

ACCEPTED MANUSCRIPT • OPEN ACCESS

Probabilistic clinical target definition with nearest neighbor correlation

To cite this article before publication: Luciano Rivetti *et al* 2025 *Phys. Med. Biol.* in press <https://doi.org/10.1088/1361-6560/ae2aa1>

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 © 2025 The Author(s). Published on behalf of Institute of Physics and Engineering in Medicine by IOP Publishing Ltd.



As the Version of Record of this article is going to be / has been published on a gold open access basis under a CC BY 4.0 licence, this Accepted Manuscript is available for reuse under a CC BY 4.0 licence immediately.

Everyone is permitted to use all or part of the original content in this article, provided that they adhere to all the terms of the licence <https://creativecommons.org/licenses/by/4.0>

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 may be required. All third party content is fully copyright protected and is not published on a gold open access basis under a CC BY licence, unless that is specifically stated in the figure caption in the Version of Record.

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

Probabilistic Clinical Target Definition with Nearest Neighbor Correlation

Rivetti L. ^{*1,2}, Buti G.², Amoudruz L.³, Ajdari A.², Sharp G.², Studen A.^{1,4}, Jeraj R^{1,5,4}, and Bortfeld T.²

¹Faculty of Mathematics and Physics, University of Ljubljana, Ljubljana, Slovenia

²Massachusetts General Hospital and Harvard Medical School, Boston, United States of America

³School of Engineering and Applied Sciences, Harvard University, Cambridge, United States of America

⁴Jožef Stefan Institute, Ljubljana, Slovenia

⁵University of Wisconsin-Madison, Madison, United States of America

Abstract

Objective: The delineation of the clinical target volume (CTV) in radiotherapy is fundamentally uncertain due to the invisibility of microscopic disease on medical images. The ICRU 83 report acknowledges this by proposing a probabilistic interpretation of the CTV, but it does not define how to compute the probability of microscopic tumor presence (MTP) in tissue. This work addresses this gap by introducing a novel stochastic model that estimates the probability of MTP at the voxel level based on local spatial correlations in the voxels' neighborhood.

Approach: We developed two first-principles stochastic models to simulate MTP under different assumptions, incorporating spatial correlation between neighboring voxels. The constant marginal probability (CMP) model assumes spatially uniform MTP and is suited for tumors without radial dependence on the distance from the gross tumor volume (GTV). The variable marginal probability (VMP) model introduces radial dependence, modeling decreasing MTP with distance from the GTV. The CMP model was evaluated on prostate cancer data, while the VMP model was assessed using breast and lung cancer data.

Results: Both models accurately reproduced the fraction of times that MTP is present. In the prostate case, the CMP model estimated a marginal probability of MTP of 0.03, consistent with a literature report that indicates an average total microscopic tumor volume of approximately 583 mm³ across patients. The VMP model successfully replicated the radial distribution of tumor islets, achieving mean absolute errors of 0.01 mm and 0.011 mm for breast and lung cancer distance distributions, respectively. However, not all MTP characteristics could be fully captured by the models, and in some cases discrepancies with population based tumor characteristics remain.

Significance: This work introduces a statistically consistent framework that enables a probabilistic definition of the CTV. The proposed models provide a new way to capture key aspects of microscopic disease spread by introducing local voxel correlations.

*Corresponding author: Luciano.Rivetti@fmf.uni-lj.si

1 34 1 Introduction

5 35 Accurate delineation of the clinical target volume (CTV) is a critical step in radiation therapy, directly
6 36 influencing treatment outcomes [21, 11, 27]. CTV delineation is very challenging, largely due to the big
7 37 uncertainties inherent in identifying microscopic disease spread [29, 30, 20]. Despite advancements in imaging,
8 38 inconsistencies in CTV contouring remain a persistent issue, especially in complex cases where inter-observer
9 39 variability can reach up to one centimeter or more [29]. This uncertainty, often an order of magnitude greater
10 40 than the millimeter-level precision achieved in dose delivery, underscores a fundamental limitation within
11 41 current CTV practices.

12 42 To address the inherent variability of manual CTV delineation, a probability-based definition of the
13 43 CTV driven by clinical data is desirable. The ICRU 83 report introduced the concept of probabilities in
14 44 the CTV definition, stating: “The CTV is a volume of tissue that contains a demonstrable gross tumor
15 45 volume (GTV) and/or subclinical malignant disease with a certain *probability of occurrence* [i.e. presence]
16 46 considered relevant for therapy.” It further specifies that “typically, a *probability of occult disease* higher
17 47 than 5% to 10% is assumed to require treatment” [16]. However, this definition raises key questions: How
18 48 should the “probability of subclinical malignant disease occurrence” or “probability of occult disease” in a
19 49 tissue be statistically defined? On what clinical evidence should these probabilities be based? And how
20 50 can probabilities assigned to tissues be translated into discretized grids used in treatment planning for dose
21 51 calculation and optimization?

22 52 For instance, consider modeling a tissue as a collection of voxels, where we assume that the probability
23 53 of a tumor cell being present in each voxel is known. In this scenario, the probability of finding at least one
24 54 tumor cell in a larger region, composed of multiple voxels, will likely increase with the size of that region.
25 55 However, this is true if and only if the probabilities of finding tumor cells in the voxels do not directly
26 56 influence each other. Consider the opposite scenario where the tumor cells found in the voxels within the
27 57 same region are known to be perfectly influenced by each other—if one voxel contains tumor, then all the
28 58 others in the region do as well, and vice versa. In this case, the probability of finding a tumor cell in the
29 59 entire region remains constant, regardless of its size. This simple example underscores the critical role of the
30 60 assumed level of influence that voxels have on each other in determining whether the larger region composed
31 61 of those voxels should be targeted for treatment, as outlined in the ICRU report. We refer to this problem
32 62 as “voxel correlation”. It emphasizes the need for a deeper understanding of these definitions, the influence
33 63 of voxel correlations, and the development of models that accurately reflect clinical data.

34 64 The first studies that introduce explicit representations of probabilities into radiation target volumes
35 65 were presented by authors that aim to create a fully probabilistic definition of the CTV, the so-called clinical
36 66 target map (CTM). In the CTM, each voxel outside of the GTV is assigned a probability of microscopic
37 67 tumor presence (MTP) which describes the probability of finding at least one microscopic tumor cell [6,
38 68 23]. This MTP probability arises from introducing a binary random variable of MTP to a voxel which takes
39 69 a positive state if the voxel contains tumor and a zero state otherwise. The MTP probabilities used for
40 70 the CTM are subsequently obtained as the expected value of the MTP random variable calculated in each
41 71 voxel. This definition of tumor probability addresses one of the ambiguities in the original ICRU definition;
42 72 the CTM probabilities are *marginal* probabilities, meaning it expresses the probability of MTP in a voxel,
43 73 without considering the state of MTP in any other voxel.

44 74 For example, the MTP model used in Shusharina et al. states that the state of MTP in one voxel is
45 75 independent of the state of MTP in any other voxel. While this model is statistically interpretable and
46 76 computationally simple, it has notable limitations [23]. First, the independence assumption does not fit well
47 77 with microscopic tumor patterns reported in the literature, such as [15, 31], where it is observed that MTP
48 78 often manifests as clusters or “islets” of millimeter size around the GTV. This suggests that the presence of
49 79 tumor in one location conditions the probability of tumor in its neighborhood.

50 80 A different MTP model, referred to as *contiguous circumferential growth model*, was proposed by Buti
51 81 and Bortfeld et al. [7, 5]. Here it was assumed that the state of MTP in a voxel fully determines the
52 82 state of MTP of all voxels at a distance less than or equal to that voxel. This model was inspired by
53 83 the way CTV is defined in clinical guidelines; i.e., the CTV margin includes a percentage of *maximum*
54 84 distances of MTP identified on pre-treated histopathological specimens for a disease site [13, 18]. Buti et al.

1
2
3 85 modeled the radial dependence of MTP using the maximum microscopic tumor distance probability density
4 function (PDF) reported in histopathologic studies such as [17]. While this model has practical appeal, it
5 cannot reproduce aforementioned microscopic tumor islets observed and therefore overestimates the marginal
6 probability of MTP in a sub-volume. Overestimating MTP probabilities could lead to overly conservative
7 radiotherapy treatment plans that potentially over-dose healthy tissue.
8

9 In this study, we present a stochastic model for MTP that quantifies the dependence of MTP in a voxel on
10 its local neighborhood. This dependency is introduced through the use of a statistically well-known concept,
11 the correlation factor, which imposes discrete rules akin to those in the design of a cellular automaton
12 [4]. It will be demonstrated that the correlation parameter regulates the number and size of microscopic
13 tumor islets appearing outside of the GTV. The model parameters are calibrated using MTP characteristics
14 extracted from histo-pathological findings reported for a population of patients of three specific cancer types:
15 prostate, breast, and lung. The model is evaluated based on its ability to replicate observed microscopic
16 tumor features, thereby providing insights into its potential applicability for clinical use in radiotherapy
17 planning.
18

19 99 2 Material and Methods

20
21 100 In this section, we first introduce the terminology and general statistical framework used to model tumor
22 probabilities throughout this study. Second, we show that existing statistical models, i.e. the *independent*
23 model first proposed by Shusharina et al. and the contiguous circumferential growth model proposed by
24 Buti and Bortfeld et al. [7, 5], fail to reproduce the observed clinical data. Third, we propose two new
25 models based on first principles that allow the generation of more realistic patterns. Fourth, we describe the
26 literature data we used to fit our models. Fifth, we describe the simulations we made and the metrics we
27 evaluated.
28

29 107 2.1 Terminology

30
31 108 We follow the terminology introduced by the ICRU by referring to any malignant disease outside of the
32 GTV as “microscopic tumor presence” (MTP). We apply this term irrespective of the fact that the disease
33 extended from the GTV or originated at a specific location. MTP implies that at least one tumor cell is
34 present at that specific location, whereas no MTP means that that location is tumor-free.
35

36 112 2.2 General statistical framework

37
38 113 Consider a 3D grid with N voxels. Each voxel $i \in \{1, \dots, N\}$ is associated with a binary random variable
39 C_i , where $C_i = 1$ indicates the presence of MTP, and $C_i = 0$ indicates its absence. We use lowercase c_i to
40 denote a specific realization (observed value) of the random variable C_i . The *marginal probability* of MTP
41 in voxel i is denoted p_i , and corresponds to the expected value of the random variable: $\mathbb{E}[C_i] = p_i$.
42

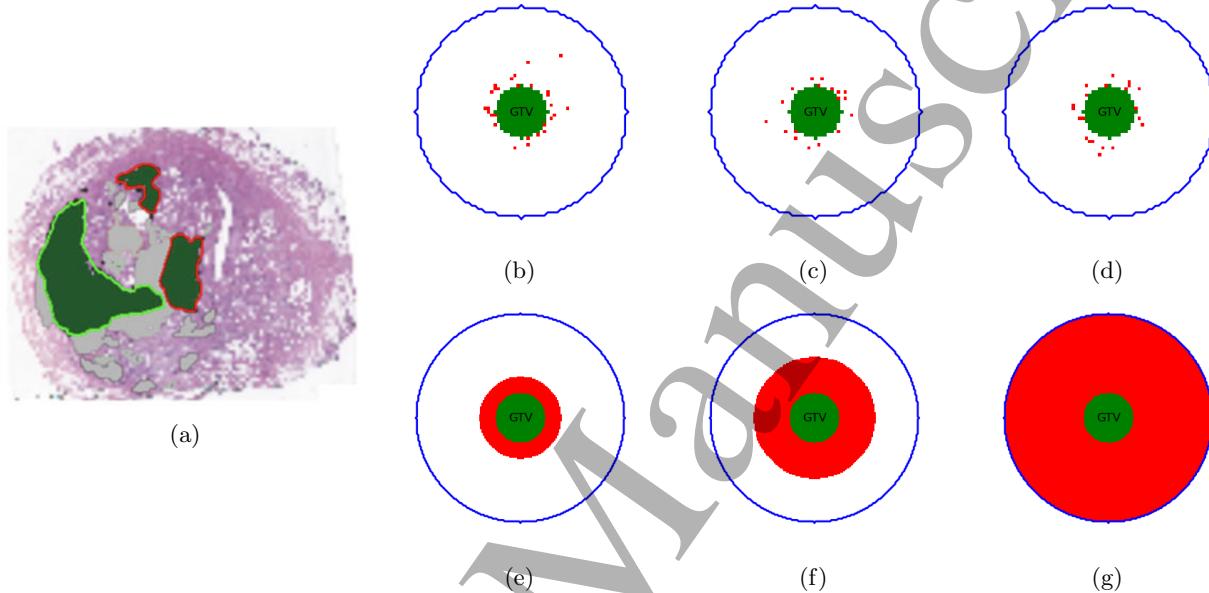
43 For the sake of simplicity, we will refer to the state or probability of the MTP random variable in voxel
44 i simply as the state or probability of voxel i . We assume that the probability of all voxels in the GTV is
45 equal to 1, $p_i = 1 \forall i \in \text{GTV}$. Once the states of all voxels in the grid have been assigned, we can visualize
46 the combined GTV and MTP regions as a single realization, referred to as the *tumor configuration*. The
47 probability of observing a specific tumor configuration is given by the corresponding *joint probability* over
48 the entire grid. This joint probability quantifies the likelihood of the system as a whole adopting a particular
49 configuration, rather than considering voxel states independently.
50

51 124 2.3 Limitations of the current models

52
53 125 The independent model of Shusharina et al. assumes that there is no correlation between the states of any
54 two voxels C_i and C_j , $i \neq j$ [23]. Therefore, tumor configurations can be generated from this model by
55 randomly sampling the state of each voxel i of the grid independently with probability p_i . Simulations in a
56

1
2
3 circular geometry, assuming a constant marginal probability throughout the prostate, show that a “poppy
4 seed” pattern of MTP is generated (Fig. 1, b-d). This pattern differs from the characteristic islet formations
5 of MTP observed in histopathologic images (Fig. 1a) [15, 31].
6
7

8 The contiguous circumferential growth model from Buti and Bortfeld et al. imposes a direct correlation
9 among the voxels. Specifically, if the state of a voxel at distance d from the GTV is randomly assigned as
10 MTP, then all voxels within a distance smaller or equal to d are also classified as MTP. Figs. 1, e-g show
11 that tumor configurations produce an “onion” pattern with the MTP appearing as being thick, isotropic
12 expansions of the GTV surface, thereby failing to capture the islet formations.
13
14



32 Figure 1: Histopathology and simulation of MTP. The green contours (or circle) indicates the GTV, the
33 red contours/structures represent MRI-invisible tumor lesions (i.e., MTP in our framework), and the blue
34 contours depict an anatomical barrier impermeable to tumor cells, such as the prostate wall. (a) Real
35 histopathology slice from a prostate cancer patient, adapted from van Houdt et al. [15] with their permission.
36 (b-d) 2D simulations generated using the independent model proposed by Shusharina et al. [23], assuming
37 a constant marginal probability throughout the prostate (2 shell model), with a circular GTV. Each dot
38 represents a voxel. (e-g) 2D simulations performed using the contiguous circumferential growth model
39 proposed by Buti et al. [7], also with a circular GTV.
40

41 It is evident that these two models represent extreme cases with rigid assumptions: either complete voxel
42 independence or full correlation. We propose that a more realistic model can be achieved by introducing
43 a parameter to model the degree of voxel correlation, bridging the gap between these extremes and better
44 reflecting observed MTP patterns.
45

46 2.4 Nearest neighbor correlation with constant marginal probability of MTP 47

48 Our approach differs from the existing methods by introducing pairwise correlations between voxels and
49 by keeping the marginal probability in every voxel within a specified volume the same, here denoted as p .
50 Similar to the design of a cellular automaton, fixed rules are defined that determine the state of each voxel
51 based on the voxels in its neighborhood. The algorithm begins by selecting a random voxel $i = 1$ from
52 the grid and sampling its state $c_1 \sim \text{Bernoulli}(p)$. Then, the states of its nearest neighbor (NN) voxels
53 are sequentially sampled according to a predefined visiting order. The state c_i of a newly visited voxel i is
54 determined in one of two ways:
55
56
57
58
59
60

- 1
2
3 148 1. with probability q , its state is copied from a randomly selected previously visited nearest neighbor
4 149 $k \in \text{vNN}_i$, i.e., $c_i = c_k$, where $\text{vNN}_i = \{k \mid k \text{ was visited before } i\}$;
5
6 150 2. with probability $1 - q$, a new value is sampled independently: $c_i \sim \text{Bernoulli}(p)$.

7 This process ensures that the state of each newly visited voxel i is partially inherited from its visited
8 nearest neighbors.
9

10 The visiting algorithm follows a structured traversal pattern to systematically explore the voxel grid.
11 The process starts with a randomly selected voxel, after which its NN voxels are visited in a cross-like
12 pattern, moving sequentially up, left, down, and right (Fig. 2b–e). Once all first-order NNs of the initial
13 voxel have been visited, the algorithm continues with the NNs of the next visited voxel, ensuring that no
14 voxel is revisited (Fig. 2f–h). This structured approach is iteratively applied until all voxels in the grid have
15 been assigned a state.
16

17 Following these rules, the *conditional probability* of C_i given its visited nearest neighbors can be written
18 as:
19

$$\mathbb{P}(C_i = 0 \mid \{C_k = c_k\}_{k \in \text{vNN}_i}) = q \left(1 - \frac{1}{n} \sum_{k \in \text{vNN}_i} c_k \right) + (1 - q)(1 - p), \quad n = |\text{vNN}_i|, \quad (1)$$

20 Using the conditional and marginal probabilities defined above, we can calculate the joint probability of
21 a given state using Bayes' theorem:
22

$$\mathbb{P}(C_i, \{C_k\}_{k \in \text{vNN}_i}) = \mathbb{P}(C_i \mid \{C_k\}_{k \in \text{vNN}_i}) \mathbb{P}(\{C_k\}_{k \in \text{vNN}_i}). \quad (2)$$

23 Thus, the probability of the configuration b) in figure 2 is equal to $p(q + (1 - q)p)$. With this algorithm,
24 if we assign $\mathbb{E}[C_0] = p$ to the first voxel, it is easy to show that
25

$$\mathbb{E}[C_i] = \mathbb{E}[C_k] = p \quad \forall i, k, \quad (3)$$

26 as intended. This result holds for any value of the correlation coefficient q , highlighting that p and q are
27 independent parameters of the model..
28

29 The statistical interpretation of q can be clarified by evaluating the Pearson correlation coefficient ρ
30 between a random variable C_i and its nearest neighbor C_k on a 1D grid. The correlation coefficient ρ is
31 given by:
32

$$\rho_{C_i, C_k} = \frac{\text{COV}_{C_i, C_k}}{\sigma_{C_i} \sigma_{C_k}} = \frac{\mathbb{E}[C_i C_k] - \mathbb{E}[C_i] \mathbb{E}[C_k]}{\sigma_{C_i} \sigma_{C_k}}. \quad (4)$$

33 Since the expected values of C_i and C_k are equal to p , and we are dealing with binary random variables,
34 the standard deviations are $\sigma_{C_i} = \sigma_{C_k} = \sqrt{p - p^2}$. Because in a 1D grid there is only one NN, we have
35 $\mathbb{E}[C_i C_k] = qp + (1 - q)p^2$. This yields
36

$$\rho = \frac{qp + (1 - q)p^2 - p^2}{p - p^2} = q. \quad (5)$$

37 Therefore, in the case where the marginal probability is the same in the whole grid, the design parameter
38 q corresponds to the Pearson correlation coefficient ρ between two voxels in a 1D grid. However, in higher-
39 dimensional grids (2D or 3D), the Pearson correlation coefficient is given by $\rho = \sum_{n=1}^{\infty} a_n q^{2n-1}$ where the
40 coefficients a_n satisfy $\sum_{n=1}^{\infty} a_n = 1$ (Appendix A). The values of n and a_n depend on the dimension and size
41 of the grid.

42 An important question is the probability of finding no tumor cells in a group of voxels. We can calculate
43 the probability of having no tumor cells in an entire grid with N voxels using equation 1 and 2 as:
44

$$\mathbb{P}(C_0 = 0, \dots, C_N = 0) = P_0(p, q, N) = (1 - p) \left(q + (1 - q)(1 - p) \right)^{N-1}. \quad (6)$$

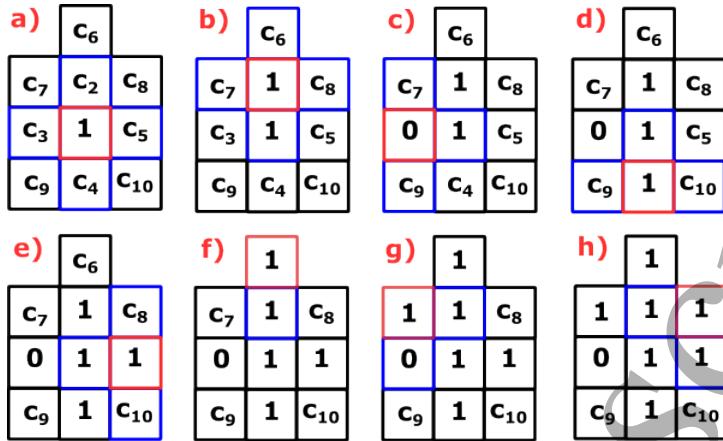


Figure 2: Illustration of the sampling algorithm in a 2D grid. The red squares indicate the last visited voxel at each step and the blue squares shows their NN. (a) The algorithm begins by randomly selecting a voxel in the grid (C_1) and sampling its state with probability p . (b-e) Then, all the NN of the first sampled voxel are visited incorporating correlation effects for each neighboring voxel. (f-h) Once all NN of C_1 are visited, the algorithm proceeds by sampling voxels adjacent to the second visited voxel (C_2) without revisiting voxels.

However, this result is specific to the entire grid and cannot be directly applied to subsets of voxels. In subsets, it is not guaranteed that every voxel has a nearest neighbor within the subset, which could affect the calculation.

Additionally, P_0 can be interpreted as the fraction of patients who do not exhibit MTP. To correctly reproduce P_0 , the parameters p and q must satisfy equation 6. Since $q \in [0, 1]$, this relationship constrains the admissible range of the marginal probability p : its maximum value, $1 - P_0$, occurs when $q = 1$, while its minimum value, $1 - P_0^{1/N}$, occurs when $q = 0$.

Interpreting how p varies with changes in P_0 , q , or N is challenging based on equation 6, especially since N appears as an exponent in the equation. To provide a clearer understanding of equation 6, we approximated this equation using Taylor series, assuming a large number of voxels in the grid (Appendix B). By discarding the quadratic and higher-order terms, we simplify equation 6 to:

$$\frac{P_0}{1-p} \approx e^{-N(1-q)p} = e^{-\frac{p}{T(q,N)}} \quad (7)$$

where $T(q, N) = \frac{1}{N(1-q)}$. Equation 7 shows that the marginal probability of the voxels p only depends on P_0 and $T(q, N)$. Fig 3 shows the relationship between p and $T(q, N)$ for different values of P_0 and using a fixed $N = 75735$. This value corresponds to the number of voxels in the sampling space of the prostate patient that we will describe in the evaluation subsection. We observe that $T(q, N)$ serves as a "transition factor" between two tumor configuration states, fully correlated and uncorrelated. As we will later show, when $T \rightarrow \infty$, the tumor configuration exhibits clustering, meaning the voxel states become more correlated. In contrast, when $T = \frac{1}{N} \approx 0$, the tumor configuration generates isolated islets, where voxel states are essentially uncorrelated. From figure 3, P_0 can be interpreted as influencing the slope of the transition between different tumor configurations. The definition of $T(q, N)$ ensures that our model remains resolution-independent: for a given N , we can adjust q to achieve a desired value of T , thereby making the model adaptable to different grid resolutions.

Not only p is bounded in this model but also the number of unconnected components or tumor islets. If we consider the case when $T \approx 0$, the voxels are uncorrelated and the number of unconnected components islets are maximal. Then, because in this regime $p \approx 1 - P_0^{1/N} \approx 0$, we can assume that the probability of having two tumor voxels connected is low and the maximum average number of tumor islets can be calculated

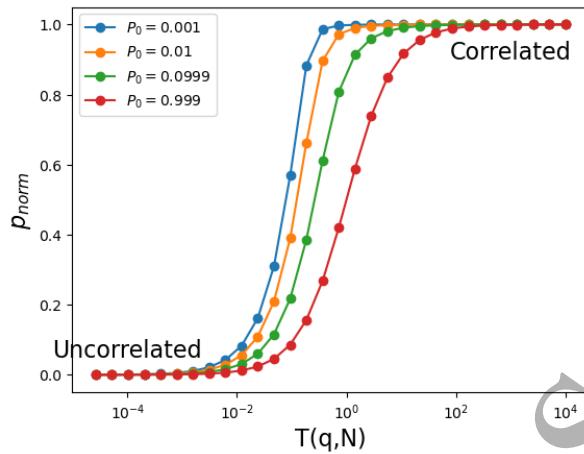


Figure 3: Normalized marginal probability versus transition factor $T(q, N)$ calculated for $N = 75735$ at different P_0 . The probability p was normalized such that $p_{\text{norm}} = \frac{p - p_{\min}}{p_{\max} - p_{\min}}$, where $p_{\min} = 1 - P_0^{1/N}$ and $p_{\max} = 1 - P_0$. The plot illustrates how the system transitions from the correlated state to the uncorrelated state, with different colors representing transitions for various values of P_0 . The minimum marginal probability is achieved when the system is uncorrelated ($T \rightarrow 0$ or $q = 0$), while the maximum occurs when the system is fully correlated ($T \rightarrow \infty$ or $q = 1$). The x -axis is displayed in logarithmic scale.

as Np . Using equation 7 we can write:

$$Np \approx \ln\left(\frac{1}{P_0}\right) \quad (8)$$

Conversely, in the opposite limit when $T \rightarrow \infty$, the voxels become fully correlated. In this regime the system can only generate two possible outcomes: either the entire admissible volume is tumor, or no tumor is present at all. Since the probability of no tumor is P_0 , the complementary probability of having tumor is $1 - P_0$. In this case, at most a single connected islet can be formed with probability $1 - P_0$, and therefore

$$\mathbb{E}[\#\text{islets}] = \bar{N}_I = 1 \cdot (1 - P_0) = 1 - P_0.$$

This expression defines the lower bound for the expected number of islets.

2.5 Nearest neighbor correlation with variable marginal probability of MTP

Several studies have shown that in certain tumors, the MTP decreases as the distance from the GTV border increases [25, 12, 8, 22, 28]. Our previous model cannot capture this behavior, as it assumes constant marginal probabilities across voxels. To address this limitation and incorporate radial dependence, we introduce variable marginal probabilities.

The visiting voxel strategy in this model is the same as the one described in the previous section. However, the main differences between the models is how we calculate the posterior probability of a new voxel given its visited nearest neighbor states. Because in this model the two NN voxels can have different marginal probabilities we should contemplate this case in the equations. To clarify this point, lets first consider a 1D grid (Fig 4a).

Inspired by the model proposed in our previous work [5, 6, 7], we adopted the *no-tunneling assumption*, which requires that in the fully correlated case, tumor configurations form continuous regions—i.e., isolated tumor islets do not appear.

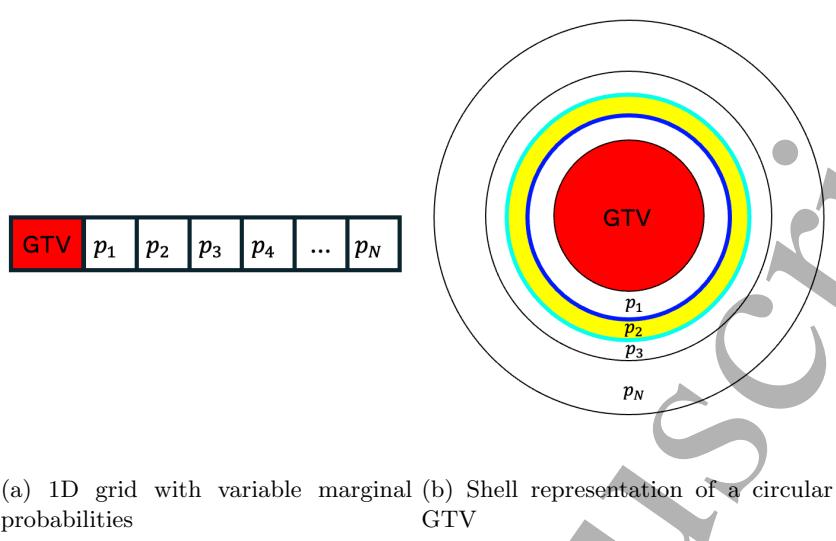


Figure 4: Variable marginal probability model representation. a) Shows that the marginal probabilities in this model are higher closer to the GTV. b) Shows the shell representation used in Shusharina et al. for 2D or 3D simulations. The second shell is outlined in blue and the third shell in cyan; the yellow region between them corresponds to the second layer. This model assumes that all the voxels in the same layer have the same marginal probability.

Let C_i and C_k be two neighboring voxels such that $p_i \geq p_k$. Under the non-tunneling assumption, the following conditional probabilities must hold:

$$\begin{aligned} \mathbb{P}(C_k = 1 | C_i = 0) &= (1 - q)p_k \\ \mathbb{P}(C_i = 1 | C_k = 1) &= q + (1 - q)p_i. \end{aligned} \quad (9)$$

These expressions enforce spatial continuity when $q = 1$. Specifically, if voxel i is in the tumor state, then all voxels from i up to the GTV boundary must also be in the tumor state. This is ensured when: $\mathbb{P}(C_i = 1 | C_k = 1) = 1$. Conversely, if voxel i is non-tumor, all subsequent voxels up to the outer boundary N must also be non-tumor. This is captured by the condition: $\mathbb{P}(C_k = 1 | C_i = 0) = 0$.

To ensure consistency with the prescribed marginal probabilities $\mathbb{E}[C_i] = p_i$ and $\mathbb{E}[C_k] = p_k$, two additional conditions must be satisfied:

$$\begin{aligned} \mathbb{P}(C_k = 1 | C_i = 1) &= \frac{p_k}{p_i}q + (1 - q)p_k \\ \mathbb{P}(C_i = 1 | C_k = 0) &= \frac{p_i - p_k}{1 - p_k}q + (1 - q)p_i. \end{aligned} \quad (10)$$

These expressions ensure that the model reproduces the correct marginal expectations. For the mathematical derivation look at Appendix C. Note that, when $p_i = p_k$, the previous equations simplifies to Equation (1), recovering the constant marginal probability model described in the previous section.

To extend this model to 2D and 3D, we adopted the shell representation proposed by Shusharina et al. [23]. In this framework, the CTM is modeled using a finite set of iso-probability surfaces, called *shells*, each associated with a probability value (Fig 4b). This probability represents the likelihood that the MTP extends beyond the corresponding shell within a given patient population. The innermost shell coincides with the GTV contour, while the outermost shell delineates the region beyond which tumor presence is considered negligible. The region between two adjacent shells is defined as a *layer*. This model also assumes that all voxels within the same layer share the same marginal probability.

The state of a new voxel C_i can be computed by considering the states of its visited nearest neighbors vNN_i along with their marginal probabilities. Assume voxel i has n visited neighbors, where the first l neighbors satisfy $p_k \geq p_i$ for all $k \in \{1, 2, \dots, l\}$, and the remaining $n - l$ neighbors satisfy $p_k < p_i$ for all $k \in \{l + 1, \dots, n\}$. Then, using equations 9 and 10, we can express the conditional probability of voxel i being in state 1, given the neighboring realizations $\{c_k\}_{k \in \text{vNN}_i}$, as:

$$\mathbb{P}(C_i = 1 | \{C_k = c_k\}_{k \in \text{vNN}_i}) = \frac{q}{n} \left(\sum_{k=1}^l \frac{p_i}{p_k} c_k + \sum_{k=l+1}^n \left[c_k + (1 - c_k) \frac{p_i - p_k}{1 - p_k} \right] \right) + (1 - q)p_i, \quad n = |\text{vNN}_i|. \quad (11)$$

In this model, we define N' marginal probabilities, denoted as $(p_1, p_2, \dots, p_{N'})$, where N' represents the number of distinct voxel classes given by the number of layers. To fully determine these $N' + 1$ parameters, we require a corresponding set of $N' + 1$ independent equations.

Empirical studies have reported the probability of finding MTP beyond a given distance r from the nearest point on the GTV, denoted here as $w(r)$ (see Section 2.6 for more details) [25, 22, 28]. We used these empirical constraints to estimate the marginal probabilities as follows:

$$1 - \mu_i(q, p_1, \dots, p_{N'}) = w(r_i), \quad \forall i \in \{1, \dots, N\}, \quad (12)$$

where r_i is the distance to the i -th shell and $\mu_i(q, p_1, \dots, p_{N'})$ is the probability of finding no tumor in all the voxels outside the i -th shell. Note that equation 12 must hold for any selected value of q in the model.

However, obtaining an analytical expression for $\mu_i(q, p_1, \dots, p_{N'})$ is generally intractable. Therefore, we estimate this value using Monte Carlo simulations. Specifically, we approximate $\mu_i(q, p_1, \dots, p_{N'})$ as the fraction of realizations that do not contain tumor voxels beyond the i -th shell. Given a fixed value of q , we can determine the optimal values $\{\hat{p}_1, \dots, \hat{p}_{N'}\}$ by solving the following optimization problem:

$$\{\hat{p}_1, \dots, \hat{p}_{N'}\} = \underset{\{p_1, \dots, p_{N'}\}}{\operatorname{argmin}} \sum_{i=1}^N (1 - \mu_i(q, p_1, \dots, p_{N'}) - w(r_i))^2. \quad (13)$$

2.6 Data

Population-based clinical data from the literature were used to fit the parameters of the proposed models for three different cancer types: prostate, breast, and lung. These datasets provide information on microscopic disease patterns, enabling a quantitative evaluation of the model's flexibility in capturing real-world tumor behavior. For each cancer type, relevant studies reporting histopathological findings were selected to extract parameters such as the probability of MTP presence, average islet volume, and MTP radial dependence. The data are summarized in Table 1. Definitions of all model parameters and metrics used in this study are provided in Table 2.

Prostate Data

Clinical data for prostate cancer were obtained from the article written by Bajgira et al., encompassing 518 patients who underwent multiparametric MRI followed by histological examination post-resection [19]. The study found that 76% of patients had multifocal microscopic islets detectable only through histopathology (not in MRI), corresponding to $P_0 = 0.24$. The average number of islets \bar{N}_I considering the 518 patients was 0.87 per patient. The reported average islet radius was 7.82 ± 5.74 mm. Assuming a truncated Gaussian distribution for the radius and spherical islet shape, the estimated average islet volume (\bar{V}_I) was 513 mm^3 , with the 5th and 95th percentiles at 1 mm^3 and 2024 mm^3 , respectively. The average total microscopic tumor volume (\bar{V}_T) was estimated by multiplying \bar{V}_I by the average number of islets considering the patients which presented MTP. This resulted in an \bar{V}_T equal to 583 mm^3 , with the 5th and 95th percentiles at 1 mm^3 and 2300 mm^3 . A detailed description of how these data and calculations were obtained is provided in Appendix D.

1
2
3 **280 Breast Data**

4
5 Breast cancer data were sourced from the publication by Stroom et al., where 60 patients underwent contrast-
6 enhanced MRI followed by histopathological examination [25]. The study reported that 72% of patients had
7 multifocal islets detectable only through histopathology ($P_0 = 0.28$). Among those patients with islets, the
8 reported average number of islets was 4.2. When averaged over the entire cohort of 60 patients, including
9 those without islets, this corresponds to an average number of islets of approximately 3 per patient. The
10 mean islet radius was 2.1 mm, and assuming spherical islets, \bar{V}_I can be approximated to 39 mm^3 . Similarly
11 \bar{V}_T can be approximated to 162 mm^3 . Additionally, the authors analyzed the spatial distribution of islets
12 relative to the GTV, reporting that it follows a normal distribution with a mean of zero and a variance of 7
13 mm^2 .

14 Additionally, a related study [22] provided probability data for detecting MTP beyond a given distance
15 r from the closest point on the GTV, denoted as $w(r)$. Note that $w(0) = 1 - P_0$ and in this case this value
16 was 0.72 (Fig 5). Since the study does not report the distribution of MTP in relation to polar coordinates,
17 we assume an isotropic MTP distribution.

18
19 **294 Lung Data**

20
21 The lung cancer data were obtained from the publication written by Stroom et al., comprising 34 patients
22 who underwent CT imaging followed by histopathological analysis post-resection [25]. Multifocal islets,
23 detectable only through histopathology, were identified in 50% of cases. The reported \bar{N}_I calculated over the
24 entire cohort of 34 patients was 2.47. The mean islet radius was reported as 1.6 mm, and assuming spherical
25 islets, the estimated \bar{V}_I was 18 mm^3 and \bar{V}_T was estimated to be 85 mm^3 . Additionally, the authors analyzed
26 the spatial distribution of islets relative to the GTV, reporting that it follows a normal distribution with a
27 mean of zero and a variance of 10 mm^2 .

28 A related study [28] provided $w(r)$, with $w(0) = 1 - P_0 = 0.5$ (Fig 5). Since the MTP distribution data
29 in polar coordinates were not reported, an isotropic distribution was assumed.

Cancer Type	P_0	\bar{N}_I	$\bar{V}_I (\text{mm}^3)$	$\bar{V}_T (\text{mm}^3)$	$w(r)$
Prostate	0.24	0.87	513	583	No
Breast	0.28	3	39	162	Yes
Lung	0.50	2.47	18	85	Yes

30
31 Table 1: Summary of microscopic tumor parameters for prostate, breast, and lung cancers. P_0 : proportion
32 of patients with MTP; \bar{N}_I : average number of islets per patient; \bar{V}_I : average islet volume per patient; \bar{V}_T :
33 average total microscopic tumor volume among patients with microscopic tumor presence; $w(r)$: indicates
34 whether data on the cumulative probability distribution of detecting MTP outside a given distance from the
35 GTV are available in the dataset (Yes/No).

Parameter/Metric	Description
q	Correlation factor representing the probability of copying a neighboring state
p	Marginal probability of a voxel
T	Transition factor controlling the change between MTP states
P_0	Proportion of patients without MTP
$w(r)$	Cumulative probability of finding MTP beyond distance r
\bar{N}_I	Mean number of MTP islets per patient
\bar{V}_I	Mean volume of MTP islets per patient (mm^3)
\bar{V}_T	Mean total MTP volume per patient (mm^3)

53 Table 2: Summary of model parameters and metrics used in this study.

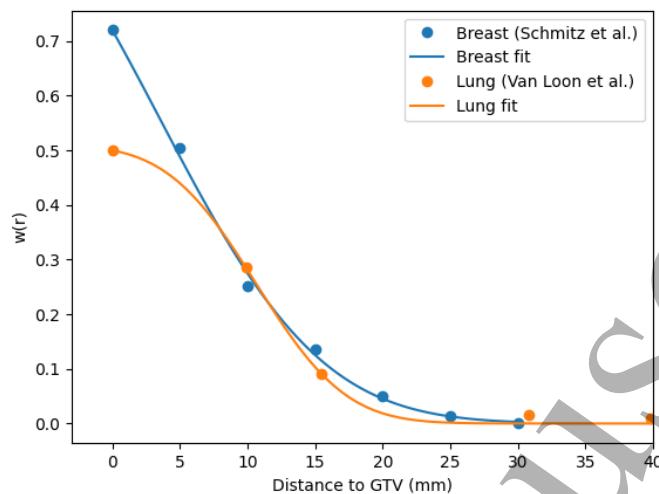


Figure 5: Probability of finding MTP beyond a given distance r from the nearest point on the GTV denoted as $w(r)$. The blue points represent the data extracted from Schmitz et al.[22] for a breast cancer population and the line shows the best fit. The orange points represent the data extracted from Van Loon et al.[28] for a lung cancer population and the line shows the best fit.

2.7 Simulation design

Several simulations were ran on patient-specific anatomical data to study the properties of the models developed in this work and to find the optimal parameters that fit the literature data.

Prostate Simulations

The prostate simulation was conducted using the constant marginal probability model described in Section 2.4. This model was chosen because, outside the GTV, it assumes that the marginal probability of MTP is uniform throughout the prostate tissue, meaning that it does not vary with distance from the visible lesion. This assumption was adopted because to the best of our knowledge there is no evidence in the literature suggesting a spatial correlation between invisible lesions and the GTV. Additionally, in clinical practice, the entire prostate is typically included within the CTV and irradiated uniformly outside the visible lesion. This widespread treatment approach further supports the notion that there is limited information on the spatial correlation between the GTV and microscopic tumor presence.

The simulations were performed on the anatomy of the third patient from the publicly available QIN PROSTATE dataset, which includes physician-delineated contours of the GTV and prostate [10]. The simulation space was resampled to a 1 mm resolution and defined as the volume outside the GTV but within the prostate (Fig. 8a).

The study was conducted by generating 5000 simulations for each value of $T(q, N)$, which was varied from 1×10^{-4} to 1×10^3 , while grid resolutions were adjusted between 0.5 mm and 2 mm. The parameter P_0 was fixed to match the proportion of patients without MTP reported on Bajgira et al [19]. From these simulations, \bar{N}_I , \bar{V}_I and \bar{V}_T were computed. Each islet was defined as the independent components of the simulated MTP. Additionally, the distribution of islet volumes for different T values was analyzed, and the 3D MTP marginal probability, along with different MTP realizations, were visually examined.

1
2
3 **326 Breast and Lung Simulations**
4
5

6 The breast and lung simulations were performed using the variable marginal probability model described in
7 Section 2.5. The breast simulations were based on the anatomy of the fifth patient from the Advanced-MRI-
8 Breast-Lesions dataset, which contains physician-delineated contours of the GTV [9]. The lung simulations
9 were conducted using the anatomy of the first patient from the NSCLC-Radiomics dataset, which includes
10 physician-delineated contours of the GTV and lungs [1].

11 For both cancer types, the simulation space was resampled to a 2 mm resolution grid and defined as the
12 volume outside the GTV, extending up to the distance where the probability of detecting MTP falls below
13 1%, while accounting for anatomical barriers (Figs. 11a and 12a).

14 The sampling space for both cancer types was discretized into 12 layers that were equally spaced con-
15 sidering $w(r)$. The distance d_i to i -th shell is equal to $d_i = w^{-1}(w_1(13 - i)/12))$, where $w_1 = 1 - P_0$ is the
16 probability of finding MTP outside the GTV. The parameter space was explored for different values of q
17 ranging from 0 to 1. The parameters p_1, \dots, p_{12} were determined such that $1 - \mu_i(q, p_1, \dots, p_{12}) = w_i$ for each
18 value of q . This was achieved by solving the optimization problem defined in Equation 13 using Bayesian
19 optimization as implemented in the `scikit-optimize` library [14]. At each step of the optimization, 10,000
20 simulations were conducted to approximate μ_i .

21 Using the optimal model parameters, we evaluated \bar{N}_I , \bar{V}_I , and \bar{V}_T . Additionally, we analyzed the
22 radial distribution of the simulated islets relative to the GTV boundary. A visual assessment of the CTM,
23 the isoprobability shell 0.05 derived from the simulations, and representative MTP realizations was also
24 performed.

25 **346 2.8 Implementation**
26

27 The constant and variable marginal probability models were implemented in Python. Simulations of the MTP
28 were accelerated using parallel computing techniques to handle the high computational load. In particular,
29 we used CUDA library to distribute independent simulations across available GPU, significantly reducing
30 runtime *.
31

32 **351 3 Results**
33

34 In this section, we investigate the behavior of the model under different parameter settings. We first explore
35 how the model responds to changes in key parameters, analyzing its statistical properties and consistency
36 with theoretical expectations. We then assess the model's ability to fit empirical MTP data reported in the
37 literature and evaluate its capacity to reproduce observed patterns.
38

39 **356 3.1 Prostate Simulations**
40

41 The simulations generated in the prostate patient showed that \bar{N}_I transition with $T(q, N)$. Notably, this
42 characteristic does not depend on the grid resolution selected for the simulation (Fig 6).

43 The maximum \bar{N}_I obtained from the simulations occurs when voxel states are nearly uncorrelated ($T \approx 0$),
44 aligning with the theoretical value of $\ln(1/P_0) = 1.43$. Conversely, the minimum \bar{N}_I is observed when voxel
45 states are almost fully correlated ($T \rightarrow \infty$), yielding $1 - P_0 = 0.76$. Notably, $T(q, N) = 0.2$ appears to
46 define the midpoint of this transition, while the \bar{N}_I reported in the literature closely matches our model's
47 predictions at $T(q, N) = 0.6$.

48 The distribution of islet volumes for three different values of T exhibited distinct patterns (Fig 7c). When
49 T is in the region where voxel states are nearly uncorrelated ($T \leq 10^{-4}$), the generated islets remain small
50 ($\sim 100 \text{ mm}^3$) as the model fails to cluster many voxels together. As T approaches the middle of the transition
51 ($T \approx 10^{-1}$), voxel clustering increases, leading to larger islets. In this regime, a fraction of the islets can
52 cluster to fill the entire sampling space, corresponding to the volume of the prostate. When T approximates

53 * The source code related to this project is available at <https://github.com/RivettiLuciano/MicroTumorModel>

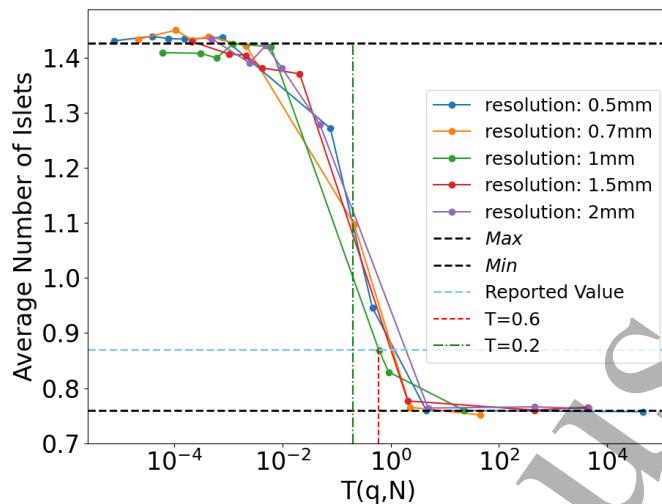


Figure 6: Average number of islets per patient (\bar{N}_I) for different grid resolutions and values of $T(q, N)$. The scatter plot, with different colors, represents results for various grid resolutions. The black dashed line at the top corresponds to the maximum theoretical \bar{N}_I given by $\ln(\frac{1}{P_0})$, while the black dashed line at the bottom represents the minimum theoretical value $1 - P_0$. The sky blue line indicates the \bar{N}_I reported by Bajgira et al. [19]. The red dashed line marks the value of T that reproduces the reported \bar{N}_I at a resolution of 1 mm. The green dot-dashed line represents the value of T that defines the midpoint of the transition between the two distinct tumor configuration states. The x-axis is displayed on a logarithmic scale.

the regime where voxel states are fully correlated ($T \geq 10^1$), the islets predominantly cluster, filling the entire sampling space.

\bar{V}_T exhibited a direct relationship with $T(q, N)$ and remained independent of grid resolution (Fig. 7d). When $T(q, N) \in [10^{-4}, 10^{-1}]$, the model produced \bar{V}_T values within the expected range reported by Bajgira et al. [19]. Notably, at $T(q, N) = 0.02$, the simulations yielded an \bar{V}_T of 583 mm^3 , matching the literature estimate. Similar trends were observed for \bar{V}_I (Appendix E).

Since \bar{V}_T incorporates information from both N_I and \bar{V}_I , we examined tumor configurations at $T(q, N) = 0.02$. At this value, the marginal probability of MTP was $p = 0.028$ (Fig. 8a), and simulations indicated that MTP was absent in 23.8% of cases. Although the marginal probability was low, the simulations still produced cases with significant MTP volumes. Figure 8 illustrates a range of representative tumor configurations:

- (b) A realization with a single islet of 429 mm^3 , approximately close to the literature-reported mean V_T .
- (c) A realization with a small islet of 1 mm^3 , near the 5th percentile of the reported V_T distribution.
- (d) A realization with an islet of 1849 mm^3 , near the 95th percentile of the reported V_T distribution.
- (e) An outlier case in which the tumor has infiltrated the entire prostate.

Furthermore, the simulated tumor islets did not exhibit any preferred orientation relative to the GTV boundaries (Fig. 8b–e).

3.2 Breast Simulations

For different values of q , the optimization process determined the marginal probabilities of the layers, $\{\hat{p}_1, \dots, \hat{p}_{12}\}$, that generated the $1 - \mu_i$ that best matched the w_i values reported in the literature (Fig. 9a).

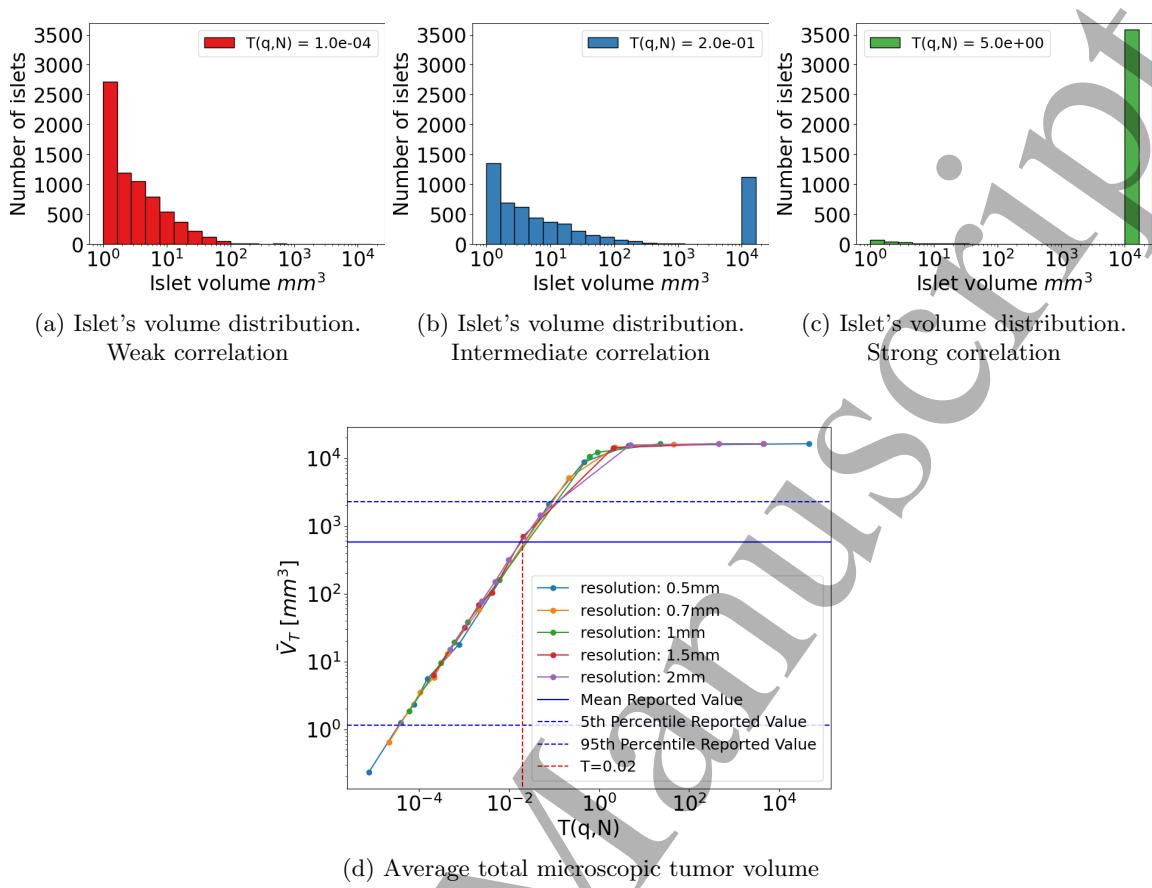


Figure 7: Dependence of MTP volume on $T(q, N)$. (a–c) Distributions of islet volumes for three representative T values at 1 mm resolution: (a) $T = 10^{-4}$ (weak correlation, red), (b) $T = 0.2$ (moderate correlation, blue), and (c) $T = 5$ (strong correlation, red). (d) Simulated \bar{V}_T across varying grid resolutions and T values. The solid blue line denotes the \bar{V}_T reported by Bajgira *et al.*, while the dashed blue lines indicate the 5th and 95th percentiles of the average total microscopic tumor volume per patient reported in [19]. All axes are shown on a logarithmic scale.

The results show that the marginal probabilities satisfy $\hat{p}_1 \geq \hat{p}_2 \geq \dots \geq \hat{p}_{12}$ for all values of q except when q is approaching zero. Similar to the prostate case, the marginal probability decreases when $q \rightarrow 0$ (Fig 9b).

The simulations performed on the breast cancer patient showed that our model can not generate the \bar{N}_I reported in Stroom *et al.* [25] for any value of q . The maximum \bar{N}_I in our simulations was 1.3 and was achieved at $q = 0.999$. On the other hand, the estimated \bar{V}_T obtained from Stroom *et al.* can be generated with the simulations when $q = 0.975$ (Appendix F) [25]. The relation between \bar{V}_T and q follows a similar trend to the one observed in the prostate case, the microscopic tumor volume increases when the voxel states are more correlated. However, the microscopic tumor volume generated were one order of magnitude smaller when compared to the prostate case. Similar results were also observed for \bar{V}_I .

Since \bar{V}_T encapsulates information about both \bar{N}_I and \bar{V}_I , we analyzed the properties of the tumor configurations obtained from simulations with $q = 0.997$. These simulations revealed that the mean distance between the geometric center of the islets and the nearest GTV boundary followed a distribution similar to gaussian distribution with mean 0 mm and variance 7 mm^2 reported in Stroom *et al.* (Fig. 10a) [25]. The R^2 between our simulations and the function reported in the paper is 0.84 and the mean absolute error (MAE) is 0.01.

The CTM demonstrated a radial dependence on distance from the GTV border, with probabilities de-

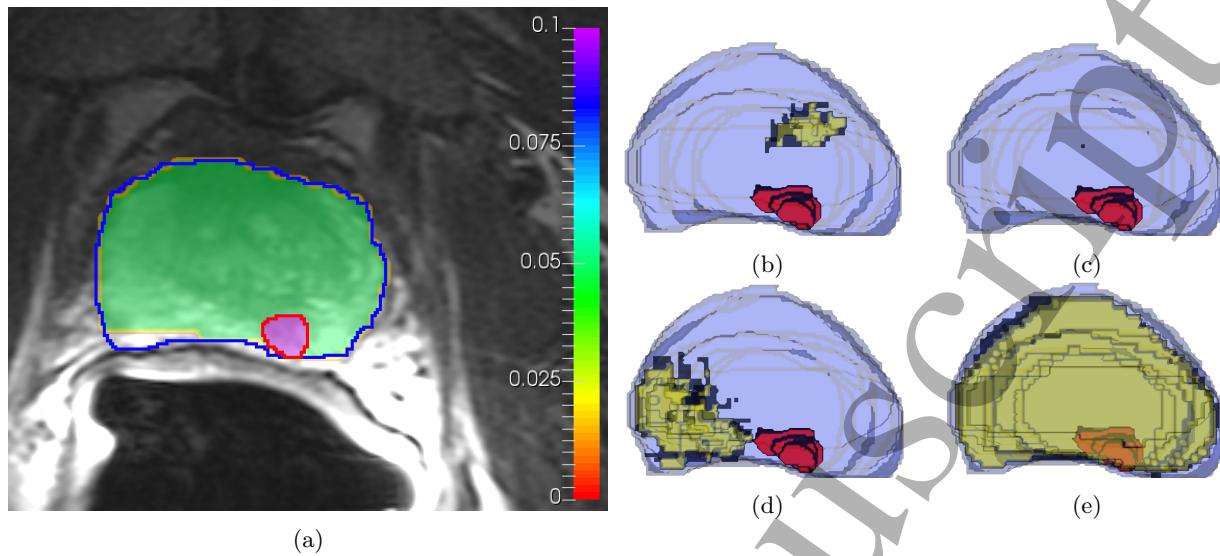


Figure 8: Microscopic tumor realizations generated for a prostate cancer patient. (a) The color map represents the CTM computed using the constant marginal probability model described in Section 2.4, overlaid on the diagnostic MRI. The blue contour delineates the anatomical barrier considered in the MTP simulation—the prostate—while the red contour outlines the GTV. (b–e) Different 3D realizations of MTP generated using a 1 mm resolution grid with $T(q, N) = 0.02$. The blue structure represents the prostate, the reddish/purple structure corresponds to the GTV, and the yellowish/greenish regions depict the microscopic tumor simulated with our model. Specifically, (b) a realization with a single islet of 429 mm^3 , approximately close to the literature-reported mean V_T ; (c) a realization with a small islet of 1 mm^3 , near the 5th percentile of the reported V_T distribution; (d) a realization with an islet of 1849 mm^3 , near the 95th percentile of the reported V_T distribution; and (e) an outlier case in which the tumor has infiltrated the entire prostate.

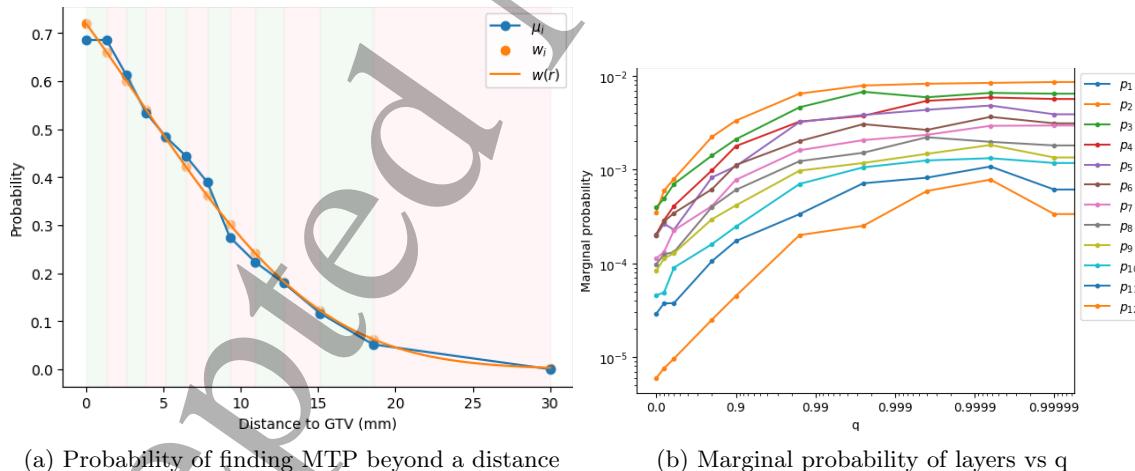


Figure 9: Estimation of marginal probabilities by solving the optimization problem described in Equation 13. (a) A representative example ($q = 0.99999$) illustrating how well μ_i approximates w_i for the optimized marginal probabilities. Blue dots represent μ_i , while orange dots correspond to w_i . The background shading highlights the layer boundaries, with green and red indicating the extent of each layer, starting with layer 1 on the left. (b) Marginal probability estimated for each layer across different values of q .

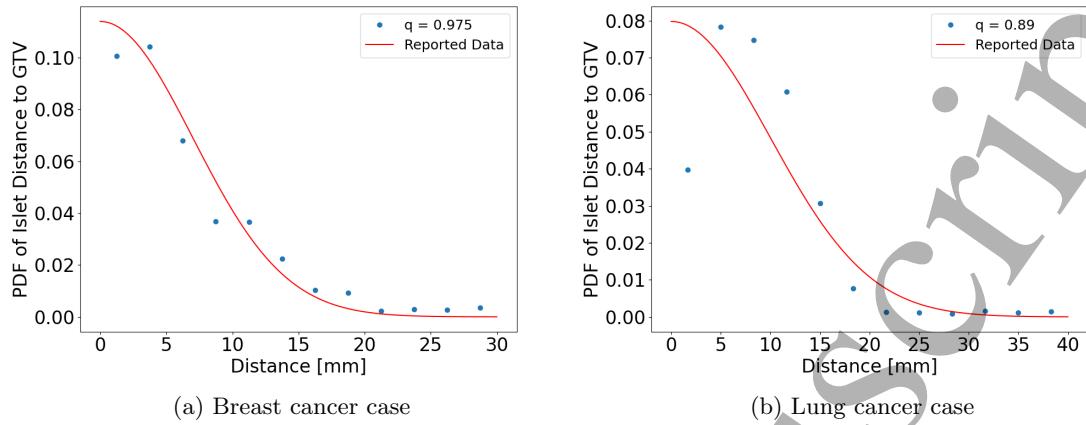


Figure 10: Probability density function (PDF) of the radial positions of microscopic tumor islets measured from the boundary of the GTV for breast and lung cancer cases. The scatter plots represent simulation results generated using our variable marginal probability model. The red line corresponds to the empirical function fitted by Stroom et al. based on their observations of the radial distribution of microscopic tumor islets [25].

creasing for voxels located farther from the target (Fig. 11a). The maximum marginal probability in the CTM was 0.007, occurring in the first layer. Using μ , we determined the distance beyond which the probability of finding any tumor voxel drops below 0.05, which was 20 mm. The isotropic expansion of the GTV with this distance represents the isoprobability shell 0.05. It can be noticed that this isoprobability shell is smaller than CTM.

Because the marginal probability decreases with increasing distance from the GTV, the simulations indicated that tumor islets preferentially clustered near the GTV boundary. Figure 11b–e presents representative tumor configurations:

- (b) A realization with a volume approximately close to the mean of the V_T distribution, with $V_T = 162 \text{ mm}^3$.
- (c) A realization near the 5th percentile of the V_T distribution, with $V_T = 8 \text{ mm}^3$.
- (d) A realization near the 95th percentile of the V_T distribution, with $V_T = 552 \text{ mm}^3$.
- (e) An outlier case with an exceptionally large tumor volume ($V_T = 1504 \text{ mm}^3$).

The simulations also showed that MTP was absent in 28.7% of the cases, and that the tumor islets frequently exhibited irregular, non-spherical shapes. In some realizations, streak-like MTP patterns were observed, similar to those illustrated in Figure 11d. Moreover, no tested value of q resulted in an MTP distribution that fully occupied the sampling domain.

3.3 Lung simulations

Similar to the breast cancer patient, the marginal probabilities across layers satisfy $p_1 \geq p_2 \geq \dots \geq p_{12}$ when voxel correlations are present, with probabilities decreasing as q decreases. It was also observed that the model was unable to replicate the \bar{N}_I reported by Stroom et al. [25]. While their study reported an \bar{N}_I of 2.47, the highest value obtained in our simulations was only 0.70.

Nevertheless, when setting $q = 0.89$, the simulated \bar{V}_T matched the estimate from Stroom et al. (Appendix G) [25]. Under this condition, 10,000 simulations showed that the distribution of islets as a function of their

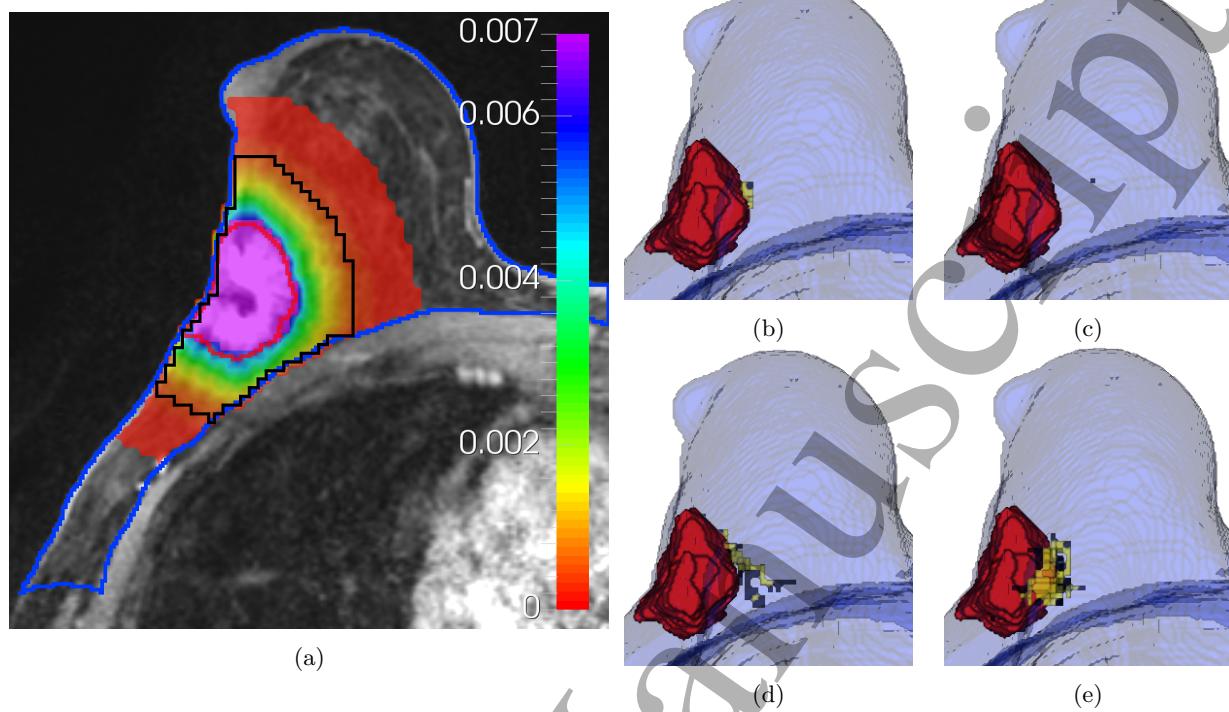


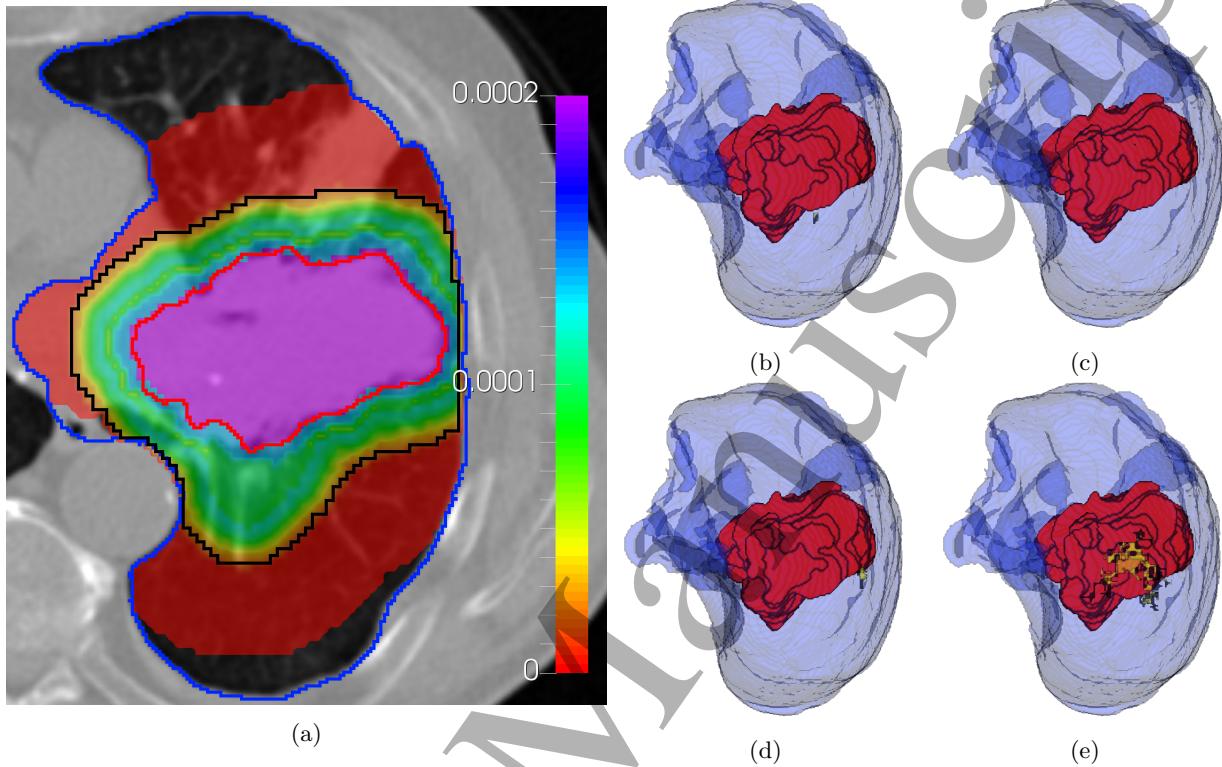
Figure 11: CTM and microscopic tumor realizations generated for a breast cancer patient. (a) The color map displays the MTP marginal probability calculated using the variable marginal probability model described in Section 2.5, overlaid on the diagnostic MRI. The blue contour represents the simulation space, accounting for anatomical barriers, the red contour outlines the GTV, and the black contour outlines the isoprobability shell 0.05. (b–e) Different 3D realizations of MTP generated using a 2 mm resolution grid and a setting of $q = 0.975$. The blue structure represents the sampling space, the reddish/purple structure corresponds to the GTV, and the yellowish/greenish regions depict the microscopic tumor realizations simulated with our model. Specifically, (b) a realization with a volume approximately close to the mean of the V_T distribution, with $V_T = 162 \text{ mm}^3$; (c) a single voxel corresponding to a realization near the 5th percentile of the V_T distribution, with $V_T = 8 \text{ mm}^3$; (d) a realization near the 95th percentile of the V_T distribution, with $V_T = 552 \text{ mm}^3$; and (e) an outlier case with an exceptionally large tumor volume ($V_T = 1504 \text{ mm}^3$).

distance from the GTV followed a trend similar to that reported in their study. However, the agreement was less strong than in the breast cancer case (Fig. 10b), as indicated by an R^2 of 0.63 and a MAE of 0.011 between our simulations and the function used in the paper.

The CTM reached a maximum marginal probability of 0.00014 and gradually decreased beyond the GTV boundary, with several layers exhibiting similar marginal probabilities (Fig. 12a). Despite the low marginal probabilities, the 0.05 isoprobability shell derived from the simulations extended 17 mm beyond the GTV boundary. The simulations also revealed the presence of small islets clustered near the GTV border. Figure 12b–e shows representative tumor configurations:

- (b) A realization with a volume approximately close to the mean of the V_T distribution, with $V_T = 82 \text{ mm}^3$.
- (c) A realization near the 5th percentile of the V_T distribution, with $V_T = 8 \text{ mm}^3$.
- (d) A realization near the 95th percentile of the V_T distribution, with $V_T = 320 \text{ mm}^3$.
- (e) An outlier case with an exceptionally large tumor volume ($V_T = 6616 \text{ mm}^3$).

1
2
3 The simulated islet volumes were consistently smaller than the GTV, and their shapes often displayed
4 irregular, non-spherical shapes. In some realizations, streak-like MTP patterns were observed, similar to
5 those illustrated in Figure 12e.
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 Figure 12: CTM and microscopic tumor realizations generated for a lung cancer patient. (a) The color map
33 displays the MTP marginal probability calculated using the variable marginal probability model described
34 in Section 2.5, overlaid on the diagnostic CT. The blue contour represents the simulation space, accounting
35 for anatomical barriers, the red contour outlines the GTV, and the black contour outlines the isoproba-
36 bility shell 0.05 generated from the simulations. (b–e) Different 3D realizations of MTP generated using
37 a 2 mm resolution grid and a setting of $q = 0.89$. The blue structure represents the sampling space, the
38 reddish/purple structure corresponds to the GTV, and the yellowish/greenish regions depict the microscop-
39 ic tumor realizations simulated with our model. Specifically, (b) a realization with a volume approxi-
40 mately close to the mean of the V_T distribution ($V_T = 82 \text{ mm}^3$); (c) a single voxel that corresponds to a realization
41 near the 5th percentile of the V_T distribution ($V_T = 8 \text{ mm}^3$); (d) a realization near the 95th percentile of
42 the V_T distribution ($V_T = 320 \text{ mm}^3$); and (e) an outlier case with an exceptionally large tumor volume
43 ($V_T = 6616 \text{ mm}^3$).
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

4 Discussion

The questions raised in the ICRU 83 report, as discussed in the introduction, highlight the necessity of clear statistical definitions and practical methods for incorporating probabilities into target definition and radiotherapy treatment planning [16]. In particular, the first question—how to define the “probability of occurrence” of tumor infiltration in a given tissue from a statistical perspective—can be directly addressed using the stochastic model presented in this work.

Our model provides a rigorous framework to estimate the probability of specific tumor configurations by computing the joint probability distribution or by employing Monte Carlo simulations. Through these simulations, various realizations of microscopic tumor presence (MTP) can be generated, and the probability of tumor presence in a given volume can be estimated as the fraction of simulations in which MTP occurs within that region.

The second question—on what clinical evidence these probabilities should be based—remains a challenge. In this study, available data consisted of population-based analyses that provide averages and standard deviations of microscopic tumor characteristics, such as islet radii, the percentage of patients with MTP, and the number of islets [28, 22, 10, 25]. While this data allow for the calibration of our stochastic model, more detailed information, including tumor staging and three-dimensional histopathological reconstructions, could enable more precise calibration. In particular, 3D histological reconstructions of MTP samples could be used to optimize the model parameters through maximum likelihood estimation, providing a more robust statistical foundation for our approach.

Finally, the third question—how to translate probabilities assigned to tissues into discretized grids for treatment planning—has been partially addressed in the literature. For instance, Bortfeld et al. demonstrated how microscopic tumor probabilities can be integrated with TCP in the treatment planning process [5]. Although their approach is computationally intensive, even for simplified models, it demonstrates the feasibility of incorporating probabilistic information into treatment optimization. However, to date, all models that incorporate probability during optimization rely solely on the clinical target map (CTM), neglecting voxel correlation. These models incorporated target uncertainty into treatment planning by using CTMs as weighting factors in the optimization loss function [3, 7, 24], this approach may result in underdosing regions where the actual risk of MTP exceeds 0.05. For instance, applying this method with the CTMs generated in our simulations would likely underdose the entire region outside the GTV, since the marginal voxel-level probabilities were all below 0.028 in the three cancer patients. However, the isoprobability surface defined at 0.05 corresponds to a much larger combined probability—where the true probability of microscopic tumor presence can reach values as high as 0.95. Consequently, relying only on marginal probabilities would lead to insufficient dose coverage. This underscores the importance of accounting for voxel correlations to ensure adequate target coverage during treatment optimization. It is also worth noting that, despite the relatively small marginal probabilities predicted by our models, the isoprobability shell at 0.05 still extends a considerable distance from the GTV border in the breast and lung cases. This occurs because the shell reflects the probability of detecting at least one tumor cell outside of the entire volume it encloses, rather than the probability at any single voxel.

Building on this motivation, a further way in which the probabilities generated by our stochastic model can support treatment planning is through the definition of non-isotropic CTVs in Euclidean space. If the CTV is defined as the volume that encompasses a given percentage of MTP realizations, our approach enables the exploration of all CTV shapes that satisfy this criterion by analyzing a large ensemble of simulations. Although these candidate CTVs correspond to isotropic expansions in the MTP probability space, their spatial manifestations in Euclidean space can become non-isotropic. This feature could be exploited in situations where conventional margin approaches, which typically rely on binary targets, are desired but the proximity of OARs to the GTV creates trade-offs between target coverage and OAR sparing. In such cases, one would prefer CTVs that carve out the OAR for preferential sparing, while extending farther in directions where compromise is less severe. We note that similar anisotropic solutions naturally arise in the dose space under probabilistic optimization, where the optimizer concentrates dose in regions with higher MTP probability while sacrificing other small, lower probability regions, in order to achieve the OAR constraints [6, 5].

The models developed in this work provide a framework for capturing voxel correlations while generating probabilistic results consistent with clinical observations reported in the literature. In the constant marginal probability model, this consistency is enforced by imposing a hard constraint that the parameter P_0 matches the observed probability of detecting MTP in a patient. In the variable marginal probability model, we ensure consistency by requiring that one minus the probability of all voxels outside the i -th shell being tumor-free, $1 - \mu_i$, corresponds to the probability of detecting MTP beyond the i -th shell, denoted as w_i . Moreover, our models explain how variations in voxel correlation influence the selection of the marginal probabilities assigned to individual voxels, thereby affecting the likelihood of different tumor configuration states. This behavior could have important implications for treatment planning. By enabling the calculation of joint probabilities, our models offer a principled way to better guide the target–OARs compromise through the assessment of the probability of finding MTP near an OAR. They can also support more accurate probabilistic robust evaluations by allowing the sampling of spatially realistic tumor configurations.

In addition to these practical applications, the models offer structural advantages that further enhance their usability and relevance. One of the most significant benefits is the limited number of tunable parameters, each grounded in sound statistical reasoning. This structure facilitates straightforward calibration using clinical data. The constant marginal probability model, for example, is defined by only two parameters, P_0 and T . The parameter P_0 is directly linked to the proportion of patients without MTP, allowing for immediate initialization based on available data. This, in turn, simplifies the tuning of T , as it remains the only free parameter in the model. A similar advantage is observed in the variable marginal probability model, where the probability of tumor presence beyond a given distance from the GTV, $w(r)$, is used to determine $p_1, \dots, p_{N'}$. The remaining parameter, q , is then calibrated to align with clinical observations, maintaining the model's adaptability while keeping it interpretable. The calibrated parameters that resulted from our study were $T = 0.02$, $q = 0.997$, and $q = 0.89$ for the prostate, breast, and lung cancer cases, respectively. Notably, these values lie well within the admissible parameter ranges, which explains why the resulting realizations display intermediate clustering behavior—neither isolated voxels nor entirely solid regions of MTP.

The parameter T reported in our study is not patient-specific, as it is defined independently of the total number of voxels in the prostate and calibrated to reproduce population-based MTP metrics. In contrast, the variable marginal probability model requires patient-specific calibration of q , since the probability of finding MTP in a given layer depends on the number of voxels within that layer. As detailed in the Methods section, T and q regulate the extent to which tumor voxels cluster together. From a biological perspective, these parameters can be interpreted as indicators of tumor aggressiveness, as they determine both the number and size of microscopic tumor islets (Figures 6 and 7b).

Our models naturally converge to previously established MTP models under specific conditions. The variable marginal probability model described in this work reduces to the model presented by Shusharina et al. as $q \rightarrow 0$ [23]. In the 1D case, it also converges to the model proposed by Buti et al. as $q \rightarrow 1$ [7], since Buti's model is fundamentally one-dimensional, assuming isotropy in the radial direction. Unlike these earlier models, our approach can also capture intermediate regimes, where voxels are neither fully independent nor fully correlated, i.e., a mixed correlation state. This was observed as the realizations in all the three cancer cases were clusters of many voxels of MTP rather than a single voxel. Beyond parameter calibration, the simplicity of our model's growing pattern also facilitates the computation of the joint probability of different tumor configurations. As outlined in the Methods section, the joint probability of a specific tumor configuration can be determined using Bayes' theorem, provided that the starting growing voxel is known.

While a model with few parameters offers certain advantages, it also has inherent limitations, particularly its reduced flexibility in simultaneously capturing multiple microscopic tumor properties. This is evident in the three patient cases: for prostate simulations, $T = 0.6$ generated realizations that matched the reported \bar{N}_I but not \bar{V}_T , whereas $T = 0.02$ matched \bar{V}_T but not the \bar{N}_I . Similar discrepancies were observed in the breast and lung simulations. These examples illustrate that, although the model provides a reasonable first-order approximation, it cannot perfectly reproduce all tumor metrics simultaneously. Introducing additional parameters could potentially improve the simultaneous fitting of different microscopic tumor features; however, the limited availability of reliable MTP data from histopathology studies poses significant challenges

1
2
3 547 for calibrating more complex, multi-parameter models.
4

5 548 Another limitation of our approach lies in its dependence on the specific sampling algorithm illustrated
6 549 in Figure 2. The algorithm tends to favor growth along the principal axes (x , y , z) originating from the
7 550 randomly selected seed voxel. This can lead to the appearance of “streaks” structures in some realizations
8 551 (as observed, for example, in Figs. 11d & 12e), resulting from stronger voxel-to-voxel Pearson correlations
9 552 along the grid axes compared to off-axis directions. Consequently, alternative sampling strategies, such as
10 553 incorporating second-order neighborhoods, adopting non-axial growth geometries, or varying the correlation
11 554 parameter across voxels, could produce realizations with different spatial patterns or morphologies.

12 555 It is important to note that our model was primarily designed to reproduce population-level characteristics
13 556 of microscopic tumor presence reported in the literature, such as the total microscopic tumor volume and
14 557 the number of tumor islets, both of which are independent of the specific tumor shape. Generating more
15 558 anatomically realistic spatial patterns would require voxel-level tumor infiltration data from histopathology
16 559 to more accurately calibrate spatial correlations. However, such detailed data are currently very difficult to
17 560 obtain and were not available for this study.

18 561 In comparison with the work of Stroom et al., who proposed a CTV margin calculation formula based
19 562 on the standard deviation of the islet distribution and a correction factor derived from radiotherapy sim-
20 563 ulations, our approach provides an alternative, data-driven perspective. Although we do not introduce an
21 564 explicit analytical expression for margin estimation, the CTV margin naturally emerges from our model as
22 565 the isotropic expansion of the GTV that encloses 95% of the simulated MTP realizations. Applying this
23 566 approach yields CTV margins of 20 mm for breast cancer and 17 mm for lung cancer, which are in the
24 567 order of magnitude to the margins predicted by the Stroom formula (15.4 mm and 22.6 mm, respectively).
25 568 The remaining differences of approximately 5 mm can reasonably be attributed to the distinct underlying
26 569 assumptions of the two methods.

27 570 Future work will focus on integrating both the CTM and voxel correlation information into the treatment
28 571 planning process to enable probabilistically robust treatment strategies. Additionally, we aim to investigate
29 572 more physically grounded models, such as the Ising model and its variants, to further refine the representation
30 573 of microscopic tumor spread and improve the accuracy of treatment optimization [26, 2].

31 574 **5 Conclusion**

32 575 In this study, we developed two first-principles stochastic models to simulate microscopic tumor presence
33 576 (MTP) under different assumptions, incorporating a correlation factor to account for spatial dependencies
34 577 between neighboring voxels. This approach enables a more realistic representation of tumor infiltration
35 578 patterns compared to existing methods. By leveraging these models, we can estimate both the marginal
36 579 probability of MTP and the joint probability of different tumor configurations, providing a complete proba-
37 580 bilistic description of MTP. Our simulations reproduced key features of microscopic tumor extension reported
38 581 for prostate, breast, and lung cancers, effectively capturing their distinct spatial patterns.

39 582 Importantly, our findings highlight the limitations of using marginal probabilities alone as weights in
40 583 treatment optimization and emphasize the necessity of incorporating voxel correlations to properly account
41 584 for spatial uncertainty. This work lays the foundation for integrating probabilistic modeling into the descrip-
42 585 tion of MTP, with potential implications for improving radiotherapy treatment planning.

43 586 **Acknowledgements**

44 587 This project has received funding from the European Union’s Horizon 2020 Marie Skłodowska-Curie Actions
45 588 under Grant Agreement No. 955956. The work was also supported by the National Cancer Institute of the
46 589 United States under grant number R01CA266275, and by the Therapy Imaging Program (TIP) funded by
47 590 the Federal Share of program income earned by Massachusetts General Hospital on C06CA059267, Proton
48 591 Therapy Research and Treatment Center.

592 References

- [1] Hugo J. W. L. Aerts et al. *Data From NSCLC-Radiomics*. <https://doi.org/10.7937/K9/TCIA.2015.PF0M9REI>. Version 4. 2014.
- [2] L. Amoudruz, G. Buti, L. Rivetti, A. Ajdari, G. Sharp, P. Koumoutsakos, S. Spohn, A. L. Grosu, and T. Bortfeld. “Ising energy model for the stochastic prediction of tumor islets”. In: *arXiv preprint arXiv:2508.20804v2* (Sept. 2025). PMID: 41001577; PMCID: PMC12458594.
- [3] Ivar Bengtsson, Anders Forsgren, and Albin Fredriksson. “Implications of using the clinical target distribution as voxel-weights in radiation therapy optimization”. In: *Physics in Medicine & Biology* 68.9 (Apr. 2023), p. 095005. DOI: 10.1088/1361-6560/acc77b. URL: <https://dx.doi.org/10.1088/1361-6560/acc77b>.
- [4] Kamalika Bhattacharjee, Nazma Naskar, Souvik Roy, and Sukanta Das. “A survey of cellular automata: types, dynamics, non-uniformity and applications”. In: *Natural Computing* 19.2 (2020), pp. 433–461. ISSN: 1572-9796. DOI: 10.1007/s11047-018-9696-8. URL: <https://doi.org/10.1007/s11047-018-9696-8>.
- [5] Thomas Bortfeld, Nadya Shusharina, and David Craft. “Probabilistic definition of the clinical target volume—implications for tumor control probability modeling and optimization”. In: *Physics in Medicine & Biology* 66.1 (Jan. 2021), 01NT01. DOI: 10.1088/1361-6560/abcad8. URL: <https://dx.doi.org/10.1088/1361-6560/abcad8>.
- [6] Gregory Buti, Nadya Shusharina, Ali Ajdari, Edmond Sterpin, and Thomas Bortfeld. “Exploring trade-offs in treatment planning for brain tumor cases with a probabilistic definition of the clinical target volume”. In: *Medical Physics* 50.1 (2023), pp. 410–423. DOI: <https://doi.org/10.1002/mp.16097>. eprint: <https://aapm.onlinelibrary.wiley.com/doi/pdf/10.1002/mp.16097>. URL: <https://aapm.onlinelibrary.wiley.com/doi/abs/10.1002/mp.16097>.
- [7] Gregory Buti, Kevin Souris, Ana Maria Barragán Montero, John Aldo Lee, and Edmond Sterpin. “Introducing a probabilistic definition of the target in a robust treatment planning framework”. In: *Physics in Medicine & Biology* 66.15 (July 2021), p. 155008. DOI: 10.1088/1361-6560/ac1265. URL: <https://dx.doi.org/10.1088/1361-6560/ac1265>.
- [8] Sorcha Campbell et al. “Evaluation of Microscopic Disease in Oral Tongue Cancer Using Whole-Mount Histopathologic Techniques: Implications for the Management of Head-and-Neck Cancers”. In: *International Journal of Radiation Oncology*Biology*Physics* 82.2 (2012), pp. 574–581. ISSN: 0360-3016. DOI: <https://doi.org/10.1016/j.ijrobp.2010.09.038>. URL: <https://www.sciencedirect.com/science/article/pii/S036030161003347X>.
- [9] Danna Daniels, Dorit Last, Katya Cohen, Yael Mardor, and Michal Sklair-Levy. *Standard and Delayed Contrast-Enhanced MRI of Malignant and Benign Breast Lesions with Histological and Clinical Supporting Data (Advanced-MRI-Breast-Lesions)*. <https://doi.org/10.7937/C7X1-YN57>. Version 2. 2024.
- [10] Andriy Fedorov, Clare Tempany, Robert Mulkern, and Fiona Fennessy. *Data From QIN PROSTATE*. <https://doi.org/10.7937/K9/TCIA.2016.fADs26kG>. 2016. DOI: 10.7937/K9/TCIA.2016.fADs26kG.
- [11] Claudio Fiorino, Robert Jeraj, Catharine H. Clark, Cristina Garibaldi, Dietmar Georg, Ludvig Muren, Wouter van Elmpt, Thomas Bortfeld, and Nuria Jornet. “Grand challenges for medical physics in radiation oncology”. In: *Radiotherapy and Oncology* 153 (Dec. 2020), pp. 7–14. ISSN: 0167-8140. DOI: 10.1016/j.radonc.2020.10.001. URL: <https://doi.org/10.1016/j.radonc.2020.10.001>.
- [12] B. Fleury, J. Thariat, R. Barnoud, G. Buiret, F. Lebreton, B. Bancel, M. Poupart, and M. Devouassoux-Shisheboran. “Approche anatomopathologique de l’extension microscopique des carcinomes épidermoïdes ORL : implications pour la définition du volume cible anatomoclinique”. In: *Cancer/Radiothérapie* 18.7 (2014), pp. 666–671. ISSN: 1278-3218. DOI: <https://doi.org/10.1016/j.canrad.2014.04.006>. URL: <https://www.sciencedirect.com/science/article/pii/S1278321814000742>.

- 1
2
3 [639] [13] Vincent Grégoire et al. "Delineation of the primary tumour Clinical Target Volumes (CTV-P) in
4 laryngeal, hypopharyngeal, oropharyngeal and oral cavity squamous cell carcinoma: AIRO, CACA,
5 DAHANCA, EORTC, GEORCC, GORTEC, HKNPCSG, HNCIG, IAG-KHT, LPRHHT, NCIC CTG,
6 NCRI, NRG Oncology, PHNS, SBRT, SOMERA, SRO, SSHNO, TROG consensus guidelines". In:
7 *Radiotherapy and Oncology* 126.1 (Jan. 2018), pp. 3–24. ISSN: 0167-8140. DOI: 10.1016/j.radonc.
8 2017.10.016. URL: <http://dx.doi.org/10.1016/j.radonc.2017.10.016>.
- 9 [640] [14] Tim Head, Manoj Kumar, Holger Nahrstaedt, Gilles Louppe, and Jaroslav Shcherbatyi. *scikit-optimize/scikit-*
10 *optimize: v0.8.1.* Version v0.8.1. Sept. 2020. DOI: 10.5281/zenodo.4014775. URL: <https://doi.org/10.5281/zenodo.4014775>.
- 11 [641] [15] Petra J. van Houdt et al. "Histopathological Features of MRI-Invisible Regions of Prostate Cancer
12 Lesions". In: *Journal of Magnetic Resonance Imaging* 51.4 (2020), pp. 1235–1246. DOI: <https://doi.org/10.1002/jmri.26933>. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/jmri.26933>. URL: <https://onlinelibrary.wiley.com/doi/abs/10.1002/jmri.26933>.
- 13 [642] [16] International Commission on Radiation Units and Measurements. *Prescribing, Recording, and Re-*
14 *porting Intensity-Modulated Photon-Beam Therapy (IMRT)*. Tech. rep. Report 83. ICRU, 2010. URL:
15 <https://www.icru.org/report/prescribing-recording-and-reporting-intensity-modulated-photon-beam-therapy-imrticru-report-83/>.
- 16 [643] [17] Xue Meng, Xindong Sun, Dianbin Mu, Ligang Xing, Li Ma, Baijiang Zhang, Shuqiang Zhao, Guoren
17 Yang, Feng-Ming (Spring) Kong, and Jinming Yu. "Noninvasive Evaluation of Microscopic Tumor
18 Extensions Using Standardized Uptake Value and Metabolic Tumor Volume in Non-Small-Cell Lung
19 Cancer". In: *International Journal of Radiation Oncology*Biology*Physics* 82.2 (Feb. 2012), pp. 960–
20 966. ISSN: 0360-3016. DOI: 10.1016/j.ijrobp.2010.10.064. URL: <http://dx.doi.org/10.1016/j.ijrobp.2010.10.064>.
- 21 [644] [18] Leyla Moghaddasi, Eva Bezak, and Loredana G. Marcu. "Current challenges in clinical target vol-
22 ume definition: Tumour margins and microscopic extensions". In: *Acta Oncologica* 51.8 (Sept. 2012),
23 pp. 984–995. ISSN: 1651-226X. DOI: 10.3109/0284186X.2012.720381. URL: <http://dx.doi.org/10.3109/0284186X.2012.720381>.
- 24 [645] [19] Ali Mohammadian Bajgiran, Sara Afshari Mirak, Shirin Shakeri, et al. "Characteristics of missed
25 prostate cancer lesions on 3T multiparametric-MRI in 518 patients: based on PI-RADSv2 and using
26 whole-mount histopathology reference". In: *Abdominal Radiology* 44.3 (2019), pp. 1052–1061. DOI:
27 10.1007/s00261-018-1823-6.
- 28 [646] [20] Christopher F. Njeh. "Tumor delineation: The weakest link in the search for accuracy in radiotherapy".
29 In: *Journal of Medical Physics / Association of Medical Physicists of India* 33 (2008), pp. 136–140.
30 URL: <https://api.semanticscholar.org/CorpusID:5262423>.
- 31 [647] [21] Nina N. Sanford et al. "Individualization of Clinical Target Volume Delineation Based on Stepwise
32 Spread of Nasopharyngeal Carcinoma: Outcome of More Than a Decade of Clinical Experience". In:
33 *International Journal of Radiation Oncology, Biology, Physics* 103.3 (Mar. 2019), pp. 654–668. ISSN:
34 0360-3016. DOI: 10.1016/j.ijrobp.2018.10.006. URL: <https://doi.org/10.1016/j.ijrobp.2018.10.006>.
- 35 [648] [22] Annemarie C. Schmitz, Maurice A.A.J. van den Bosch, Claudette E. Loo, Willem P.Th.M. Mali,
36 Harry Bartelink, Maria Gertenbach, Roland Holland, Johannes L. Peterse, Emiel J.Th. Rutgers, and
37 Kenneth G. Gilhuijs. "Precise correlation between MRI and histopathology – Exploring treatment
38 margins for MRI-guided localized breast cancer therapy". In: *Radiotherapy and Oncology* 97.2 (2010),
39 pp. 225–232. ISSN: 0167-8140. DOI: <https://doi.org/10.1016/j.radonc.2010.07.025>. URL:
40 <https://www.sciencedirect.com/science/article/pii/S0167814010004780>.
- 41 [649] [23] Nadya Shusharina, David Craft, Yen-Lin Chen, Helen Shih, and Thomas Bortfeld. "The clinical target
42 distribution: a probabilistic alternative to the clinical target volume". In: 63.15 (July 2018), p. 155001.
43 DOI: 10.1088/1361-6560/aacfb4. URL: <https://dx.doi.org/10.1088/1361-6560/aacfb4>.

- 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
- [24] A Smolders, I Bengtsson, A Forsgren, A Lomax, D C Weber, A Fredriksson, and F Albertini. "Robust optimization strategies for contour uncertainties in online adaptive radiation therapy". In: *Physics in Medicine & Biology* 69.16 (July 2024), p. 165001. DOI: 10.1088/1361-6560/ad6526. URL: <https://dx.doi.org/10.1088/1361-6560/ad6526>.
- [25] Joep Stroom, Kenneth Gilhuijs, Sandra Vieira, Wei Chen, Javier Salguero, Elizabeth Moser, and Jan-Jakob Sonke. "Combined Recipe for Clinical Target Volume and Planning Target Volume Margins". In: *International Journal of Radiation Oncology*Biology*Physics* 88.3 (2014), pp. 708–714. ISSN: 0360-3016. DOI: <https://doi.org/10.1016/j.ijrobp.2013.08.028>. URL: <https://www.sciencedirect.com/science/article/pii/S0360301613030344>.
- [26] Salvatore Torquato. "Toward an Ising model of cancer and beyond". In: *Physical Biology* 8.1 (Feb. 2011). Epub 2011 Feb 7, p. 015017. DOI: 10.1088/1478-3975/8/1/015017.
- [27] Jan Unkelbach et al. "The role of computational methods for automating and improving clinical target volume definition". In: *Radiotherapy and Oncology* 153 (Dec. 2020), pp. 15–25. ISSN: 0167-8140. DOI: 10.1016/j.radonc.2020.10.002. URL: <https://doi.org/10.1016/j.radonc.2020.10.002>.
- [28] Judith van Loon et al. "Microscopic Disease Extension in Three Dimensions for Non-Small-Cell Lung Cancer: Development of a Prediction Model Using Pathology-Validated Positron Emission Tomography and Computed Tomography Features". In: *International Journal of Radiation Oncology*Biology*Physics* 82.1 (2012), pp. 448–456. ISSN: 0360-3016. DOI: <https://doi.org/10.1016/j.ijrobp.2010.09.001>. URL: <https://www.sciencedirect.com/science/article/pii/S0360301610032384>.
- [29] Shalini K. Vinod, Michael G. Jameson, Myo Min, and Lois C. Holloway. "Uncertainties in volume delineation in radiation oncology: A systematic review and recommendations for future studies". In: *Radiotherapy and Oncology* 121.2 (2016), pp. 169–179. ISSN: 0167-8140. DOI: <https://doi.org/10.1016/j.radonc.2016.09.009>. URL: <https://www.sciencedirect.com/science/article/pii/S0167814016343316>.
- [30] E. Weiss and C. Hess. "The Impact of Gross Tumor Volume (GTV) and Clinical Target Volume (CTV) Definition on the Total Accuracy in Radiotherapy". In: *Strahlentherapie und Onkologie* 179 (2003), pp. 21–30. DOI: 10.1007/s00066-003-0976-5. URL: <https://doi.org/10.1007/s00066-003-0976-5>.
- [31] Constantinos Zamboglou et al. "Uncovering the invisible—prevalence, characteristics, and radiomics feature-based detection of visually undetectable intraprostatic tumor lesions in 68GaPSMA-11 PET images of patients with primary prostate cancer". In: *European Journal of Nuclear Medicine and Molecular Imaging* 48 (2020), pp. 1987–1997. URL: <https://api.semanticscholar.org/CorpusID:227034057>.

1 2 3 A Appendix: Relation of Pearson Correlation Coefficient with 4 Parameter q 5

6
7 722 Given a 2D grid with 9 elements, where each voxel has the same marginal probability p (Fig 13), the Pearson
8 correlation coefficient between two neighbor voxels, C_8 and C_9 , is defined as:
9

$$\rho_{C_8, C_9} = \frac{\mathbb{E}[C_8 C_9] - \mathbb{E}[C_8]\mathbb{E}[C_9]}{\sqrt{\mathbb{E}[C_8^2] - (\mathbb{E}[C_8])^2}\sqrt{\mathbb{E}[C_9^2] - (\mathbb{E}[C_9])^2}}. \quad (14)$$

10
11 724 Because the state of each voxel is Bernoulli distributed, $\mathbb{E}[C_8] = \mathbb{E}[C_8^2] = \mathbb{E}[C_9] = \mathbb{E}[C_9^2] = p$. Therefore,
12
13

$$\rho_{C_8, C_9} = \frac{\mathbb{E}[C_8 C_9] - p^2}{p - p^2}. \quad (15)$$

c_1	c_2	c_3
c_4	c_5	c_6
c_7	c_8	c_9

p	p	p
p	p	p
p	p	p

20
21
22
23
24
25
26
27
28 Figure 13: 2D 9-voxel grid with a uniform marginal probability distribution. The left image labels each voxel
29 with its corresponding random variable, while the right image illustrates the marginal probability assigned
30 to each voxel.
31

32
33 725 The key challenge is determining $\mathbb{E}[C_8 C_9]$. Since our method involves randomly selecting a voxel in the
34 grid as the initial growth point, the expectation is computed by averaging over all possible starting voxels.
35 Let $\mathbb{E}_i[C_8 C_9]$ denote the expected value given that the growth starts at voxel i . Then, the overall expectation
36 is:
37

$$\mathbb{E}[C_8 C_9] = \frac{1}{9} \sum_{i=1}^9 \mathbb{E}_i[C_8 C_9] = \frac{1}{9} \sum_{i=1}^9 \sum_{\text{all states}} c_8 c_9 P_i(\text{state}), \quad (16)$$

40
41 729 where $P_i(\text{state})$ represents the joint probability of a particular state given that voxel i was the starting
42 point. By introducing equation (16) in equation (15), we obtain:
43
44

$$\rho_{C_8, C_9} = \frac{1}{9} \sum_{i=1}^9 (\rho_{C_8, C_9})_i, \quad (17)$$

45
46 731 where $(\rho_{C_8, C_9})_i$ is the Pearson correlation coefficient between C_8 and C_9 , considering voxel i as the
47 starting point. Each $(\rho_{C_8, C_9})_i$ can be expressed as a weighted sum of different powers of the correlation
48 factor q . For instance:
49
50
51
52
53
54
55
56
57
58
59
60

$$\begin{aligned}
 (\rho_{C_8, C_9})_7 &= q, \\
 (\rho_{C_8, C_9})_5 &= \frac{1}{2}(q + q^3), \\
 (\rho_{C_8, C_9})_4 &= \frac{1}{2} \left(q + \frac{1}{2}q^3 + \frac{1}{2}q^5 \right), \\
 (\rho_{C_8, C_9})_1 &= \frac{1}{2} \left(q + \frac{1}{3}q^3 + \frac{1}{3}q^5 + \frac{1}{3}q^7 \right).
 \end{aligned}$$

A key observation is that the closer the starting voxel is to the voxels whose correlation is being assessed, the smaller the exponent of q . This behavior arises from the design of the visiting algorithm. Additionally, in each case, the coefficients multiplying q sum to 1, ensuring that when $q = 1$, the Pearson correlation coefficient equals 1, as expected.

Although this example is presented in 2D, the same visiting rules apply in 3D, making these results generalizable. Then, we can express the correlation coefficient between two neighboring voxels as:

$$\rho = \sum_{n=1}^{\infty} a_n q^{2n-1}, \quad (18)$$

where a_n are coefficients that satisfy:

$$\sum_{n=1}^{\infty} a_n = 1. \quad (19)$$

Note that the magnitude of a_n will depend on the size of the grid and the pair of voxels that we want to assess the correlation.

B Appendix: Taylor Series Approximation of Equation 6

To approximate Equation 6 using a Taylor series, we first rewrite it as:

$$\left(\frac{P_0}{1-p} \right)^{\frac{1}{N-1}} = 1 - p + pq. \quad (20)$$

Assuming that the number of voxels N is large and thus $x = \frac{1}{N-1}$ is small, we expand the term $\left(\frac{P_0}{1-p} \right)^x$ in a Taylor series of x around $x = 0$:

$$\sum_{n=0}^{\infty} \left[\left(\ln \frac{P_0}{1-p} \right)^n \frac{1}{n!} x^n \right] = 1 - p + pq. \quad (21)$$

By neglecting terms of the series with orders $n = 2$ and higher, we obtain:

$$1 + \left(\ln \frac{P_0}{1-p} \right) x \approx 1 - p + pq. \quad (22)$$

After re-inserting $x = \frac{1}{N-1}$, this leads to

$$\frac{P_0}{1-p} \approx e^{-(N-1)(1-q)p}. \quad (23)$$

Finally, since N is assumed to be large, we approximate $N - 1$ as N in the exponential term, resulting in:

$$\frac{P_0}{1-p} \approx e^{-N(1-q)p}. \quad (24)$$

1 2 3 751 C Appendix: Conditional probability variable marginal probabili- 4 752 ty model 5

6
7 753 Consider a 1D grid (Fig. 4a) where each voxel follows the *non-tunneling assumption*. For two neighboring
8 754 voxels, C_i and C_k , with marginal probabilities p_i and p_k such that $p_i \geq p_k$, the conditional probabilities are
9 755 defined as:

$$10 \quad \begin{aligned} \mathbb{P}(C_k = 1 | C_i = 1) &= aq + (1 - q)p_k, \\ 11 \quad \mathbb{P}(C_k = 1 | C_i = 0) &= (1 - q)p_k, \\ 12 \quad \mathbb{P}(C_i = 1 | C_k = 1) &= q + (1 - q)p_i, \\ 13 \quad \mathbb{P}(C_i = 1 | C_k = 0) &= bq + (1 - q)p_i, \end{aligned} \quad (25)$$

14
15 756 Here, a and b are scalar factors that modulate the probability of copying or flipping the state of the
16 757 neighboring voxel.

17
18 758 To determine a and b , we use the expectation constraints:

$$19 \quad \mathbb{E}[C_i] = p_i, \quad \mathbb{E}[C_k] = p_k. \quad (26)$$

20
21 759 Expanding these expectations:

$$22 \quad p_k = \frac{1}{2} \left(p_k + p_i(aq + (1 - q)p_k) + (1 - p_i)(1 - q)p_k \right), \quad (27)$$

$$23 \quad p_i = \frac{1}{2} \left(p_i + p_k(q + (1 - q)p_i) + (1 - p_k)(bq + (1 - q)p_i) \right). \quad (28)$$

24
25 760 Solving for a and b , we obtain:

$$26 \quad a = \frac{p_k}{p_i}, \quad b = \frac{p_i - p_k}{1 - p_k}. \quad (29)$$

27
28 761 These values ensure that the marginal probabilities remain consistent across the system.

29 762 D Appendix: Data Calculated from Bajgira et al.

30
31 763 In this section, we describe how the prostate cancer population data from Bajgira et al. [19] was obtained
32 764 and processed for our analysis.

33 765 Patient and Tumor Classification

34
35 766 Table 1 of their study reports histologic findings for 518 patients who underwent multiparametric MRI
36 767 followed by prostatectomy. Tumors in their paper can be classified into:

- 37 768 1. Index tumors (i): the main lesion in each patient, defined as the tumor with the highest Gleason Score.
38 769 For the purpose of our analysis, these were identified as the gross tumor volumes (GTVs).
- 39 770 2. Islet tumors (j): additional lesions with lower Gleason Score, considered multifocal satellites of the
40 771 index tumor.
- 41 772 3. Total tumors (T): the full set of index and islet tumors.

42
43 773 The study also distinguishes between lesions detected by MRI and those that were missed but later
44 774 identified by histopathology.

1
2
3 **Notation and Data Summary**

4
5 **775** For clarity, we define below the notation and numerical values extracted from Bajgira et al. [19], which are
6 **777** used throughout this appendix.

Symbol	Definition	Value
N_{patients}	Total number of patients	518
N_T	Total number of missed tumors	563
N_i	Number of missed index tumors	111
N_j	Number of missed islet tumors ($N_T - N_i$)	452
μ_T	Mean radius of missed tumors	9.4 mm
μ_i	Mean radius of missed index tumors	15.8 mm
μ_j	Mean radius of missed islet tumors	to be calculated
σ_T	Std. dev. of missed tumors	0.72 mm
σ_i	Std. dev. of missed index tumors	8.8 mm
σ_j	Std. dev. of missed islet tumors	to be calculated

19 **20** Table 3: Summary of variables and values used in the appendix.
21
22

23 **778 Proportion of Patients Without Microscopic Disease**

24
25 **779** Among the 518 patients, 396 presented multifocal disease (i.e., one or more islet tumors in addition to
26 **780** the index tumor). If we assume that every multifocal case includes at least one MRI-invisible lesion, the
27 **781** probability of a patient having no microscopic multifocal disease is estimated as:

$$P_0 = 1 - \frac{396}{518} = 0.28.$$

28
29 **782** This assumption is supported by the observation that of the
30
31

$$567 = \underbrace{1085}_{\text{Total tumors}} - \underbrace{518}_{\text{Index tumors}}$$

32
33 **783** islet tumors, only
34

$$115 = \underbrace{522}_{\text{Detected tumors}} - \underbrace{407}_{\text{Detected index tumors}}$$

35
36 **784** were detected by MRI, meaning the majority were invisible.
37
38

39 **785 Average Number of Missed Islet Tumors**

40
41 **786** The average number of missed islet tumors per patient is
42

$$\frac{N_T - N_i}{N_{\text{patients}}}.$$

43
44 **787** Substituting the reported values:
45

$$\frac{563 - 111}{518} = 0.87.$$

1
2
3 **Average Radius of Missed Islet Tumors**
4
5

6 788 The mean radius of missed islet tumors is obtained from the weighted difference:
7
8

$$\mu_j = \frac{\mu_T \cdot N_T - \mu_i \cdot N_i}{N_j}.$$

9 790 Substituting the reported values:
10
11

$$\mu_j = \frac{9.4 \text{ mm} \cdot 563 - 15.8 \text{ mm} \cdot 111}{452} = 7.82 \text{ mm.}$$

12
13 **Standard deviation of missed islet tumors**
14
15

16 791 We derive σ_j from first principles, starting with sums of squares.
17
18

19 792 Let the observed radii of the N_T missed tumors be x_1, \dots, x_{N_T} . Partition the indices so that the first N_i
20 are the missed index tumors and the remaining N_j are the missed islet tumors (with $N_T = N_i + N_j$). Then
21 the total sum of squares splits as
22

$$\sum_{n=1}^{N_T} x_n^2 = \sum_{n \in i} x_n^2 + \sum_{n \in j} x_n^2.$$

23 793 Divide both sides by N_T to obtain the total second moment:
24
25

$$\frac{1}{N_T} \sum_{n=1}^{N_T} x_n^2 = \frac{1}{N_T} \sum_{n \in i} x_n^2 + \frac{1}{N_T} \sum_{n \in j} x_n^2.$$

26
27 794 Recognizing subgroup averages,
28
29

$$\mathbb{E}[X^2]_T = \frac{N_i}{N_T} \left(\frac{1}{N_i} \sum_{n \in i} x_n^2 \right) + \frac{N_j}{N_T} \left(\frac{1}{N_j} \sum_{n \in j} x_n^2 \right),$$

30
31 795 SO
32

$$\boxed{\mathbb{E}[X^2]_T = \frac{N_i \mathbb{E}[X^2]_i + N_j \mathbb{E}[X^2]_j}{N_T}}.$$

33
34 796 Using $\mathbb{E}[X^2]_k = \sigma_k^2 + \mu_k^2$ for $k \in \{T, i, j\}$ gives
35
36

$$\sigma_T^2 + \mu_T^2 = \frac{N_i(\sigma_i^2 + \mu_i^2) + N_j(\sigma_j^2 + \mu_j^2)}{N_T}.$$

37
38 797 Solving for σ_j yields the population-standard deviation expression
39
40

$$\boxed{\sigma_j = \sqrt{\frac{N_T(\sigma_T^2 + \mu_T^2) - N_i(\sigma_i^2 + \mu_i^2)}{N_j} - \mu_j^2}}.$$

41
42 801 by replacing the values from table 3 we obtained $\sigma_j = 5.76 \text{ mm}$
43
44
45

46
47 **Average islet volume calculation**
48
49

50 802 To estimate the average volume of tumor islets (\bar{V}_I), we modeled their radius as a random variable
51
52

$$r \sim \mathcal{N}(\mu_j, \sigma_j),$$

1
2
3 where μ_j and σ_j are the mean and standard deviation of the islet radius obtained in the previous subsection.
4 Since tumor radii cannot be negative, we applied a truncated normal distribution restricted to $r > 0$.
5 Assuming spherical geometry, the volume of an islet is
6

$$V = \frac{4}{3}\pi r^3.$$

7 We generated a large sample of radii from the truncated normal distribution, computed the corresponding
8 spherical volumes, and summarized the resulting distribution. The estimated population mean volume was
9

$$\mathbb{E}[V] \approx 513 \text{ mm}^3,$$

10 with the 5th and 95th percentiles at 1 mm^3 and 2024 mm^3 , respectively, providing an empirical 90% interval
11 for the islet volumes.
12

17 811 Average Total Microscopic Tumor Volume

18 812 The Average Total Microscopic Tumor Volume (\bar{V}_T) per multifocal patient was calculated by multiplying \bar{V}_I
19 by the average number of islets per patient:
20

$$21 \quad \bar{V}_T = \bar{V}_I \times \frac{N_j}{N_{\text{multifocal patients}}} = 513 \text{ mm}^3 \times \frac{452}{396} \approx 584 \text{ mm}^3.$$

22 Similarly, the 5th and 95th percentiles of the \bar{V}_T distribution were obtained by scaling the corresponding
23 percentiles of \bar{V}_I distribution by the same factor:
24

$$25 \quad V_{5\%} \approx 1 \text{ mm}^3 \times 1.14 \approx 1.14 \text{ mm}^3, \quad V_{95\%} \approx 2024 \text{ mm}^3 \times 1.14 \approx 2307 \text{ mm}^3.$$

26 Thus, the \bar{V}_T provides an estimate of the total microscopic tumor burden from islets in multifocal patients,
27 with an empirical 90% interval of approximately $1\text{--}2307 \text{ mm}^3$.
28

32 818 E Prostate Simulation Results

33 819 In this section we show the simulations of the average islet volume (\bar{V}_I) for the prostate patient (Fig 14).
34

37 820 F Breast Simulation Results

38 821 In this section we show the simulations of the average total microscopic tumor volume (\bar{V}_T) for the breast
39 patient (Fig 15).
40

43 823 G Lung Simulations Results

44 824 In this section we show the simulations of the average total microscopic tumor volume (\bar{V}_T) for the lung
45 patient (Fig 16).
46

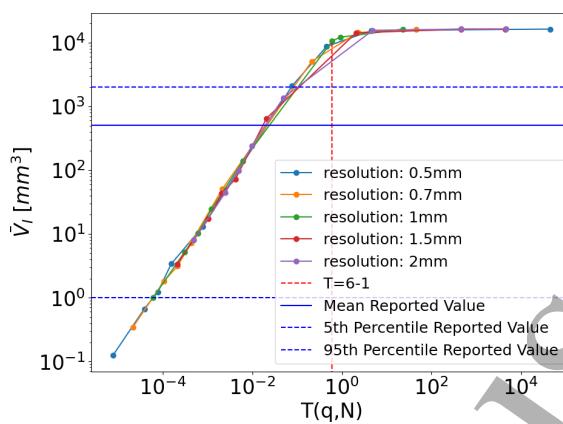


Figure 14: Simulated \bar{V}_I for the prostate case across different grid resolutions and values of $T(q, N)$. Each color represents a different grid resolution. The solid blue line indicates the \bar{V}_I reported by Bajgira et al., while the dashed black lines correspond to the 5th and 95th quantiles estimated from [19]. The red dashed line marks the value of T that produces the \bar{V}_I reported in Bajgira et al. [19]. Both axes are in logarithmic scale.

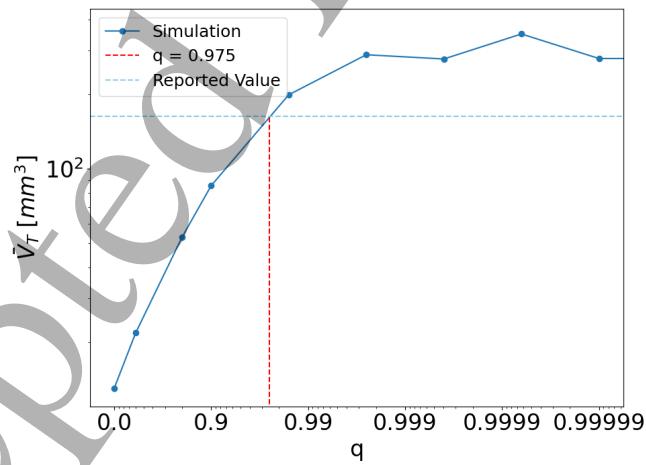


Figure 15: \bar{V}_T simulated using a grid resolution of 2 mm for the breast cancer case. The skyblue dash line shows the $\bar{V}_T = 162 \text{ mm}^3$ estimated from Stroom et al. and the red dashed line shows the value of q needed to replicate that result by doing simulations [25].

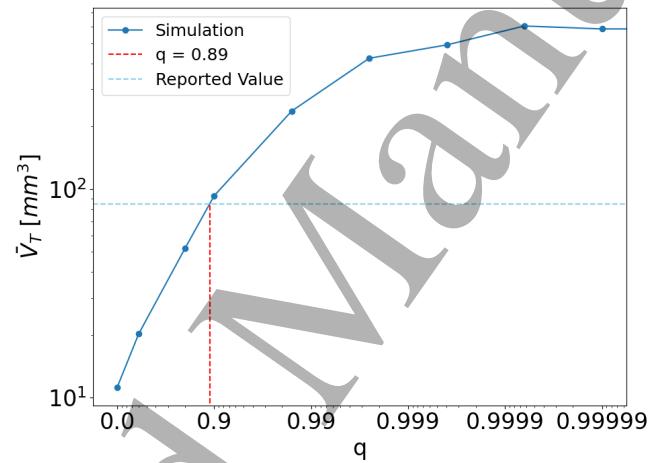


Figure 16: \bar{V}_T simulated in the lung case using a grid resolution of 2 mm. The skyblue dash line shows the $\bar{V}_T = 82 mm^3$ estimated from Stroom et al. and the red dashed line shows the value of q needed to replicate that result by doing simulations [25].