good practice rules, typically used in current clinical practice, is overall pessimistic for target coverage and optimistic for organs-at-risk sparing, with a relatively low confidence level (81%). Exploring the possible scenarios in the error space in a statistically consistent fashion enables a larger and more familiar confidence level (90%), but at the cost of conservative evaluations of worst-case DVH metrics.

Another approach would be to select scenarios in the dosimetric space, i.e. to select the best dose distributions according to *a priori* defined clinical criteria. Focusing on target coverage provides considerably more optimistic target coverage metrics (and mildly pessimistic OAR sparing). This would probably be a good approach when OAR sparing is easily achievable, and one wants to deliver the most conformal dose possible to achieve target coverage for a given confidence level. A more balanced approach would be to classify the best dose distributions according to the value of the objective function accepted by the radiation oncologist. In such case, a good balance is obtained between the reported worst-case target coverage and OAR sparing. Such approach could be easily implemented in existing commercial solutions.

## VI. Acknowledgements

This work is partially inspired from discussions within the European Particle Therapy Network working group 5. Kevin Souris is funded by the Walloon region (MECATECH/BLOWIN, grant number 8090). Sara T Rivas is supported by the Walloon region ("Convention hors pôles ProTherWal", grant number 7289). J. A. Lee is a Senior Research Associate with the Belgian fund of scientific research

(F.R.S.-FNRS). Ben George is supported by a Cancer Research UK Centres Network Accelerator Award Grant (A21993) to the ART-NET consortium. Edmond Sterpin's research is partially supported by "Fonds Baillet-Latour".

## VII. References

- 1 F. Albertini, E.B. Hug, and A.J. Lomax, Is it necessary to plan with safety margins for actively scanned proton therapy?, *Phys. Med. Biol.* **56**(14), 4399–4413 (2011).
- 2 A. Fredriksson, A characterization of robust radiation therapy treatment planning methods— from expected value to worst case optimization, *Med. Phys.* **39**(August), 5169 (2012).
- 3 M. Lowe, F. Albertini, A. Aitkenhead, A.J. Lomax, and R.I. Mackay, Incorporating the effect of fractionation in the evaluation of proton plan robustness to setup errors, *Phys. Med. Biol.* **61**(1), 413–429 (2015).
- 4 M. Casiraghi, F. Albertini, and A.J. Lomax, Advantages and limitations of the “worst case scenario” approach in IMPT treatment planning., *Phys. Med. Biol.* **58**(5), 1323–1339 (2013).
- 5 R. Malyapa, M. Lowe, A. Bolsi, A.J. Lomax, D.C. Weber, and F. Albertini, Evaluation of Robustness to Setup and Range Uncertainties for Head and Neck Patients Treated with Pencil Beam Scanning Proton Therapy, *Int. J. Radiat. Oncol. Biol. Phys.* **95**(1), 154–162 (2016).
- 6 A. Fredriksson, Automated improvement of radiation therapy treatment plans by optimization under reference dose constraints., *Phys. Med. Biol.* **57**(23), 7799–7811 (2012).
- 7 R. Bokrantz and A. Fredriksson, Controlling Robustness and Conservativeness in Multicriteria Intensity-Modulated Proton Therapy Optimization Under Uncertainty, (2013).
- 8 W. Liu, X. Zhang, Y. Li, *et al.*, Robust optimization of intensity modulated proton therapy, **1079**(2012), (2014).
- 9 J.Y. Chang, H. Li, X.R. Zhu, *et al.*, Clinical implementation of intensity modulated proton therapy for thoracic malignancies., *Int. J. Radiat. Oncol. Biol. Phys.* **90**(4), 809–18 (2014).
- 10 J. Unkelbach, M. Alber, M. Bangert, *et al.*, Robust radiotherapy planning, *Phys. Med. Biol.* **63**(22), (2018).
- 11 D. De Ruysscher, E. Sterpin, K. Haustermans, and T. Depuydt, Tumour movement in proton therapy: Solutions and remaining questions: A review, *Cancers (Basel)*. **7**(3), 1143–1153 (2015).
- 12 S. Ge, Z. Liao, J. Yang, *et al.*, Potential for Improvements in Robustness and Optimality of Intensity-Modulated Proton Therapy for Lung Cancer with 4-Dimensional Robust Optimization, *Cancers (Basel)*. **11**(1), 35 (2019).
- 13 E.W. Korevaar, S.J.M. Habraken, D. Scandurra, *et al.*, Practical robustness evaluation in radiotherapy – A photon and proton-proof alternative to PTV-based plan evaluation, *Radiother.*

- 1  
2  
3 Oncol. **141**, 267–274 (2019).  
4  
5 14 M. van Herk, P. Remeijer, C. Rasch, and J. V Lebesque, The probability of correct target dosage:  
6 dose-population histograms for deriving treatment margins in radiotherapy., *Int. J. Radiat.*  
7 *Oncol. Biol. Phys.* **47**(4), 1121–35 (2000).  
8  
9  
10 15 M. Stuschke, A. Kaiser, C. Pöttgen, W. Lübcke, and J. Farr, Potentials of robust intensity  
11 modulated scanning proton plans for locally advanced lung cancer in comparison to intensity  
12 modulated photon plans., *Radiother. Oncol.* **104**(1), 45–51 (2012).  
13  
14  
15 16 W. Liu, S.E. Schild, J.Y. Chang, *et al.*, Exploratory Study of 4D versus 3D Robust Optimization in  
16 Intensity Modulated Proton Therapy for Lung Cancer, *Int. J. Radiat. Oncol. Biol. Phys.* **95**(1),  
17 523–533 (2016).  
18  
19  
20 17 W. Liu, S.J. Frank, X. Li, *et al.*, Effectiveness of robust optimization in intensity-modulated proton  
21 therapy planning for head and neck cancers., *Med. Phys.* **40**(5), 051711 (2013).  
22  
23  
24 18 K. Souris, A. Barragan Montero, G. Janssens, D. Di Perri, E. Sterpin, and J.A. Lee, Technical Note:  
25 Monte Carlo methods to comprehensively evaluate the robustness of 4D treatments in proton  
26 therapy, *Med. Phys.* **46**(10), 4676–4684 (2019).  
27  
28  
29 19 Z. Perkó, S.R. Van Der Voort, S. Van De Water, C.M.H. Hartman, M. Hoogeman, and D.  
30 Lathouwers, Fast and accurate sensitivity analysis of IMPT treatment plans using Polynomial  
31 Chaos Expansion, *Phys. Med. Biol.* **61**(12), 4646–4664 (2016).  
32  
33  
34 20 S. Van Der Voort, S. Van De Water, Z. Perkó, B. Heijmen, D. Lathouwers, and M. Hoogeman,  
35 Robustness Recipes for Minimax Robust Optimization in Intensity Modulated Proton Therapy  
36 for Oropharyngeal Cancer Patients, *Int. J. Radiat. Oncol. Biol. Phys.* **95**(1), 163–170 (2016).  
37  
38  
39 21 K. Souris, J.A. Lee, and E. Sterpin, Fast multi-purpose Monte Carlo simulation for proton therapy  
40 using multi- and many-core CPU architectures, *Med. Phys.* **1700**, 1–23 (2016).  
41  
42  
43 22 U. Oelfke and T. Bortfeld, Inverse planning for photon and proton beams, *Med. Dosim.* **26**(2),  
44 113–124 (2001).  
45  
46  
47 23 M.G. Witte, J.J. Sonke, J. Siebers, J.O. Deasy, and M. Van Herk, Beyond the margin recipe: The  
48 probability of correct target dosage and tumor control in the presence of a dose limiting  
49 structure, *Phys. Med. Biol.* **62**(19), 7874–7888 (2017).  
50  
51  
52 24 J. Ödén, K. Eriksson, and I. Toma-Dasu, Incorporation of relative biological effectiveness  
53 uncertainties into proton plan robustness evaluation, *Acta Oncol. (Madr)*. **56**(6), 769–778  
54 (2017).  
55  
56  
57 25 S. van de Water, I. van Dam, D.R. Schaart, A. Al-Mamgani, B.J.M. Heijmen, and M.S. Hoogeman,  
58 The price of robustness; impact of worst-case optimization on organ-at-risk dose and  
59 complication probability in intensity-modulated proton therapy for oropharyngeal cancer  
60 patients, *Radiother. Oncol.* **120**(1), 56–62 (2016).

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Accepted Manuscript