

Improving the Delivery Characteristics in Volumetric Modulated Arc Therapy (VMAT) and Tomotherapy for Cancer Treatment

by

Wilmer Henao

A dissertation submitted in partial fulfillment
of the requirements for the degree of
Doctor of Philosophy
(Industrial and Operations Engineering)
in the University of Michigan
2019

Doctoral Committee:

Professor Marina A. Epelman, Chair

Professor Amy E.M. Cohn

Associate Professor Martha M. Matuszak

Professor H. Edwin Romeijn, Georgia Institute of Technology

Wilmer E. Henao

wilmer@umich.edu

ORCID iD: 0000-0002-7441-8536

©Wilmer E. Henao 2019

Dedication

I dedicate this dissertation to my mother.

Acknowledgments

I want to sincerely thank my Ph.D. advisor Professor Marina Epelman, for helping me grow as a researcher, for the guidance, support, and trust during the past five years. I want to thank her for being an inspiration towards research and for suggesting the best ideas. I hope to one day achieve her level of mathematical vision. I would also like to thank the members of the committee, Prof. H. Edwin Romeijn, for helping me during my first year, the most difficult. I would like to thank Prof. Martha Matuszak for all her support with the data and providing me with her point of view whenever I had questions. I am grateful to Prof. Amy Cohn for letting me get to know the world of research and the world of teaching.

Special thanks to Prof. Brian Denton, who was always there with a smile, and to Professors Shen, Lavieri, Keyserling, Shi, and Jiang. They were always accommodating and supportive whenever I needed them for anything, offering tips and advice.

I would like to offer my appreciation to all the people I worked with, in particular the people at the Radiation Oncology department of the University of Michigan Hospital: Dan Polan, Carlos Anderson, Kelly Paradis. Big thanks to the people at UT Southwestern in Dallas: Weiguo Lu, Ming Li, and Mason Anders. Big thanks to Prof. Mike Overton for spending so much time working with me at NYU.

I want to thank the members of my research team, Troy Long, Ilbin Lee, and Victor Wu, who were very supporting technically as well as academically, especially during the first years when I needed a lot of their help.

A big thank you to my funding sources. I got funding through NIH P01-CA059827 and the

University of Michigan Rackham Merit Fellowship.

I'm very grateful to my dear friends in IOE, to my cohort, especially Donald, Lauren, and Nima. Friends from previous years and friends from the next. I am also grateful for the time spent with my friends in Ann Arbor: Luis Baldomero, Ivan, Luis Guevara, Abby, Nabil, Pamela. Big thanks to Sarah for never disconnecting the computer when I was running stuff.

As usual, a big thanks to Dr. Ian Malcolm, the man who got me interested in mathematics, and to the Alberto Gonzalez futbol club.

Big thanks to my brother and to my cousin Victor, for believing in me, and to my mother for making me think that everything was possible.

TABLE OF CONTENTS

Dedication	ii
Acknowledgments	iii
List of Figures	viii
List of Tables	xii
List of Abbreviations	xiii
Abstract	xv
Chapter	
1 Introduction	1
1.1 Radiation Therapy	2
1.2 Treatment Planning	3
1.3 VMAT Delivery Systems	7
1.4 Tomotherapy Delivery Systems	10
2 A Proposal for Tomotherapy Treatment Delivery and Planning Enhancement	13
2.1 Introduction	13
2.2 The Conventional Tomotherapy Delivery Approach and its Flaws	17
2.3 Optimization Models for Tomotherapy Treatment Planning	20
2.3.1 FMO-style Treatment Planning Model	21
2.3.2 New Delivery and Treatment Planning Paradigms	23
2.3.3 Detailed Model	26
2.4 Implementation	30
2.4.1 Objective Function	30
2.4.2 Solver Options	32

2.4.3	Standard Solver Parameters	33
2.4.4	Branching Priorities	33
2.4.5	The Partition Heuristic	34
2.4.6	Hints	35
2.4.7	Warm Start	35
2.5	Experiments and Results	36
2.5.1	Treatment Plan Evaluation Tools	37
2.5.2	Prostate Case Results	40
2.5.3	Running Times	50
2.6	Discussion	51
2.6.1	Modulation Factor	51
2.6.2	Resolution Increase Capabilities	53
2.7	Conclusions	54
3	VMAT with Aperture Control	57
3.1	Introduction	57
3.2	Methods	60
3.2.1	A New Aperture-Edge Penalty	60
3.2.2	Treatment Planning Problem Formulation	62
3.2.3	The Restricted Master Problem	66
3.2.4	The Pricing Problem	66
3.2.5	Solving the Pricing Problem	71
3.2.6	Aperture Selection/Refinement Heuristic for (MP)	78
3.3	Experiments and Results	83
3.3.1	Test Cases and Implementation Details	83
3.3.2	Calibration of Parameters in Metric P	85
3.3.3	Aperture Refinement	88
3.3.4	The Lung Case	90
3.3.5	The Spine Case	95
3.3.6	The Head and Neck Case	99
3.4	Multi-Arc VMAT	101
3.4.1	The Brain Case	102
3.5	Translation to the Clinical System	105
3.6	Expected Reduction of Dosimetric Discrepancies	106
3.7	Conclusions	107
4	Conclusions and Future Research Suggestions	108

4.1 Tomotherapy	108
4.2 VMAT	109
Bibliography	111

LIST OF FIGURES

Figure

1.1	An illustration of the LINAC delivery system and an MLC.	3
1.2	A patient positioned on the couch, ready for treatment. (Source: Varian Medical Systems, Inc.)	4
1.3	Discretization in treatment planning: discretization of the patient’s body and discretization of the gantry path.	5
1.4	A typical LINAC system for VMAT	8
1.5	Comparison of tomotherapy and VMAT MLCs.	9
1.6	Gantry rotation and couch movement in tomotherapy.	10
1.7	Aperture modulation with a binary colimator.	11
2.1	Gantry rotation and couch movement in tomotherapy create a helicoidal delivery pattern.	14
2.2	Tomotherapy delivery: gantry rotation and aperture modulation with a binary collimator.	15
2.3	Conventional tomotherapy treatment planning requires each leaf to open and close at every projection. Transition times in purple shown to scale.	18
2.4	An illustration of a typical term of function $F(z)$ of (2.5) that includes penalties for both over- and under-dosing a voxel in a target structure.	31
2.5	DVH plots of planned doses resulting from the FMO model for the prostate case with 51 projections per gantry rotation.	38
2.6	LOT histogram of prostate case plan with 51 projections per rotation resulting from the FMO model with $T^M = 20$ milliseconds and $T^A = 0$	41
2.7	Dose-Volume Histogram comparison of prostate case plans with 51 projections per rotation, resulting from the FMO (continuous curves) and Simple (dotted curves) models, both with $T^M = 20$ milliseconds and $T^A = 170$ milliseconds in the Simple model.	42
2.8	Leaf Control Sinogram comparison of prostate case plans with 51 projections per rotation, resulting from the Simple and Detailed models with $T^A = 170$ milliseconds and $T^M = 20$ milliseconds. Blue color represents LOTs in the Simple model, red color represents LOTs in the Detailed model, and purple color represents the regions with LOTs common to both models.	43

2.9	Dose-Volume Histogram comparison of prostate case plans with 51 projections per rotation, resulting from the Simple (continuous curves) and Detailed (dotted curves) models with $T^A = 170$ milliseconds and $T^M = 20$ milliseconds. Despite close similarities, there are, in fact, small differences between the two sets of DVH plots.	44
2.10	LOT histograms of prostate case plans with 51 projections per rotation, resulting from the Simple and Detailed models with $T^A = 170$ milliseconds and $T^M = 20$ milliseconds; (a) Simple model, (b) Detailed model.	45
2.11	DVH comparison of prostate case with 51 projections per rotation, resulting from the Detailed model with $T^M = 20$ milliseconds (dotted curve) and $T^M = 40$ milliseconds (continuous curve). Both instances used $T^A = 170$ milliseconds and were solved to achieve a 0.01% optimality gap.	47
2.12	LOT histograms of prostate case plans with 51 projections per rotation, resulting from the Detailed model with $T^A = 170$ milliseconds and (a) $T^M = 20$ milliseconds and (b) $T^M = 40$ milliseconds.	47
2.13	DVH comparison of prostate case with 153 projections per rotation, resulting from the Simple (continuous curves) and Detailed (dotted curves) models with $T^M = 20$ milliseconds and $T^A = 170$ milliseconds. Both instances were solved to achieve a 0.01% optimality gap. Some differences in the DVH plots are more prominent, especially for the rectum.	48
2.14	LOT histograms of prostate case plans with 153 projections per rotation, resulting from the Simple and Detailed models with $T^M = 20$ milliseconds and $T^A = 170$ milliseconds; (a) Simple model, (b) Detailed model.	49
2.15	Leaf Control Sinogram comparison of solutions to instances of the Simple model with $T^M = 20$ milliseconds and $T^A = 170$ with low (red) and high (blue) voxel resolutions.	52
3.1	Comparison of an irregular and a rounded aperture shape.	58
3.2	Example of an aperture illustrating vertical (stipples) and horizontal (dash-dot-dot) components of leaf edges contributing to the perimeter of the aperture.	60
3.3	The compatibility constraint (3.7d) ensures that neighboring apertures remain “reachable” given the gantry rotation speed and constraints on the speed of movement of the MLC leaves.	65
3.4	Illustration of a network constructed to solve the pricing problem for the case $M = 3$ rows in the MLC.	72
3.5	Comparison of values of $P(A)$ (horizontal axis) and $P_Y(A)$ (vertical axis) for 180 apertures in the plan obtained by applying the initial aperture selection phase of the algorithm to the lung case with $\xi = 0.75$ and $C = 0.0001$. Each dot associated with aperture A has coordinates $(P(A), P_Y(A))$. The red line represent the best linear fit for this data set, achieving $R^2 = 0.77$	87

3.6	Comparison of values of $P(A)$ (horizontal axis) and $P_Y(A)$ (vertical axis). Dots of different shapes and colors correspond to apertures in plans obtained with different values of C , as shown in the legends, after the initial aperture selection phase of the algorithm.	89
3.7	DVH plots of plans for the spine case with 36 equispaced control points around the circular arc. (a) Benchmark FMO plan; (b) Output of the initial aperture selection phase with scaling parameter $C = 0.0001$; (c) Output of several subsequent iterations of the aperture refinement phase.	91
3.8	Comparison of Dose Volume Histograms (DVHs) for the lung case plans obtained with scaling parameter $C = 0$, i.e., without aperture shape penalty (dotted lines) and with $C = 1$ (solid lines). The structures shown correspond to the target and the most important OARs.	92
3.9	The edge metric penalty $\sum_{k=1}^K P_Y(A_k)\delta_k y_k$ (top) and the quality function $F(z)$ (bottom) of treatment plans for the lung case obtained for different values of the scaling parameter C . In the top graph, the values are scaled by the edge metric penalty corresponding to the plan obtained with $C = 0$	94
3.10	The modified edge metric penalty $\sum_{k=1}^K P(A_k)\delta_k y_k$ of treatment plans for the lung case obtained for different values of C . The values are scaled by the modified edge metric penalty corresponding to the plan obtained with $C = 0$	95
3.11	Comparison of aperture shapes in lung case plans obtained using values $C = 0.0$ (left column) and $C = 1.0$ (right column) at 5 of the control points (one in each row). Beamlets are shown as squares of different colors representing the maximum dose deposition coefficient from this beamlet to any voxel in the targets using a dark-blue to bright-yellow spectrum. The darkest blue beamlets don't deliver dose to any target, and the brightest yellow beamlets have the largest dose deposition coefficient to any target voxel from that control point. The aperture contours are shown by green outlines.	96
3.12	The edge metric penalty (top) and the modified edge metric penalty (bottom) of treatment plans for the spine case obtained for different values of the scaling parameter C . The values in each graph are scaled by the value of the corresponding penalty associated with the plan obtained with $C = 0$. The horizontal axis of both plots uses a logarithmic scale.	97
3.13	Comparison of DVH plots for the spine case plans obtained for different values of the scaling parameter C	98
3.14	Comparison of DVH plots for the spine case plans obtained with scaling parameter $C = 0$, i.e., without aperture shape penalty (dotted lines) and with $C = 0.02$ (solid lines). The structures shown correspond to the target and the most important OARs.	99
3.15	Comparison of DVH plots for the head and neck case plans obtained with scaling parameter $C = 0$, i.e., without aperture shape penalty (left) and with $C = 0.5$ (right). The structures shown correspond to the targets and the most important OARs.	101

3.16	A DVH plot for the single-arc brain case plan obtained for $C = 0$, i.e., without aperture shape penalty. The structures shown correspond to the target and the most important OARs.	103
3.17	Comparison of Dose Volume Histograms (DVHs) for the brain case plans obtained with scaling parameter $C = 0$, i.e., without aperture shape penalty (dotted lines) and with $C = 0.03$ (solid lines). The structures shown correspond to the target and the most important OARs.	104

LIST OF TABLES

Table

2.1	Goals of the prostate case according to RTOG requirements.	37
2.2	Number of leaf pulsation events in the plans with 51 projections per rotation.	45
2.3	Solution times of instances of Simple and Detailed models, in cumulative seconds since the beginning of the root relaxation phase. Here, $mLOT = T^M$, $aLOT = T^A$, $v =$ number of voxels, and all times are shown in seconds	50
3.1	Summary of case sizes in VMAT experiments.	84
3.2	Treatment goals for the lung case.	92
3.3	The times (in seconds) spent in the initial aperture generation phase (first row) and in each “pass” through the aperture refinement loop (subsequent rows) for the lung case.	93
3.4	Treatment goals for the spine case.	96
3.5	Treatment goals for the head and neck case.	100
3.6	Treatment goals for the brain case.	102

LIST OF ABBREVIATIONS

AAA Analytical Anisotropic Algorithm

CCC Collapsed Cone Convolution/Superposition Algorithm

CT Computed Tomography

DICOM Digital Image and Communications in Medicine

DVH Dose Volume Histogram

DAO Direct Aperture Optimization

EM Edge Metric

FMO Fluence-Map Optimization

GTV Gross Target Volume

IMRT Intensity-Modulated Radiation Therapy

LINAC Linear Accelerator

LOT Leaf-Opening Time

MIP Mixed-Integer Programming

MLC Multi-Leaf Collimator

MRI Magnetic Resonance Imaging

MVCT Megavoltage Computed Tomography

MU Monitor Units

OAR Organ at Risk

RTOG Radiation Therapy Oncology Group

QA Quality Assurance

PTV Planning Target Volume

SBRT Stereotactic Body Radiation Therapy

TPS Treatment Planning System

VMAT Volumetric Modulated Arc Therapy

ABSTRACT

Radiation therapy treatments for cancer aim to deliver a toxic dose of radiation to malignant tumors, while controlling the dose to healthy tissues and organs. In external beam therapy, radiation is delivered by a linear accelerator, and a Multileaf Collimator (MLC) is used to change the shape of the beam opening, or aperture. A treatment plan for each patient is designed by specifying the apertures and timing and source intensity decisions for each beam angle used in the treatment. Mathematical optimization models are commonly used in modern radiation treatment planning, and rely on mathematical models of the dose distribution delivered to the patient by the portions of the beams exposed by the apertures.

Volumetric Modulated Arc Therapy (VMAT) and Tomotherapy are two of the highly utilized forms of external-beam radiation therapy, each with unique features of the MLCs. For both modalities, the current standard treatment planning methodology is prone to creating a discrepancy between the doses that are intended, and the doses that are actually delivered to the patient. The goal of our research is to develop improved treatment planning strategies that reduce this divergence.

In tomotherapy, the MLC consists of binary leaves alternating between fully open and closed states, while the beam traces a helicoidal trajectory around the patient. Dosimetric discrepancies in tomotherapy have been attributed to the lack of accurate models of leaf motion and dose delivery during leaf transitions between states, with their impact exacerbated by short leaf open times (LOTs). Moreover, the discretization of beam motion currently used for dose calculations is relatively coarse, which also contributes to the dosimetric errors. We propose a new treatment

planning delivery and modeling paradigm for tomotherapy that allows us to impose lower bounds on minimum and average LOTs, while allowing for arbitrarily finer discretization of beam motion — both features absent from existing approaches.

VMAT treatments use a rectangular MLC, with leaves that can open and close partially, creating complex two-dimensional aperture shapes, while the beam moves along pre-specified trajectories, or arcs. Current VMAT treatment planning approaches tend to create plans with complex and irregular apertures with small areas and excessive edge lengths. The dose delivery models for such apertures are less accurate than for apertures with simpler, rounder shapes; therefore discrepancies between planned and delivered dose are frequently observed in VMAT treatments as well. Our proposed optimization model for VMAT treatment planning considers a tradeoff between the quality of the planned treatment with respect to the clinical goals and a penalty on irregularly-shaped apertures. This model extends and combines a VMAT planning model with a new edge metric penalty with favorable mathematical properties compared to penalties studied in the literature. Due to the complexities of VMAT delivery, the resulting optimization models require development of heuristic solution approaches. We develop a significant extension of a heuristic algorithm for VMAT treatment planning applicable to the new model.

We test the models and algorithms proposed in this thesis on clinical cases, demonstrating improvements in treatment characteristics of the resulting treatment plans — LOTs and discretization levels in tomotherapy and aperture shape metrics in VMAT. While the precise reductions in dosimetric discrepancies resulting from these changes will need to be confirmed by dosimetric studies, these favorable characteristics should lead to clinical treatments that are more effective and safer for the patients.

CHAPTER 1

Introduction

This dissertation proposes new treatment planning methodologies for two approaches to delivery of external beam radiation for cancer treatment: Volumetric Modulated Arc Therapy (VMAT) and tomotherapy. In both delivery approaches, the current standard treatment planning methodology is prone to creating a discrepancy between the doses that are *intended*, and the doses that are actually *delivered* to the patient. The goal of our research is to develop improved treatment planning strategies that reduce this divergence.

We begin this chapter by providing an overview of radiation therapy treatments (Section 1.1) and treatment planning (Section 1.2). In the subsequent sections, we provide some additional details of VMAT and tomotherapy treatments delivery.

Following this introduction, Chapters 2 and 3 of this dissertation, dealing with our contributions to tomotherapy and VMAT treatment planning, respectively, are independent of each other. In particular, each chapter begins with a more detailed description of the radiation delivery mechanism for the specific treatment modality and its specific properties and features that can lead to the aforementioned discrepancies between planned and delivered dose distributions; we then propose new types of optimization models for treatment planning that can reduce the presence of these features. We proceed by discussing solution methods for the proposed models, and test them on clinical cases. We conclude in Chapter 4 by summarizing our contributions and proposing directions for future investigations.

1.1 Radiation Therapy

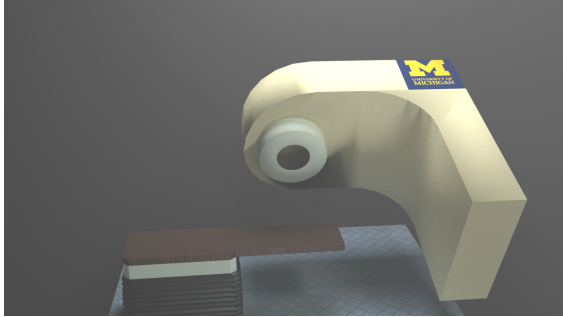
The goal of radiation therapy is to deliver a toxic dose of radiation to malignant targets, while simultaneously controlling the dose to each healthy Organ at Risk (OAR) (*Hawkins, 1994*). Healthy tissue recovers quicker than cancerous lesions, and therefore radiation oncologists break up the treatment into several *fractions*. This fractionation of the treatment allows healthy tissue the opportunity to heal and regenerate, while a toxic, destructive dose is delivered to cancerous tissue before the lesions can regrow.

Radiation therapy can be delivered either internally or externally. Internal radiation therapy is also known as *brachytherapy*.¹ It is an invasive procedure in which small radioactive sources are implanted at the treatment location or inserted via catheters. External beam photon radiation therapy is non-invasive. In this dissertation, we focus on x-ray, or photon beam, treatments involving dose delivery from a Linear Accelerator (LINAC), which distributes ionizing radiation through the relevant portion of the patient's body.² The LINAC accelerates subatomic particles and allows these particles to collide with "heavy" metals in order to produce high energy rays. These rays (beams) are formed into the desired shapes by a Multi-Leaf Collimator (MLC), which rests inside the gantry that rotates around the patient. An MLC consists of a set of pairs of "leaves" which can be moved in and out of the beam field, modulating the radiation. The opening of the beam surface created by specific leaf positions is referred to as an *aperture*. See Figure 1.1 for an illustration of the treatment gantry and an MLC.

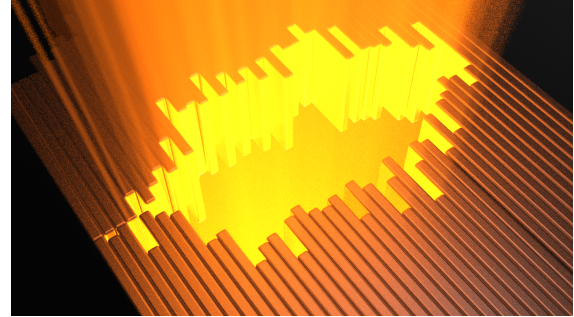
The beam is directed at the patient who is positioned on the treatment couch; see Figure 1.2. During treatment, the gantry rotates around the patient; the couch can stay stationary or it may move or rotate as well. The specifics of the movement of the gantry and the motion of the MLC leaves depend on the treatment modality. In particular, treatment machines illustrated in Figures 1.1 and 1.2 are used for Intensity-Modulated Radiation Therapy (IMRT) and VMAT. Tomotherapy is

¹*Brachys* is the Greek root for "short-distance."

²Other external beam radiation modalities include proton and gamma ray therapies.



(a) A gantry rotates around the couch where the patient is placed during treatment. The MLC is located inside the opening at the end of the gantry.



(b) An MLC is a set of leaves that close and open, i.e., move in and out of the beam opening, creating an aperture.

Figure 1.1: An illustration of the LINAC delivery system and an MLC.

a modality that involves coordinated movements of the gantry and the couch to create a helical gantry trajectory relative to the patient, and a binary MLC illustrated later on in this chapter.

1.2 Treatment Planning

A patient referred for external-beam radiation therapy begins the treatment by undergoing imaging; Computed Tomography (CT) uses x-rays to generate cross-sectional images of the body, Magnetic Resonance Imaging (MRI) uses powerful magnetic fields to achieve a similar goal. A patient can also undergo a preliminary ultrasound that uses high-frequency waves to produce images (sonograms) of the region of interest. After data acquisition, physicians use the images to identify and contour the relevant *structures*, i.e., the OARs and the targets, either manually or with the aid of software. Each target can be a visible malignant Gross Target Volume (GTV), or a Planning Target Volume (PTV), which contains the GTV and a region surrounding it to account for cancer tissue spread and geometric uncertainties. A decision is made regarding the number, individual goals, and timing of fractions, the desired doses to the targets, and safe doses to healthy structures that the radiation oncologist recommends. For example, in a spine case, the dosimetric requirements might be as follows:

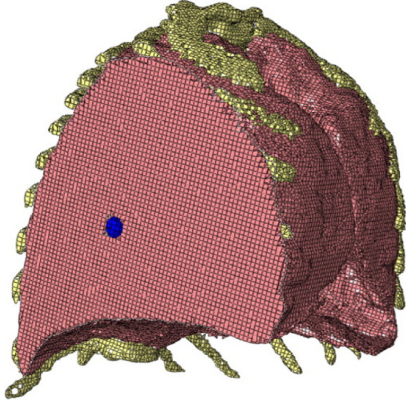


Figure 1.2: A patient positioned on the couch, ready for treatment. (Source: Varian Medical Systems, Inc.)

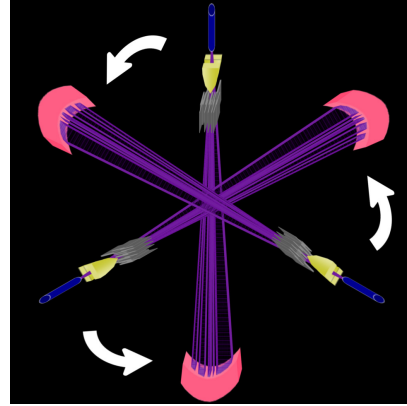
- PTV: Minimum dose ≥ 37 Gy; Maximum dose ≤ 47 Gy, D90 ≥ 44 Gy
- Spinal Cord: Maximum dose ≤ 31 Gy
- Esophagus: D50 ≤ 27 Gy
- Trachea: Maximum dose ≤ 36 Gy,

where Gray (Gy) is the unit of dose delivered, and expressions DX refer to the Dose Volume Histogram (DVH) metrics (e.g., in the above example, 90% of the PTV, by volume, should receive the dose of at least 44 Gy).

The goal of *treatment planning* is to determine a *treatment plan*, i.e., the specifics of delivery of radiation to the particular patient, that adheres to the above treatment goals and is consistent with the capabilities of the treatment modality and equipment used. Depending on the type of cancer and the treatment modality, several particular challenges may arise during treatment planning: accounting for organ function information, uncertainties in inter-fraction motion, changes



(a) A discretized organ at risk: A 3D model of a lung. (Source: *Marias et al., 2011.*)



(b) The gantry path is discretized into control points (VMAT) or projections (tomotherapy).

Figure 1.3: Discretization in treatment planning: discretization of the patient’s body and discretization of the gantry path.

in patient geometry between fractions, intra-fraction motion such as breathing, natural tissue inhomogeneities such as denser bony regions, or porous sectors in the lung, and ensuring sufficient accuracy of dose calculations during the planning process. It is the latter that we strive to address and improve in this dissertation by proposing new modeling approaches for optimization-based treatment planning. All of our proposed models address planning of the total treatment; we assume that an appropriate fractionation schedule will be determined as a post-processing step.

To estimate the delivered dose during the treatment planning process, the treatment volume of the patient’s body is discretized into small 3-dimensional *voxels* (sometimes also referred to as “points”). Each of these voxels is contained in one or more structures. The typical voxel size is about 1 or 2 cubic millimeters for cases that require extreme precision, such as brain tumors, and up to 1 cubic centimeter for simpler cases, such as prostate cancer. We show a schematic of a discretized lung in Figure 1.3a. The dose delivered to the patient is represented as a vector of doses received by the voxels, each of which is calculated as a sum of contributions to the dose from the exposed portions of the beam, with the gantry positioned at different angles with respect to the patient.

The nature of the gantry movements is dictated by the treatment modality. The more traditional

IMRT approach, for example, uses a limited number of static beam angles. The gantry remains stationary at each beam angle while the MLC leaves are repositioned to *modulate* the beam intensity by exposing different apertures for different amounts of time. A common treatment planning technique for IMRT is the so-called Fluence-Map Optimization (FMO), in which the surface of the beam is discretized into rectangular *beamlets*, and the intensity (via duration of exposure) of each beamlet at each angle is determined independently. After FMO planning is performed, a leaf-sequencing algorithm identifies several apertures that, when exposed one after the other, achieve the different beamlet intensities. (An alternative approach to IMRT treatment planning is Direct Aperture Optimization (DAO), which, as the name suggests, directly chooses specific apertures and corresponding intensities.) The width of the beamlets corresponds to the width of the MLC leaves, and the length of the beamlets is chosen based on the desired discretization density.

In IMRT, the dose is calculated as the sum of contributions of individual beamlets from each beam angle. In FMO, for each beamlet-voxel pair at a beam angle, a dose deposition coefficient, which is the dose that is deposited in the voxel from this beamlet at unit intensity, is pre-calculated using advanced Monte Carlo simulation of particles at the sub-atomic level. The contribution of the beamlet to the voxel's dose is then proportional to the total time this beamlet is exposed by the MLC, with the dose deposition coefficient serving as the coefficient of proportionality. In FMO and other aperture-based treatment planning approaches, the dose deposition coefficient for an aperture-voxel pair is often approximated by the sum of coefficients of beamlets exposed within the aperture. Although more direct, aperture-specific, methods based on Monte Carlo simulation also exist, their use during treatment planning is often too computationally demanding, and instead they are used as a Quality Assurance (QA) step afterwards.

In practice, IMRT has limitations, including long delivery times necessitated by the large number of apertures needed to produce complex fluence profiles at each beam angle, and high Monitor Units (MU) (a measure of the total output of the machine during treatment). This concern can be addressed in VMAT and tomotherapy treatments, where, unlike in IMRT, the gantry continuously

moves relative to the patient, and the MLC leaves change positions while the gantry is in motion. For the purposes of treatment planning, the path of the gantry around the patient is also discretized, as illustrated in Figure 1.3b. In VMAT, the discretization points are referred to as “control points,” while in tomotherapy they are typically called “projections.” (Following the conventions of each treatment modality, we will use the terms beam angles, control points, and projections when discussing IMRT, VMAT, and tomotherapy treatments, respectively.) A common discretization-based dose calculation approach uses the contribution to the dose from a stationary aperture at an individual control point as an approximation of the contribution made as the gantry moves from this control point to the next; as long as the discretization of the gantry trajectory is sufficiently fine, this provides a sufficiently accurate approximation.

Radiation treatment planning via mathematical optimization has gained prominence both in research and clinical practice, and some forms of it have been adopted by commercial treatment planning systems. There are many optimization modeling approaches proposed for treatment planning, chosen based on the unique characteristics of each treatment modality and the available software capabilities. A survey by *Romeijn and Dempsey (2008)* provides a comprehensive, if slightly outdated, overview of optimization models used in IMRT planning; we will review the relevant literature for VMAT and tomotherapy planning in the forthcoming sections.

In the following sections, we provide some additional details of delivery paradigms and treatment planning in VMAT and tomotherapy. We begin by discussing VMAT in Section 1.3, due to its similarities to the concepts discussed so far, and address tomotherapy in Section 1.4.

1.3 VMAT Delivery Systems

VMAT radiation therapy delivery system was initially proposed by *Yu (1995)* under the moniker of Intensity-Modulated Arc Therapy, or IMAT. VMAT is delivered by LINAC systems similar to those used in IMRT (in fact, several delivery systems on the market are capable of delivering both

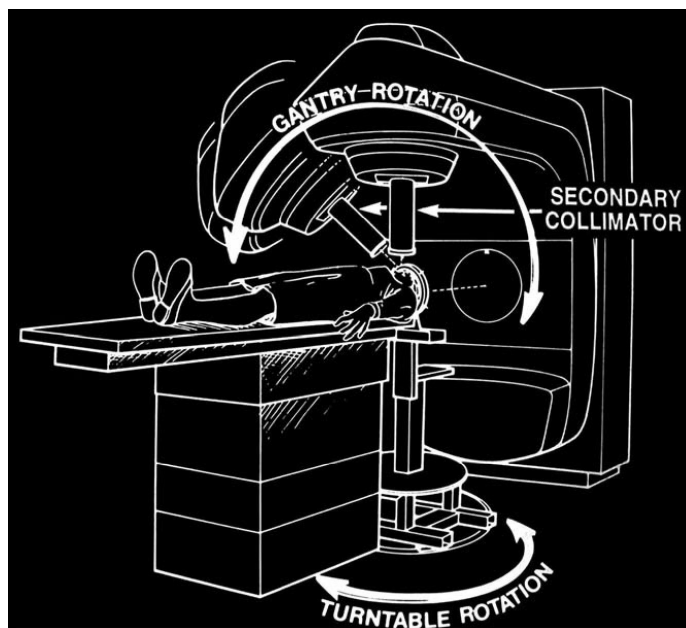
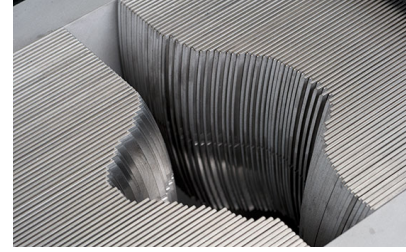


Figure 1.4: A typical LINAC system for VMAT

IMRT and VMAT treatments), but unlike the “step-and-shoot” IMRT delivery from a small number of stationary gantry angles, VMAT treatments are delivered by a continuously rotating gantry. The gantry is equipped with a rectangular MLC that can be used to change the beam aperture dynamically. The leaves can move in and out of the beam opening and are capable of creating complex two-dimensional shapes (see Figure 1.5b for another illustration). The high precision of continuously moving MLC leaves combined with continuous delivery allows for highly conformal treatments with faster treatment times than IMRT, allowing VMAT treatments to be used in the most complex cases. The VMAT couch lacks horizontal displacement, but can rotate as shown in Figure 1.4, or can be stationary for a *coplanar* treatment. As opposed to the helical trajectory outlined by the tomotherapy LINAC, the VMAT gantry moves along a trajectory lying on a sphere around the patient, centered on the target. VMAT allows the delivery from non-coplanar angles, which is an advantage over tomotherapy (see Section 1.4). We will primarily focus on coplanar single-arcs VMAT treatments, where the couch remains stationary and the gantry rotates, tracing a circular trajectory, or arc, around the patient one or multiple times. Multi-arc treatments, with



(a) Binary Multi-Leaf Collimator (tomotherapy); each leaf is either completely open or closed.



(b) Multi-Leaf Collimator (VMAT); leaves can create complex shapes.

Figure 1.5: Comparison of tomotherapy and VMAT MLCs.

a couch in a different position for each arc, are conceptually similar, and will be explored in Section 3.4.

As we already mentioned, for treatment planning purposes, the trajectory of the gantry in a VMAT treatment is discretized into control points. A treatment plan is typically given by specifying an aperture (i.e., an MLC configuration), intensity of radiation source, and gantry rotation speed (equivalently, exposure time) at each control point. The delivery machine then uses interpolation to specify behavior along the continuous treatment path.

Several optimization algorithms for VMAT treatment planning have been developed, mostly considering one of two general approaches. The first approach is to use two-stage arc-based models. In the first stage, the “ideal” beam profiles for all control points are determined, often using an IMRT-style fluence-map optimization model. In the second stage, arc-sequencing (a process similar to leaf-sequencing in IMRT, but applied to an arc) takes place in order to construct a deliverable VMAT treatment that approximates the “ideal” treatment obtained in stage one. *Cao et al. (2009)*, *Craft et al. (2012)*, *Wala et al. (2012)*, *Salari et al. (2012)*, and *Papp and Unkelbach (2014)* create different algorithms to perform this second step. The second approach to VMAT treatment planning is to use direct control point-based decisions, in which each aperture at every control point is designed to take into account the particular LINAC constraints. One of these approaches uses column generation-type heuristics that add or replace an aperture at a different control point in every iteration (see, e.g., *Men et al., 2010*; *Peng et al., 2012, 2015*). Another approach works on

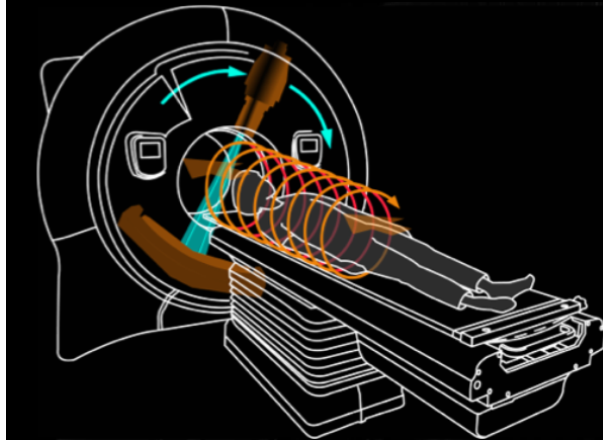


Figure 1.6: Gantry rotation and couch movement in tomotherapy.

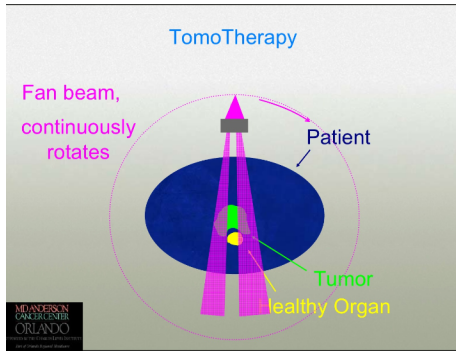
multiple control points simultaneously; this leads to a column-and-row generation approach such as *Mahnam et al. (2017)*. Our algorithmic approach presented in Chapter 3 generalizes and refines the column generation heuristic of *Peng et al. (2012)*.

1.4 Tomotherapy Delivery Systems

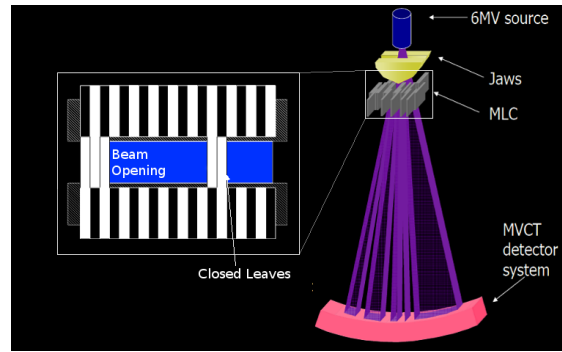
Tomotherapy combines the accuracy of computed tomography with the power of intensity-modulated radiation therapy. It shares the continuous gantry motion paradigm of VMAT, but has many unique characteristics. During treatment, the gantry rotates around a fixed axis, while the couch moves along a straight line perpendicular to the gantry's plane of movement. As a result, tomotherapy displays a characteristic helicoidal delivery pattern shown in Figure 1.6.

Tomotherapy can be a very accurate treatment modality that allows the treatment of large regions, or regions with several lesions. Its computed tomography capabilities allow the verification of the location of each tumor before treatment, reducing geometric uncertainties in treatment delivery.

Tomotherapy is affected by steep gradients in dose delivery, which complicates calculations of dose deposition coefficients. Another difference related to the beam is the source-to-axis distance,



(a) Aperture from Front of LINAC.



(b) Binary MLC shapes the “fan” beam.

Figure 1.7: Aperture modulation with a binary colimator.

at 85 cm, shorter than for IMRT and VMAT (usually 100 cm). Other distinctive features of tomotherapy include the absence of a flattening filter³ or a beam hardener,⁴ and lack of an electron stopper. Together, these features deliver higher energy production at the center of the beam. The radiation beam in tomotherapy is therefore significantly different from other treatment modalities in which flat isodoses are achieved. The features of the beam, however, have been studied extensively, among others, by *Jeraj et al. (2004)*; *Sterpin et al. (2008)*; *Seco and Verhaegen (2013)*, and we can assume that Monte Carlo simulations that model the beam fairly accurately are available.

Instead of rectangular beams and two-dimensional apertures used in IMRT and VMAT treatments, tomotherapy uses a narrow beam and a *binary* MLC modulating it. Figure 1.5 shows an illustration of the latter and provides a side-by-side comparison with the former. Figure 1.7 provides a more detailed schematic depiction of modulation of the beam by a binary MLC. As the gantry traces the helicoidal trajectory around the patient, each leaf is dynamically opened and closed. The name “binary” stems from the fact that a leaf cannot be opened or closed part-way; it can only be in one of two states, open or closed (i.e., on or off). In standard models of tomotherapy delivery, and in treatment planning approaches, changes in leaf state are assumed to happen instan-

³The flattening filter is a conic piece of equipment that flattens the dose uniformly across the whole aperture. Tomotherapy does not require a uniform beam intensity because the collimators can achieve the desired uniformity. Moreover, by not using a flattening filter, medical physicists and dosimetrists can engineer a more homogeneous energy spectrum, without any loss of intensity in the process.

⁴The beam hardener filters out lower intensities of the beam, leaving only the “hardest” components.

taneously. As we will discuss in Section 2.2, this assumption leads to dosimetric inaccuracies in tomotherapy treatment planning.

There are certain additional inherent disadvantages to tomotherapy, including an excessive number of leaf pulsations (i.e., changes in leaf status during treatment), the requirement of extra beam-on time, short leaf-opening times, and more extended delivery periods that cause the machine components to wear out (*Kampfer et al., 2011*). We believe that our approach to modeling and delivering tomotherapy treatment introduced in Chapter 2 has the potential to tackle many of these problems.

CHAPTER 2

A Proposal for Tomotherapy Treatment Delivery and Planning Enhancement

2.1 Introduction

Tomotherapy (literally, “slice therapy,” *Mackie et al., 1993*) is a technique for delivering external beam radiation that combines the power of intensity-modulated radiation therapy and the accuracy of computerized tomography, which allows concurrent accurate tomographic setup verifications. Tomotherapy enables the treatment of a wide target area due to the axial movement of the couch; this characteristic makes it ideal for the treatment of either larger targets or more significant ministrations that contain several targets.

During tomotherapy treatment, a ring gantry rotates around a fixed axis at a speed of 1 to 10 revolutions per minute (*Webb, 2001*), while the couch moves along the straight line perpendicular to the gantry’s plane of movement. Originally intended to remove the possibility of collision between the patient and the treatment unit (*Mackie et al., 1993*), tomotherapy displays the characteristic helicoidal delivery pattern shown in Figure 2.1. The beam used in tomotherapy is narrow, and it is modulated with a pneumatic *binary* Multi-Leaf Collimator, in which each leaf can be in one of two states, open or closed (i.e., on or off). Thus, in tomotherapy, beamlets have a length equal to the length of the leaves. The collimator usually consists of 64 leaves (each leaf is usually 6 mm wide

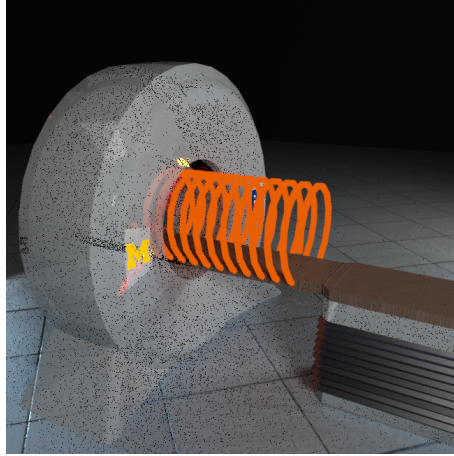
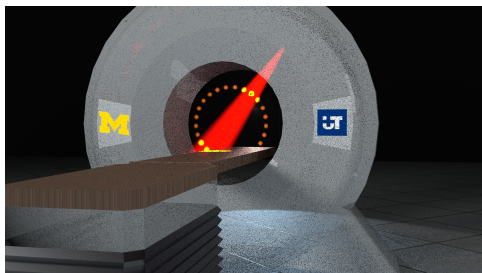


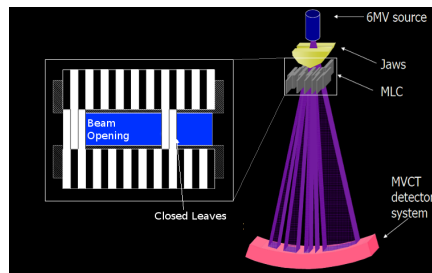
Figure 2.1: Gantry rotation and couch movement in tomotherapy create a helicoidal delivery pattern.

and anywhere from 10 to 50 mm long). While earlier versions of tomotherapy equipment only allowed leaf lengths of 1, 2.5, or 5 cm that were fixed at the start of each treatment, the latest version known as “TomoEdge” allows a dynamical change of leaf length during treatment; however, this feature is beyond the scope of our analysis, which focuses on the more traditional delivery systems. The gantry and MLC mechanism are illustrated in Figures 2.2a and 2.2b. To fix ideas, with respect to the 3D coordinate system, we say that the gantry’s plane of movement, or *the treatment plane*, is the $X-Z$ plane orthogonal to the direction of couch movement, Y . The MLC is aligned parallel to the treatment plane so that the MLC leaves travel in the Y -direction when they open and close. We also define the *slice width* as the longitudinal extent (i.e., in the Y -direction) of the portion of the treated area of the patient’s body covered by the beam emanating from the MLC with the gantry in a fixed position and the leaves opened (*Langen et al., 2010*).

A final feature in tomotherapy modeling is the *pitch*: the ratio of couch translation per rotation to the slice width. The study of *Kissick et al. (2005)* recommends pitch values in the set $\frac{0.86}{N}$ for some $N \in \mathbb{N}$, where lower values are better, since, according to *Westerly et al. (2009)*, “conventional thinking suggests that increasing the pitch may result in a loss of longitudinal resolution in the dose distribution.” A pitch taking any of the aforementioned values muffles the appearance of the



(a) Each rotation of the gantry is discretized into 51 projections.



(b) Binary MLC shapes the “fan” beam.

Figure 2.2: Tomotherapy delivery: gantry rotation and aperture modulation with a binary collimator.

so-called *thread effect*, which refers to a disadvantageous helicoidal thread-like pattern in dose delivery attributed to the mismatch of helical beams. The values $\frac{0.86}{N}$ have been further confirmed by [Chen et al. \(2011\)](#), emphasizing that correct optimization makes the thread effect disappear.

A first sketch of a practical design of the tomotherapy machine can be traced back to [Mackie et al. \(1993\)](#). It took seven years until the first tomotherapy LINAC prototype entered the clinical setting in the early 2000’s ([Sheng, 2017](#)). [Shepard et al. \(2000\)](#) suggested iterative methodologies for treatment planning optimization, and [Lu \(2010\)](#) proposed a non-voxel broad-beam (NVBB) IMRT approach.

As in other treatment modalities, traditional tomotherapy treatment planning models use the discretization of the patient’s structures into voxels. A precursor of those methodologies is the treatment in [Mackie et al. \(1985\)](#). According to [Grigorov et al. \(2003\)](#), conventional tomotherapy delivery methodologies today rely on the discretization of the helicoidal trajectory into projections. As mentioned in Chapter 1, a projection is the tomotherapy-specific term equivalent to a control point or a beam angle. In conventional approaches to tomotherapy, it is common to divide the trajectory into 51 projections per 360° rotation. Each projection has an associated set of binary leaves, and for each leaf-projection combination (i.e., beamlet), there is a specification of leaf opening and closing times. Furthermore, the projection encompasses an associated collection of pre-calculated coefficients that indicate how much dose is deposited from each leaf to each voxel;

these dose coefficients are assumed to remain constant throughout the projection. The computation of these coefficients uses the nuclear physical properties of the beam, and Monte-Carlo simulations supply these coefficient values (*Seco and Verhaegen, 2013*).

It should be noted that physical limitations of the LINAC require all leaves to remain closed for a few seconds at the beginning of the helicoidal treatment. According to *Langen et al. (2010)*, the leaves must be closed for 10 seconds at the beginning of each treatment until the LINAC output stabilizes. This can be accomplished in the treatment planning models by initiating the indexing of the projections at the first discretization point where the leaves are allowed to open, 10 seconds into the gantry's rotation.

Tomotherapy treatments have certain inherent disadvantages, including the inability to deliver from all non-coplanar directions, an excessive number of *leaf pulsations*, i.e., leaf openings and closings, the longer beam-on time, and overall wearing out of the systems (*Kampfer et al., 2011*), in spite of newer in-house LINACs that slightly reduce the wear-and-tear, by virtue of a reduction of voltage output from 6 MV to 5.2-5.7 MV (*Sheng, 2017*).¹ Standard models of tomotherapy delivery, and thus treatment planning approaches, *assume* that changes in leaf state happen instantaneously. However, this assumption is not accurate, as opening and closing each leaf pair takes about 20 milliseconds, which results in inaccurate estimates of delivered doses. We believe there is room for improvement, and our approach to modeling and delivering tomotherapy treatments has the potential to tackle some of these problems. We will expand on some of these issues in the next section, and in the following sections we propose alternative tomotherapy treatment planning approaches based on mixed integer optimization models.

¹Further reductions seem attractive but are not possible because “energies lower than 6 MV have low x-ray production per incident electron” (*Mackie et al., 1993*).

2.2 The Conventional Tomotherapy Delivery Approach and its Flaws

In the conventional tomotherapy delivery, the LINAC is set to produce a constant fluence rate of 850 cGy per min (*Sheng, 2017*). The modulation at each projection is accomplished by opening each leaf of the MLC for varying amounts of time as the gantry passes through the 7.06° ($7.06^\circ \approx \frac{360^\circ}{51}$) arc segment associated with the projection; see Figure 2.3. Thus, mathematically, treatment planning essentially uses a Fluence Map Optimization model, i.e., by controlling leaf opening times, different fluences for beamlets within each projection can be achieved. The resulting treatment is delivered at the slowest necessary constant pre-determined gantry rotation speed. The standard tomotherapy paradigm *requires* that a leaf that opens is closed during that same projection. When a centered target, large or small, is being treated, leaf-opening events at neighboring projections can be merged, e.g., the opening of a leaf at projection 1 can be timed to occur towards the end of the corresponding arc segment, while the opening of the same leaf at projection 2 can be timed to occur towards the beginning of the following arc segment; these can be combined into a single leaf opening per two segments in order to increase the overall leaf-opening-time and reduce the number of leaf pulsations. However, such merges are not recommended when treating small off-centered targets due to the “blurring effect” (*Sheng, 2017*).

The conventional approach to both treatment delivery and planning endures some limitations. As mentioned above, in the conventional approach, each leaf that opens has to close within that very same projection (unless the aforementioned pairwise merging is deployed, in which case it has to close within the next projection). As mentioned earlier, a factor contributing to the discrepancy in dose calculation is leaf speed. Despite of the commonly-made assumption to the contrary, the pneumatic mechanism that closes and opens each leaf does not change the leaf state instantaneously. As a matter of fact, each *leaf event* takes roughly 20 msec, which means that it takes a total of 40 msec to open and close a leaf in each projection. It is challenging to properly account



Figure 2.3: Conventional tomotherapy treatment planning requires each leaf to open and close at every projection. Transition times in purple shown to scale.

for dose delivery during this time, which can have a significant impact on the overall accuracy of dose calculations. Typically, each gantry rotation takes 15 seconds, which means it is spending roughly 300 milliseconds on each projection. Thus, 40 milliseconds out of 300 milliseconds, or roughly 13.3% of the time, is spent in a state that's difficult to model. If we compare the time a leaf spends in transition to the time it stays open, the ratio is even larger, especially for shorter opening times. This uncertainty could potentially translate into a discrepancy in computed dose vs. actual dose that is more than 3%, which is the acceptance criterion used for QA according to [Westerly et al. \(2009\)](#). The TomoTherapy Hi-Art II Treatment Planning System (TPS) partially accounts for this effect, but its calculations are based on assumptions of linearity (i.e., constant velocity) of leaf motion and uniformity of all leaves; neither of these assumptions hold in practice (see, for example, Figure 4 of [Westerly et al., 2009](#) for an illustration of deviations from the linear model of various leaves in the MLC of a test tomotherapy machine). In the same paper, the authors measure a discrepancy between the theoretical dose and the dose that gets delivered to several patients.² Using 3D-Megavoltage Computed Tomography (MVCT) imaging, they determine that the source of the discrepancy is *the short lengths* of time that the leaves stay open. Short Leaf-Opening Time (LOT) tends to aggravate this discrepancy for reasons mentioned above.

²The authors measured discrepancy on a water phantom with the use of ion chambers.

Note that, according to the recommended protocols from Task Group 148 (*Langen et al., 2010*), namely, “leaf opening times shorter than 20 ms are deleted from the control sinogram³ since they are too small in relationship to the actual leaf transition times,” dosimetrists must check the new dose distribution that removes the short leaf opening times, and verify if the treatment should still be approved.

In order to control the dose discrepancies, the authors in *Westerly et al. (2009)* propose having longer LOTs (if possible, longer than 100 milliseconds); the authors rely on an increase in the pitch to achieve that goal.⁴ Increasing the pitch makes the arc length of the helix described by the gantry around the patient shorter, forcing the machine to deliver more dose at every projection, which, in turn, is accomplished by keeping the leaves open for longer periods of time. Unfortunately, increasing the pitch decreases dose homogeneity (*Langen et al., 2010*) and causes loss of longitudinal resolution. Furthermore, when the dose per fraction is higher than 2 Gy, it is recommended to reduce the pitch below 0.2. This is well below the 0.287 value proposed by *Westerly et al. (2009)*, although it is partially justified by the results from *Gutiérrez et al. (2007)* which claims that increasing the pitch does not have a significant effect if the slice width is smaller than 1 cm.

Besides dose discrepancies, there exists a concern for the excessive wear-and-tear of the machine due to the numerous leaf-state changes in treatments where pulsations occur at almost every single projection. The original sketches of *Mackie et al. (1993)* estimated mean time between failures of the collimator to be about twelve months. In recent years, the wear-and-tear has been shown to decrease the expected lifetime of the linear accelerator (*Kampfer et al., 2011*). If it is possible to reduce the number of pulsations, we would be able to increase the reliability of the collimator, the lifetime of the machine, and the overall safety of the radiation treatment process. We can indirectly achieve this by increasing the average LOT, since the number of pulsations is in the denominator of this indicator. Another concern has to do with case resolution: current methodology pegs the

³See Section 2.5.1.2 for an explanation of the term.

⁴Recall that the pitch is defined as the ratio of the couch translation per rotation to the slice width.

number of projections at 51 per rotation cycle, or approximately 7° , leading to the aforementioned beamlet blurring effect. Increasing the resolution of projections per rotation would allow for more precise calculation of the dose deposition coefficients, but doing so under the prevailing paradigm of a full leaf pulsation in each projection where the leaf opens would result in an increase in the number of pulsations and shorter LOTs. A model that allows a refinement of the resolution without decreasing the length of leaf-opening events would be preferred.

We think it is essential for the practitioners to have the freedom to choose whatever pitch is considered appropriate according to the particular case recommendations. Moreover, the impact of pitch increase on LOTs is only indirect and is hard to predict. In the forthcoming proposed treatment planning models, we use integer programming techniques to explicitly control leaf events and constrain LOTs, while leaving the door open to increases in the resolution of gantry rotation discretizations without the negative consequences discussed above.

2.3 Optimization Models for Tomotherapy Treatment Planning

In this section, we present several optimization models for tomotherapy treatment planning. We first present a model reflecting the traditional treatment paradigm, and then propose two models applicable if a new delivery paradigm is adopted.

All of our optimization models share some common notation. In all of them, we use the following parameters and variables (additional notation for each model will be introduced separately):

Parameters:

- P — the number of projections in the gantry trajectory; we denote $[P] = 1, \dots, P$ and assume projection $p = 1$ is the first projection along the gantry's trajectory where the leaves can be open
- L — the number of leaves in the binary MLC; we denote $[L] = 1, \dots, L$

- \mathcal{V} — the set of voxels
- \bar{t} — time for the gantry to traverse each projection (we assume gantry rotation speed and number of projections per rotation has been fixed, and \bar{t} has been pre-calculated based on those settings)
- D_{lvp} — dose deposition coefficient for leaf l at projection p and voxel v (dose delivered per millisecond of leaf opening)
- T^A — lower bound on the average LOT
- T^M — lower bound on the minimum LOT

Variables:

- z_v — dose delivered to voxel v
- t_{lp} — time leaf l stays open at projection p
- $\beta_{lp} = 1$ if leaf l is open for all or part of projection p ; 0 otherwise

In all of our models, we will calculate the dose to each voxel as a sum of contributions from all leaf openings at all control points (*Gibbons et al., 2009*). Moreover, all of our optimization models will use the same objective function, $F(z)$, which evaluates treatment quality based on the dose distribution vector $z \in \mathbb{R}_+^{|\mathcal{V}|}$.

2.3.1 FMO-style Treatment Planning Model

In the FMO-style treatment planning model, suitable for the traditional delivery paradigm, we determine the opening time of each leaf at each projection independently. The model is as follows:

$$\underset{t_{lp}, z_v, \beta_{lp}}{\text{minimize}} F(z) \quad (2.1a)$$

$$\text{subject to } z_v = \sum_{l \in [L]} \sum_{p \in [P]} D_{lp} t_{lp}, \quad v \in \mathcal{V} \quad (2.1b)$$

$$T^M \beta_{lp} \leq t_{lp} \leq \bar{t} \beta_{lp}, \quad l \in [L], \quad p \in [P] \quad (2.1c)$$

$$\beta_{lp} \text{ binary}, \quad l \in [L], \quad p \in [P].$$

Here, constraints (2.1b) define doses as the sums of contributions from individual leafs at various projections. Each constraint in (2.1c) connects the values of t_{lp} and β_{lp} and ensures that the time a leaf stays open at a projection, if it is positive, does not exceed the time spent at the projection and is no shorter than the required per-projection minimum, T^M . If no such minimum is imposed by the planning system, i.e., $T^M = 0$, then we can eliminate all β variables from the model by setting them to 1, and the resulting model will have linear inequality and equality constraints in continuous variables only — this is the simplest of the models we consider. On the other hand, in the spirit of [Langen et al. \(2010\)](#), it may be desirable to impose a lower bound of $T^M = 20$ milliseconds on positive leaf opening time at each projection, to ensure that a leaf opening included in the proposed treatment plan does not get eliminated by the tomotherapy delivery system. The resulting mixed integer programming problem represents a slight improvement of the simplest FMO-style treatment planning model. If desired, this model can be augmented to include a lower bound on the average LOT:

$$\sum_{l \in [L]} \sum_{p \in [P]} t_{lp} \geq T^A \left(\sum_{l \in [L]} \sum_{p \in [P]} \beta_{lp} \right). \quad (2.2)$$

2.3.2 New Delivery and Treatment Planning Paradigms

Westerly et al. (2009) suggests that increasing the average LOT of a treatment is a good approach for dose discrepancy control. While the average LOT may be a good proxy for individual LOTs, it does not account for the entire distribution of LOTs used in the treatment. For example, it is entirely possible that only a few LOTs become longer while the rest of the LOTs remain short. In our new proposed tomotherapy treatment planning models, we can incorporate lower bounds on minimum and/or average LOTs. By explicitly controlling these delivery metrics, we can improve delivery characteristics of the planned treatment.

To make use of our proposed treatment planning models, a new tomotherapy delivery paradigm needs to be adopted. In particular, in our models we do away with the assumption that a leaf opened at a projection must be closed within the same projection. In the first of our models (which we refer to as the “simple model”), we explicitly assume that merging of leaf opening times at the boundary of adjacent odd-even projection pairs occurs, as discussed in Section 2.2. The lower bounds on minimum and average LOTs in the context of this model cannot exceed $2\bar{t}$, i.e., the combined time of two adjacent projections.

The latter consideration is a limitation of this model, especially if a significant increase in minimum or average LOTs is desired, and/or if a finer discretization of the gantry trajectory is being used (thus decreasing the value of \bar{t}). Therefore, we propose a second treatment paradigm, in which a leaf, once opened, can stay open for an arbitrary number of projections within the gantry trajectory, and an optimization model for treatment planning in this context, which we call the “detailed model.”

Each model uses different additional binary variables to keep track of projections in which the leaves open and/or close.

2.3.2.1 Simple Model: Odd-Even Projection Pairing

To review, the model in this section is built upon the following assumptions:

1. If a leaf is open for all or part of projection p , we assume the opening occurs at the end of the projection if p is odd, and at the beginning of the projection if p is even.
2. The duration of two projections exceeds the lower bound on the minimum LOT, i.e., $2\bar{t} \geq T^M$. This implies that minimum LOT constraints can be satisfied by $t_{lp} + t_{lp+1}$.
3. The duration of two projections also exceeds the lower bound on average LOT, i.e., $2\bar{t} \geq T^A$, for the same reason.

For simplicity of presentation, we will assume that the overall number of projections, P , is even.

The model uses the following additional set of variables:

- $\gamma_{lp} = 1$ if $\beta_{lp} = 1$, or $\beta_{lp+1} = 1$, or both; 0 otherwise, for p odd,

and is given by:

$$\begin{array}{l} \text{minimize } F(z) \\ \quad t_{lp}, z_v \\ \quad \gamma_{lp}, \beta_{lp} \end{array} \tag{2.3a}$$

$$\text{subject to } z_v = \sum_{l \in [L]} \sum_{p \in [P]} D_{lp} t_{lp}, \quad v \in \mathcal{V} \tag{2.3b}$$

$$0 \leq t_{lp} \leq \bar{t} \beta_{lp}, \quad l \in [L], p \in [P] \tag{2.3c}$$

$$\gamma_{lp} \leq \beta_{lp} + \beta_{lp+1} \leq 2\gamma_{lp}, \quad l \in [L], p \in [P] \text{ is odd} \tag{2.3d}$$

$$t_{lp} + t_{lp+1} \geq T^M \gamma_{lp}, \quad l \in [L], p \in [P] \text{ is odd} \tag{2.3e}$$

$$\sum_{l \in [L]} \sum_{p \in [P]} t_{lp} \geq T^A \left(\sum_{l \in [L]} \sum_{p \in [P] \text{ is odd}} \gamma_{lp} \right) \tag{2.3f}$$

$$\beta_{lp}, \gamma_{lp} \text{ binary}, \quad l \in [L], p \in [P].$$

Here, constraints (2.3c) ensure that the time a leaf stays open at a projection is nonnegative and does not exceed the time spent at the projection, and is positive only if $\beta_{lp} = 1$. Constraints (2.3d) connect the values of β 's and γ 's. Constraints (2.3e) ensure that the combined LOT of an odd-even projection pair meets the corresponding lower bound whenever it is nonzero. Constraint (2.3f) ensures that the average LOT satisfies the corresponding lower bound.

Note that, in the above model, the lower bound on the LOT is imposed on $t_{lp} + t_{lp+1}$, where p is odd (assuming one or both of the values of t are positive), and the average LOT is also calculated based on the above “paired” values of individual LOTs. In particular, each combination of an odd-even projection pair where the leaf is open in one or both projections is considered to be a separate leaf opening. These calculations will not be accurate in some situations where a leaf’s open time is equal to \bar{t} , i.e., the leaf remains open for an entire projection. We expect this model to perform reasonably well if \bar{t} is large relative to T^M and T^A , to the extent that makes it unlikely that many leaves will be open for entire projections. After a solution to the treatment planning optimization problem is obtained, accurate LOT statistics should be computed based on actual values of t , rather than based on projection pairings.

Another observation is that the above model does not enforce the condition that $\beta_{lp} = 1$ *only* if $t_{lp} > 0$. Indeed, the combination $\beta_{lp} = 1$ and $t_{lp} = 0$ is allowed by the model, but this is without loss of generality or optimality. To see this, suppose a solution (t, β, γ) satisfies constraints (2.3), and consider another solution, $(t, \tilde{\beta}, \tilde{\gamma})$, which has been updated by setting $\tilde{\beta}_{lp} := 1$ if and only if $t_{lp} = 1$ and $\tilde{\gamma}_{lp} := 1$ if and only if $\tilde{\beta}_{lp} = 1$ and/or $\tilde{\beta}_{lp+1} = 1$ when p is odd. This new solution satisfies constraints (2.3c) and (2.3d) by construction. Notice also that $\tilde{\beta} \leq \beta$ and $\tilde{\gamma} \leq \gamma$, and so the right-hand sides of constraints (2.3e) and (2.3f) may decrease after this update, while their left-hand sides will stay the same, i.e., these constraints will also be satisfied by the updated solution. Moreover, since the values of t are unchanged, the dose distribution resulting from this plan, and its objective value, remain the same.

2.3.3 Detailed Model

In this model, we will use additional binary decision variables b_{lp} , m_{lp} , e_{lp} , along with β_{lp} , to provide a more precise calculation of LOTs, as well as to allow for direct optimization of the timing of leaf opening within a projection. In this model, we use the following conventions and make the following assumptions:

1. In every projection, the leaf is either closed, i.e., $t_{lp} = 0$, or $t_{lp} > 0$ and one and only one of the variables b_{lp} , m_{lp} , e_{lp} is equal to 1 (we refer to the latter cases as projection having *type* b , m , or e).
2. If $m_{lp} = 1$, then the leaf is open for the entire projection (the converse, however, need not be true — see below).
3. In the projection where a leaf opens, we assume that it opens towards the end of the projection and set $e_{lp} = 1$; it is possible that the leaf opens at the very beginning of the projection, i.e., $t_{lp} = \bar{t}$. This projection might (or might not) be followed by a sequence of projections of type m , and the sequence might (or might not) conclude by one projection where the leaf is open in the beginning and closes at or before the end of the projection (again, we allow for the possibility that $t_{lp} = \bar{t}$ in this case). Schematically, “valid” sequences of leaf behavior during each opening can be: e , eb , em , emb , emm , $emmb$, $em \dots m$, $em \dots mb$. (In the following, we often refer to *leaf opening sequences* of projections, or sometimes simply to *sequences*.) Given this structure, the total number of leaf openings can be obtained by taking the sum of “ e ” variables.
4. The duration of a projection exceeds the lower bound on the minimum LOT, i.e., $\bar{t} \geq T^M$. This assumption implies that minimum LOT constraints only need to be explicitly imposed on sequences of the form e and eb ; for sequences that include one or more projections of type m the lower bound will be satisfied automatically. Notice that in this model, we do not

need to have a similar assumption regarding the lower bound on the average LOT T^A .

The model uses the following additional set of variables:

- b_{lp}, m_{lp}, e_{lp} — these variables are equal to 1 if the leaf l opens in the beginning, “middle,” and end of the projection, respectively,

and is given by:

$$\begin{array}{l} \text{minimize } F(z) \\ \text{variables } t_{lp}, z_v \\ e_{lp}, m_{lp}, b_{lp}, \beta_{lp} \end{array} \quad (2.4a)$$

$$\text{subject to } z_v = \sum_{l \in [L]} \sum_{p \in [P]} D_{lv} t_{lp}, \quad v \in \mathcal{V} \quad (2.4b)$$

$$\bar{t} m_{lp} \leq t_{lp} \leq \bar{t} \beta_{lp}, \quad l \in [L], p \in [P] \quad (2.4c)$$

$$e_{lp} + m_{lp} + b_{lp} = \beta_{lp}, \quad l \in [L], p \in [P] \quad (2.4d)$$

$$m_{lp} \leq m_{l_{p-1}} + e_{l_{p-1}}, \quad l \in [L], p \in [P] \setminus \{1\} \quad (2.4e)$$

$$b_{lp} \leq m_{l_{p-1}} + e_{l_{p-1}}, \quad l \in [L], p \in [P] \setminus \{1\} \quad (2.4f)$$

$$t_{lp} + t_{l_{p+1}} \geq T^M (e_{lp} + b_{l_{p+1}} - 1), \quad l \in [L], p \in [P - 1] \quad (2.4g)$$

$$t_{lp} \geq T^M (e_{lp} + e_{l_{p+1}} - \beta_{l_{p+1}}), \quad l \in [L], p \in [P - 1] \quad (2.4h)$$

$$\sum_{l \in [L]} \sum_{p \in [P]} t_{lp} \geq T^A \left(\sum_{l \in [L]} \sum_{p \in [P]} e_{lp} \right) \quad (2.4i)$$

$$\beta_{lp}, m_{lp}, e_{lp}, b_{lp} \text{ binary}, \quad l \in [L], p \in [P].$$

Here, constraints (2.4c) ensure that the time a leaf stays open at a projection is nonnegative and does not exceed the time spent at the projection; they incorporate forcing constraints that ensure positivity only if $\beta_{lp} = 1$ and that, if the projection has type m , the leaf stays open for the entire projection.

Constraints (2.4d) specify that each projection where a leaf is “open” must have one of three types: e , m , or b .

The next group of constraints ensures that each projection in which a leaf is open is classified as type e , m , or b according to the stated assumptions. In particular, if projection p has type m then projection $p - 1$ must have type e or m (constraints (2.4e)), and if it has type b then $p - 1$ must have type e or m (constraints (2.4f)).

Each of the constraints (2.4g) ensures that the combined LOT of the projection sequence of the form eb , if one starts at projection p , satisfies the corresponding lower bound. The right-hand side of the constraint is equal to T^M if both e_{lp} and b_{lp+1} are equal to one, and is 0 or negative otherwise.

Similarly, each of the constraints (2.4h) ensures that the LOT of each projection sequence of the form e , if one starts at projection p , satisfies the corresponding lower bound. This constraint requires a somewhat more detailed explanation. (We will drop the subscript l in this paragraph to make the explanation more concise.) First, if $e_p = 0$, then a new leaf opening sequence does not start at projection p , which can happen in two scenarios: either the leaf is closed at p , or it is open, but the opening occurred prior to p , and so projection p has type m or b rather than e . In the former scenario, either a new leaf opening starts in projection $p + 1$, in which case $\beta_{p+1} = e_{p+1} = 1$, or the leaf remains closed in projection $p + 1$, in which case $\beta_{p+1} = e_{p+1} = 0$ — in both of these cases, the right-hand side of the constraint evaluates to 0. In the latter scenario, in projection $p + 1$ the leaf may be closed ($\beta_{p+1} = e_{p+1} = 0$) or open with $\beta_{p+1} = e_{p+1} = 1$ (i.e., a new leaf opening sequence starts) or with $\beta_{p+1} = 1$ and $e_{p+1} = 0$ (i.e., the current leaf opening sequence continues) — in each case, the right-hand side of the constraint is 0 or negative. Summarizing, if $e_p = 0$, then the lower bound constraint on t_p is not enforced (correctly). On the other hand, if $e_p = 1$, then a new leaf opening sequence starts at projection p . Again, there are two scenarios to consider: the sequence does or does not continue in projection $p + 1$. In the former scenario, we have $\beta_{p+1} = 1$ and $e_{p+1} = 0$, so the right-hand side of the constraint is 0 and the lower bound on t_p is not enforced (correctly). In the latter scenario, we have $\beta_{p+1} = e_{p+1} = 0$ if the leaf is closed

in the next projection, or $\beta_{p+1} = e_{p+1} = 1$ if a new sequence starts — in both cases, the sequence that started in p indeed has form e , and the lower bound on t_p is enforced (correctly).

Constraint (2.4i) ensures that the average LOT satisfies the corresponding lower bound.

Note that the above model does not enforce the condition that $\beta_{lp} = 1$ only if $t_{lp} > 0$. Indeed, the combination $\beta_{lp} = 1$ and $t_{lp} = 0$ (i.e., a “phantom” β) is allowed by the model, but this is without loss of generality or optimality. To see this, suppose a solution (t, β, e, m, b) satisfies constraints (2.4), and consider another solution, $(t, \tilde{\beta}, \tilde{e}, \tilde{m}, \tilde{b})$, which has been updated by setting $\tilde{\beta}_{lp} := 1$ if and only if $t_{lp} > 0$, setting $\tilde{e}_{lp} := 1$ if:

1. $e_{lp} = 1$ and $\tilde{\beta}_{lp} = 1$
2. $m_{lp} = 1$ and $e_{lp} = 1$ and $\tilde{e}_{lp} = 0$ (a sequence of type $em \dots m$ is preceded by a “phantom” e)
3. $e_{lp-1} = 1$ and $t_{lp-1} = 0$ and $b_{lp} = 1$ (a sequence “ eb ” when b is preceded by a “phantom” e),

and 0 otherwise, by setting $\tilde{m}_{lp} := 1$ if $m_{lp} = 1$ and $t_{lp-1} > 0$, and $\tilde{m}_{lp} := 0$ otherwise. Finally, $\tilde{b}_{lp} := 1$ if and only if $b_{lp} = 1$ and $\tilde{\beta}_{lp} = 1$ and $\tilde{b}_{lp} := 0$ otherwise. This new solution satisfies constraints (2.4b), (2.4c) and (2.4d) by construction. Constraint (2.4e) is also satisfied because in case a leading e is removed, the next m will be replaced by a new leading e and m will always be preceded by an e or an m . Constraint (2.4f) is satisfied using a similar argument. Constraint (2.4g) will be satisfied whenever a sequence eb becomes $\tilde{e}\tilde{b}$ and it will also be satisfied in case it becomes \tilde{e} , since all contribution was attributed to the b segment (t_{lp} was zero). In addition, constraint (2.4h) will be satisfied because if we have a sequence of type \tilde{e} it means that it is derived from a sequence of type e that fulfills this constraint already, or it is constructed from a sequence of type eb , in which case the projection of type b will satisfy the constraint by itself. Constraint (2.4i) will also be satisfied because the right-hand side can only decrease (total number of openings can only decrease).

Note also that, similarly to the model (2.4), this model can underestimate the minimum or average LOT of a treatment plan; a post-processing step can be used to update the final plan metrics based on values of t .

Recall that this model assumes that $\bar{t} \geq T^M$. This supposition is not unreasonable to make with 51 projections per rotation of the gantry, where \bar{t} at the fastest rotation speeds would translate to approximately 300 milliseconds. Even a refinement of 153 projections per gantry rotation would result in \bar{t} of approximately 100 milliseconds. However, faster gantry speeds or finer grids of projections could violate this assumption. It is possible to modify the above model in case this problem arises by including additional constraints on durations of longer (in terms of the number of projections) sequences. For example, suppose $\bar{t} < T^M \leq 2\bar{t}$, i.e., a sequence of the type em or emb may be too short to satisfy the lower bound on LOT, but any sequence with 2 or more m -projections is guaranteed to be sufficiently long. Then the model should be augmented by adding constraints that enforce the lower bound for the sequences of these two types, similarly to (2.4g) and (2.4h).

2.4 Implementation

We tested the treatment planning models presented in Section 2.3 on several instances derived from a clinical prostate cancer case at various discretization resolutions. In this section, we discuss some informative details of our implementation.

2.4.1 Objective Function

The objective function $F(z)$ used in our models was implemented to be a convex, smooth piecewise quadratic function designed to guide the treatment towards satisfying physician’s goals; similar models have been used in many prior studies including [Peng et al. \(2012\)](#); [Peng \(2013\)](#); [Peng](#)

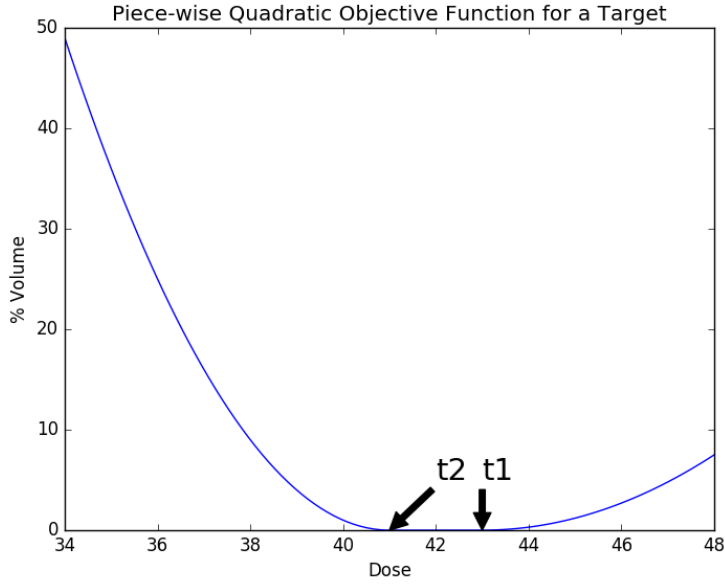


Figure 2.4: An illustration of a typical term of function $F(z)$ of (2.5) that includes penalties for both over- and under-dosing a voxel in a target structure.

et al. (2015); Long (2015). This objective function has the form:

$$F(z) = \sum_{v \in \mathcal{V}} (u_v(z_v - t_{v,1})_+^2 + o_v(t_{v,2} - z_v)_+^2). \quad (2.5)$$

Here, $0 \leq t_{v,2} \leq t_{v,1}$ are lower and upper thresholds, respectively, for the dose to voxel v , and u_v and o_v are weights associated with the (quadratic) penalty for, respectively, under-dosing and over-dosing the voxel. The values of these parameters are typically consistent for voxels inside each structure. For OARs, only over-doses are penalized, and so the lower thresholds are set to 0 (and commonly, so are the upper thresholds). For target structures, the parameters are selected to strongly penalize under-doses; for clinical reasons, over-doses are often penalized as well, but usually with $o_v < u_v$. An illustration of a typical term of the function associated with a target voxel is illustrated in Figure 2.4.

The values of coefficients in $F(z)$ are, clearly, specific to the particular treatment site and the OARs adjacent to the target. They are also often patient-specific; for example, u_v and o_v are

often linked to the number of voxels in each structure. Moreover, the particular patient geometry often demands that the weights as well as the threshold values are tweaked by trial-and-error until the solution to the optimization problem achieves an adequate treatment. After initial values are selected and the corresponding treatment is found, the treatment planner and the physician assess the treatment, and the parameters are adjusted with the goal of improving its deficient aspects; the process repeats until the physician concludes that the dose distribution is satisfactory.

2.4.2 Solver Options

Using the objective function defined by (2.5), we implemented the resulting mixed-integer quadratic programming models of Section 2.3 using Python 3.7, and used Gurobi 7.5.2 (Gurobi Optimization, LLC) as the solver.⁵

Due to the computational complexities inherent to Mixed-Integer Programming (MIP), it is important to emphasize the implementation of optimization strategies tailored for the particular problem at hand in order to achieve reasonable running times. We found that Gurobi struggled to solve larger-scale instances of our problems (especially problem (2.4)) until an appropriate combination of setting was discovered. We therefore include an overview of the options we have tried. While our discussion is focused on the specifics of Gurobi options (as discussed in *Gurobi Optimization, 2018*), other commercial solvers, such as, e.g., CPLEX, use similar concepts and paradigms for solving these types of problems, and can be similarly tuned to improve their performance; the details and specific options available would, of course, be different.

Gurobi optimizer for mixed-integer quadratic problems consists of four steps: a pre-solve, solution of the root relaxation, an application of a Branch-and-Cut algorithm, and a summary stage. The performance and impact of the pre-solve step varies greatly for each problem instance, and the root relaxation solution step does not involve many control parameters besides *CutPasses*; there-

⁵We observed that specifying the models using AMPL modeling language required slightly less computing time, but Python proved to be more conducive to experimenting with various versions of the models.

fore, most of our strategies for improving solution times using the available options were focused on the Branch-and-Cut stage. In the end, we noted the `CutPasses` parameter to be of little or no impact, and most of the run time was spent on solving the root relaxations, and so we anticipate that computational performance of the solver could have been improved if the solver allowed to incorporate an algorithm for solving the root relaxation that is specifically tailored to our problems.

2.4.3 Standard Solver Parameters

Many solver parameters are “standard,” i.e., their presence is common in most solvers for mixed integer programs. For example, empirical exploration in Gurobi was coerced to feasibility by a tightening of the default *FeasibilityTol* parameter from 10^{-6} to 10^{-7} . However, tightening the feasibility tolerance increases running time; we reduced this negative impact by modifying the high-level optimizer solution strategy parameter *MIPFocus*. After the optimizer has done some work, we can further modify the strategy, and move from finding the best lower bound to finding better feasible solutions using the parameters *ImproveStartGap* and *ImproveStartTime*. In Gurobi, these parameters trigger a change of the strategy when we achieve a particular MIP gap, or when some time has already passed. We use both of these triggers in our optimization runs.

2.4.4 Branching Priorities

The combinatorial structure of the model can sometimes be used to speed up its solution by specifying a particular branching priority strategy for its variables. For the detailed model (2.4), we were able to achieve a speedup of the Branch-and-Cut solution stage by prioritizing branching on fractional values of variables β_{lp} 's, since setting $\beta_{lp} = 0$ immediately fixes the values $e_{lp} = m_{lp} = b_{lp} = 0$, removing several potential branches from the tree. A similar argument shows that the next highest priority should be given to e_{lp} 's, since setting $e_{lp} = 0$ reduces uncertainty about the values of b_{lp+1} and m_{lp+1} for the next projection. We set these priorities by assigning

appropriate values of Gurobi parameters *BranchPriority* for each family of variables.

2.4.5 The Partition Heuristic

We found the *partition heuristic* to be very useful for both models. The partition heuristic splits the variables into several groups and runs a Large-Neighborhood Search (LNS) algorithm on each of these groups by fixing the values of the variables outside the group, and solving the resulting sub-problem for the variables in the group. This heuristic has the potential to become very expensive if not specified carefully. The following choices were instrumental in improving performance.

The number of groups G should be related to the number of threads available in the computer. Our machine had 12 threads, and therefore, it made sense to create 12 groups. $G = 24, 36, 48, \dots$ would also have been reasonable.

In defining the groups, we took advantage of the progressive nature of the helical gantry path, and grouped variables according to their natural location in the patient's body. Group 1 contained variables corresponding to projections closer to the patient's head, and group $\#G$ — closer to their feet, i.e., we assigned all binary variables corresponding to projection p to group $g = \lceil p/G \rceil$.

We assigned all z 's as shared variables to all groups. It is possible we would have achieved further speed-ups if we instead had partitioned these variables to different groups as well, by associating a portion of z variables to their "closest" group of projections, either in terms of physical location or in terms of the relative magnitude of dose deposition coefficients.

We set the parameter *partitionPlace* to 30, so the partition heuristic runs at the following steps of the optimization:

1. Before the root relaxation is solved
2. At the start of the root cut loop
3. At the end of the root cut loop

4. At the nodes of the branch-and-cut search

2.4.6 Hints

Hints bestow the model with a high-quality indication of the value that a variable might take. These hints impact the heuristics and branching decisions when the solver is exploring the search tree. High-quality hints enhance the exploration of high-quality integer solutions.

We used as our hints the values of the variables in the optimal solutions of continuous relaxations of our models, i.e., convex continuous optimization problems obtained by replacing binary restrictions on the variables by lower and upper bounds of 0 and 1, respectively. These relaxed models were easily solved (within a few seconds), and their solutions provided reasonable hints that sped up the subsequent runtimes of mixed-integer models. The nonlinear solver we selected to solve the relaxation in order to obtain the hints was different to the one used within Gurobi's MIP solver to tackle the root relaxation. Gurobi adds cuts to the root node relaxation problem; the addition of cuts speeds up the subsequent Branch-and-Cut step of the solution process, but makes the root node relaxation solution time significantly longer.

2.4.7 Warm Start

Further enhancement to solution times of our models can be obtained by providing a warm start solution, i.e., a feasible solution which can be used by the branching algorithm to prune tree nodes by bound at earlier stages of the algorithm.

To create such a warm start solution, we created a smaller instance of the problem by using a coarser discretization in the voxel space, a process referred to as downsampling, while keeping the same number of projections, and deriving the dose deposition coefficients for the downsampled model based on the original values. (For instance, in the prostate case discussed in Section 2.5, the full resolution instance consists of 16,677 voxels, and by downsampling, we created an instance

with 1,385.)

By solving the downsampled instance, we obtain feasible values of projection-indexed variables, say, $\hat{\beta}_{lp}$, \hat{m}_{lp} , \hat{e}_{lp} , \hat{b}_{lp} , \hat{t}_{lp} in the detailed model. We can then calculate the corresponding values of \hat{z}_v in the full resolution, making the combined vector a feasible solution that can be used as a warm start for the full resolution model.

Experiments indicate an improvement of more than 25% in the running time when this warm start is implemented together with the partition heuristic above. Experiments also show that the partition heuristic is necessary in order to take full advantage of the warm start, given the nonlinear nature of the problem.

Since the low-resolution problem is itself a nonlinear mixed-integer program, we can further improve the running time of the overall scheme by providing it with hints as well. In extreme cases, i.e., when solutions to even higher-resolution instances with more voxels are desired, we envision the creation of nested warm starts.

2.5 Experiments and Results

We evaluated the models proposed in the previous section on a dataset provided by our collaborators at the UT Southwestern Medical Center Department of Radiation Oncology. The preprocessing step which included data acquisition, contouring, and calculation of the dose deposition coefficients D_{lvp} was performed using UT Southwestern’s Collapsed Cone Convolution/Superposition Algorithm (CCC) algorithm on GPU’s. The dataset’s leaf length is 2.5 cm and the pitch is $0.287 = 0.86/3$. We will assume a gantry speed of 4 rotations per minute in order to mimic the experiments of *Westerly et al. (2009)*.

The dataset was based on a clinical prostate case. The number of structures in the case is 19 and the treatment goals are set to satisfy not only the Radiation Therapy Oncology Group (RTOG) requirements presented in Table 2.1, but also some stricter restrictions on the rectum dose from the

Goal	Delivered
Target Dose	78 Gy
Target Hot Spot	99% below 83 Gy
Target Cold Spot	99% above 73 Gy
Rectum	50% below 60 Gy
Bladder	50% below 45 Gy
Bladder	95% below 80 Gy
Penile Bulb	50% below 51 Gy

Table 2.1: Goals of the prostate case according to RTOG requirements.

literature. According to research by *Biegala and Hydzik (2016)*, 50% of the rectum should get a dose smaller than 35Gy, and 17% should get less than 65Gy. The rectum-associated restrictions are medically the most difficult to satisfy, and therefore, the rectum is identified in the literature as the dose-limiting organ in prostate radiotherapy (*De Meerleer et al., 2004*).

We performed experiments on two different representations of the above case. The main difference between the representations is the discretization of the gantry trajectory into 51 projections per rotation in the first instance and 153 projections per rotation in the second. The voxel discretization was also slightly different, with 16,677 voxels in the first instance and 16,686 voxels in the second.

2.5.1 Treatment Plan Evaluation Tools

Before presenting our experiments, we first describe the visual tools we will use to evaluate and compare their results, namely, DVH plots (Section 2.5.1.1), the sinograms (Section 2.5.1.2), and the LOT histograms (Section 2.5.1.3).

2.5.1.1 Dose Volume Histogram (DVH) Plots

The DVH plots are one of the most commonly used tools for displaying a quantitative summary of dose distributions in different structures. The DVH plot is a set of curves, each curve corresponding

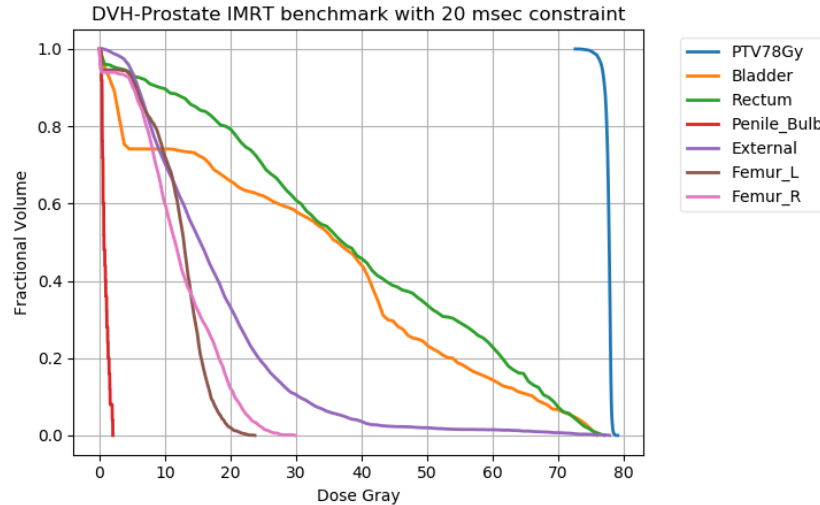


Figure 2.5: DVH plots of planned doses resulting from the FMO model for the prostate case with 51 projections per gantry rotation.

to a different structure of interest, on a plane formed by the dose (horizontal axis) and the Fractional Volume (vertical axis). For each structure, a point on the corresponding curve is the percentage of the structure receiving a specific dose or higher. DVH curves are monotonically non-increasing.

One of the goals of radiation therapy planning is to design a treatment that delivers high doses to the targets; therefore, we would like to see the DVH curves corresponding to targets that stay to the right of the OAR's DVH curves. A good treatment will generally display a rapidly decreasing, almost vertical, curve around the target dose. Another goal of radiation therapy is to deliver low doses to the OARs. Typically, OAR DVH curves will gravitate towards the left of the graph. Figure 2.5 shows an example of a DVH plot of a dose distribution for a prostate case (namely, the planned dose distribution of a treatment plan obtained with the FMO model for the case with 51 projections per gantry rotation).

We will use DVH plots to visualize results in Chapter 3 as well.

2.5.1.2 Sinograms

In tomotherapy, a sinogram is a 2D array that is used as a visualization tool of the treatment itself (rather than the dose delivered to the patient). Depending on the treatment paradigm used, the visual conventions of the sinograms have to be adjusted. To depict the conventional treatments, the Planned Fluence Sinogram is typically used, whereas to illustrate the treatments developed by our proposed approaches, we used the Leaf Control Sinogram.

Planned Fluence Sinogram Traditionally, the sinogram is used to represent the measurements collected at the LINAC's exit detector, registering the fluence delivered from every single leaf at every single projection. The sinogram's trail corresponds to the orders passed as an input to the planning system. In the academic setting, the sinogram is a 2D plot with the axes corresponding to the projections and the leaves, respectively (different authors make different choices regarding whether the projections will be displayed along the horizontal or vertical axis). In the conventional approach, the sinogram takes the form of a heatmap, where each color pixel represents the LOT of the corresponding projection-leaf combination. The sinogram plot shows the evolution of each leaf across projections, with black pixels indicating that the leaf was closed, and colored pixels indicating open leaves and the LOT of the corresponding beamlets.

Leaf Control sinogram. For the treatments delivered with our proposed approach, the above form of sinograms is not particularly useful, since we allow the leaves to stay open longer than the time it takes to traverse a projection. Instead, we will use what we termed the leaf control sinogram to visualize treatments. This sinogram has leaves plotted along one axis, and the time it takes for the gantry to traverse its trajectory — against the other. For each leaf, we use a color to visualize when it is open. Projections can provide a grid on the time axis, but otherwise the leaf control sinogram is independent of the density of projections discretization. Another useful feature of this type of sinogram is that it can be used to compare two different treatments by overlaying two sinograms using different colors in the same plot. Figure 2.15 provides an example of such a

comparison; we discuss it in detail in the following section.

2.5.1.3 Leaf Open Times Histograms

As we have already mentioned, the goal of our proposed alternative approaches to tomotherapy treatment planning is to increase leaf opening times during delivery in order to decrease dosimetric errors. We will use traditional histogram plots to visualize distributions of LOTs (across all leaves and all projections) we obtain; see Figure 2.6 for an example.

2.5.2 Prostate Case Results

Unless mentioned otherwise, we used the full voxel resolution for each of the two cases (or, more precisely, two representations of the same clinical case, with 51 and 153 projections per rotation, respectively). The parameters of the objective function $F(z)$ were chosen by manual trial-and-error iteration based on the solutions to the FMO model (2.1) with $T^M = 20$ milliseconds.

We ran our experiments on a desktop powered by an Intel Core™ i7-8700 CPU at 3.20GHz, with six cores and 12 logical processors, and 64 Gb of RAM.

2.5.2.1 Analysis of Case 1: 51 Projections per Gantry Rotation

The solution to the FMO model with $T^M = T^A = 0$ has an average leaf opening time of 125 milliseconds; this is very similar to the average of 132 milliseconds for such cases found in the literature ([Westerly et al., 2009](#)) due to attenuation in the prostatic region. The LOT histogram resulting from the FMO model with $T^M = 20$ milliseconds and $T^A = 0$ is shown in Figure 2.6. It generates a plan with the average LOT of 220 milliseconds.

We next compare the plan obtained as a solution to the Simple model (2.3) with $T^M = 20$ milliseconds and $T^A = 170$ milliseconds to the FMO solution with $T^M = 20$ milliseconds and $T^A = 0$. As depicted in Figure 2.7, the planned dose distributions produced by the two models are

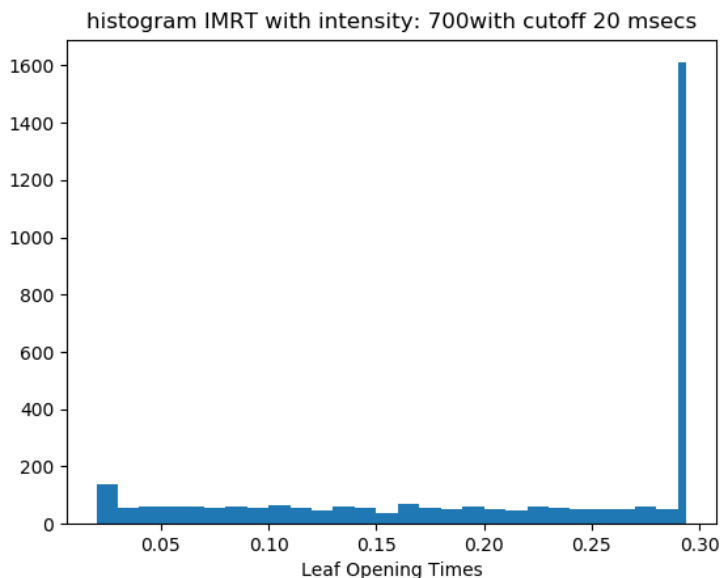


Figure 2.6: LOT histogram of prostate case plan with 51 projections per rotation resulting from the FMO model with $T^M = 20$ milliseconds and $T^A = 0$.

nearly identical. It should be emphasized that all the DVH plots presented in this chapter (as well as Chapter 3) are created using *planned* dose distributions of the proposed treatments, whereas *delivered* dose distributions will deviate from the planned ones. Moreover, since the treatments resulting from different planning and delivery paradigms will have different delivery characteristics, these dose discrepancies will be different as well; if our proposed approaches indeed produce plans that have better delivery characteristics, their delivered doses should adhere to the planned ones more closely.

In this and forthcoming experiments we used $T^A = 170$ in the constraint on average LOT in all instances of Simple and Detailed models. The main reason for this choice was our desire to explicitly require a reasonably high average LOT that was enforceable in the experiments with both 51 and 153 projections per rotation, for the sake of making comparisons. In the latter instances, the Simple model would only allow LOTs of at most 200 milliseconds, and we chose the lower bound on the average LOT to be slightly smaller than that.

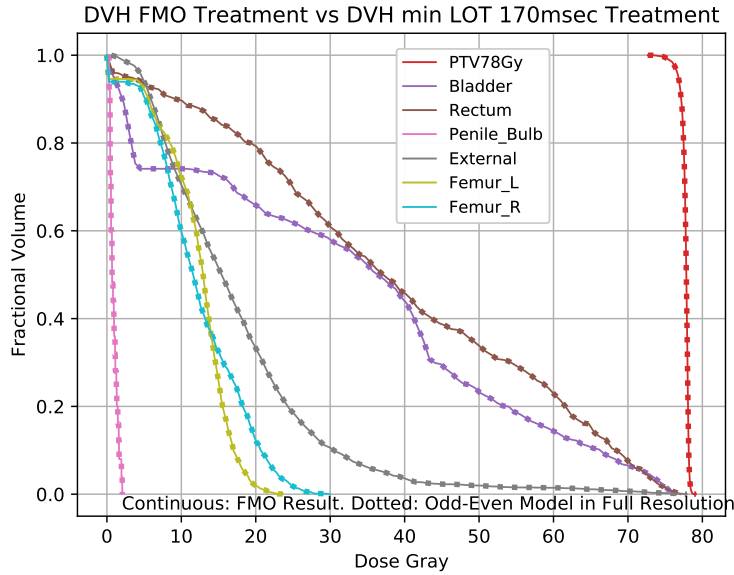


Figure 2.7: Dose-Volume Histogram comparison of prostate case plans with 51 projections per rotation, resulting from the FMO (continuous curves) and Simple (dotted curves) models, both with $T^M = 20$ milliseconds and $T^A = 170$ milliseconds in the Simple model.

Finally, we compared solutions to the Simple model and the Detailed model (2.4), both using parameter values $T^A = 170$ milliseconds and $T^M = 20$ milliseconds. We compare the resulting treatments using a leaf control sinogram in Figure 2.8: we use blue color to represent LOTs in the plan obtained by the Simple model and red color to represent LOTs in the plan obtained by the Detailed model. Consequently, purple color represents the regions where LOTs are common to both models. The sinogram comparison shows that the plans produced by the two models are fairly similar. Furthermore, the DVH plots show that both treatments produce almost identical planned dose distributions, as shown in Figure 2.9, which also closely resemble the FMO dose distribution in Figure 2.5.

In view of the similarity of the sinograms, it is unsurprising that the LOT histograms also look very similar; see Figure 2.10. The average LOTs obtained with the Simple and Detailed models were 360.8 milliseconds and 388.5 milliseconds, respectively (by comparison, recall that the average LOTs obtained with the FMO model with $T^M = 0$ and $T^M = 20$ were 125 and 220 mil-

Sinogram Comparison of Odd-Even Model vs. Detailed Model

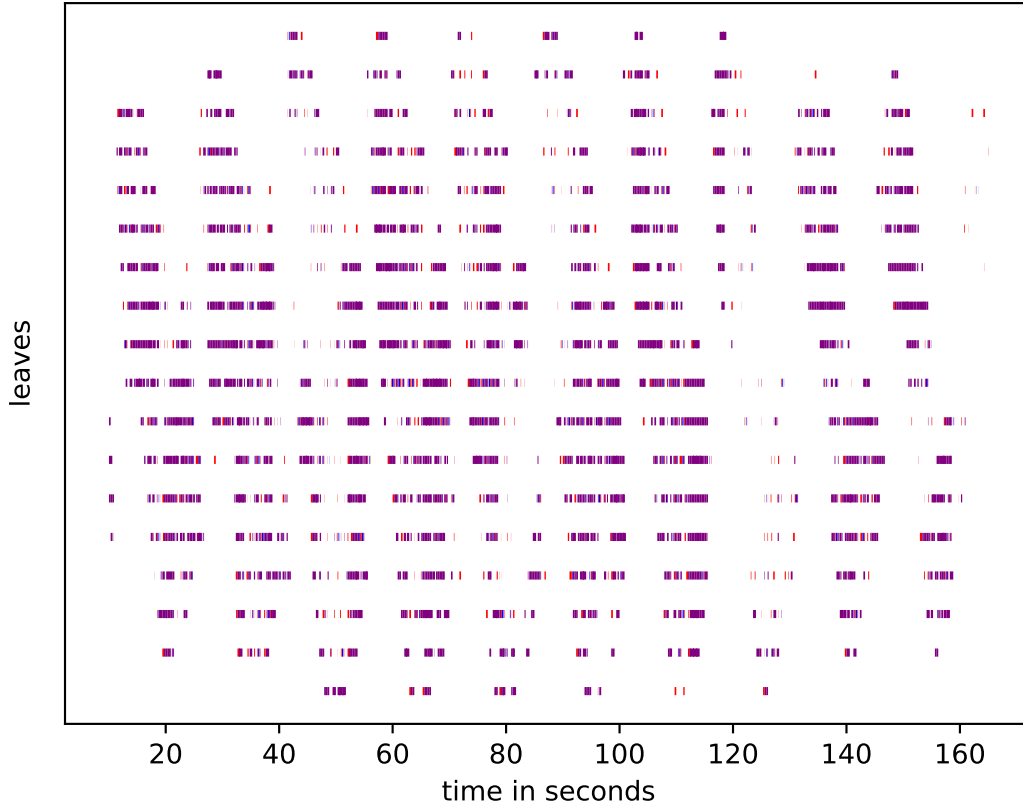


Figure 2.8: Leaf Control Sinogram comparison of prostate case plans with 51 projections per rotation, resulting from the Simple and Detailed models with $T^A = 170$ milliseconds and $T^M = 20$ milliseconds. Blue color represents LOTs in the Simple model, red color represents LOTs in the Detailed model, and purple color represents the regions with LOTs common to both models.

Comparison of DVH plots: Odd-Even model vs. Detailed model

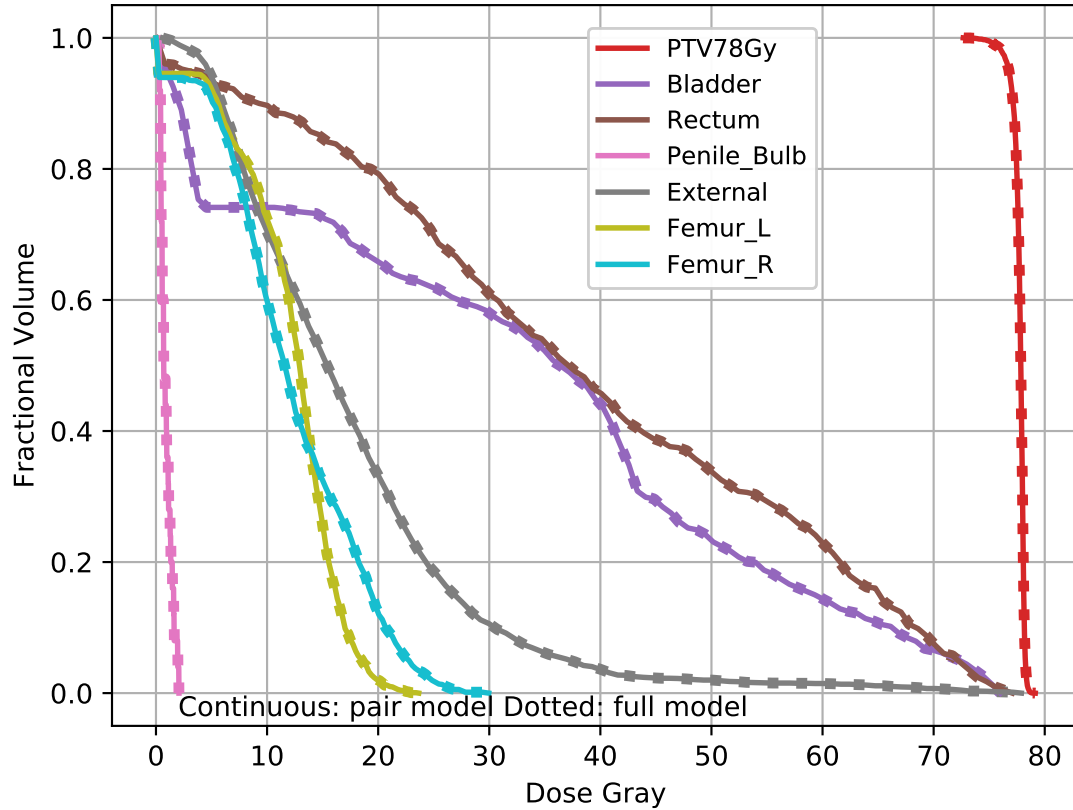
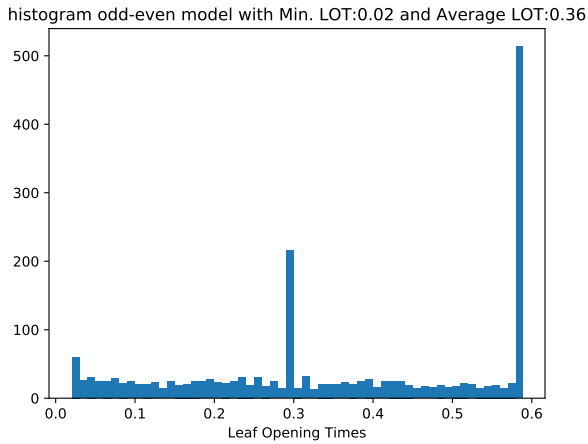
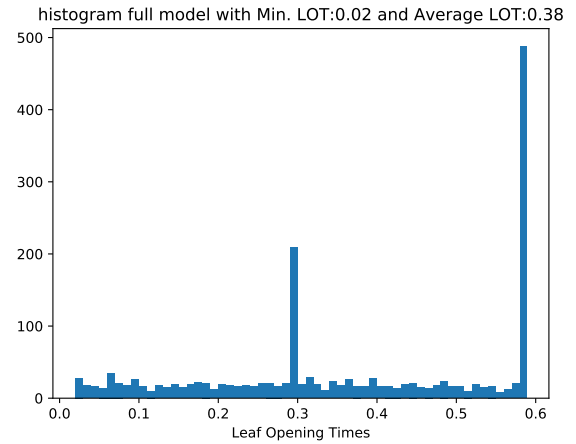


Figure 2.9: Dose-Volume Histogram comparison of prostate case plans with 51 projections per rotation, resulting from the Simple (continuous curves) and Detailed (dotted curves) models with $T^A = 170$ milliseconds and $T^M = 20$ milliseconds. Despite close similarities, there are, in fact, small differences between the two sets of DVH plots.



(a) LOT histogram, Simple model



(b) LOT histogram, Detailed model

Figure 2.10: LOT histograms of prostate case plans with 51 projections per rotation, resulting from the Simple and Detailed models with $T^A = 170$ milliseconds and $T^M = 20$ milliseconds; (a) Simple model, (b) Detailed model.

liseconds, respectively). To summarize, for the case with 51 projections per rotation, the Simple and Detailed models produce nearly-identical treatment plans, both in terms of LOT distributions and DVH plots (of planned dose distributions), the latter of which are also nearly identical to the planned DHV plot obtained from the FMO model solution.

The reduction in leaf pulsation events (i.e., the total number of times the leaves open and close during a treatment) is a welcome byproduct of increasing LOTs, since reducing the number of pulsations reduces the wear-and-tear of the multileaf collimator. Table 2.2 summarizes the number of pulsation events in the plans obtained by the four models we have discussed, showing a significant decrease obtained by the proposed models.

Exercise	# Events
FMO (No constraints)	4958
FMO $T^M = 20$ msecs, $T^A = 0$	3199
Simple model $T^M = 20$ msecs, $T^A = 170$ msecs	1952
Detailed model $T^M = 20$ msecs, $T^A = 170$ msecs	1680

Table 2.2: Number of leaf pulsation events in the plans with 51 projections per rotation.

In this instance, it took us 28,934 seconds, or about 8 hours, to solve the Simple model within a 0.001% optimality gap. The Detailed model took 163,194 seconds (over 2 days!) to achieve a 0.001% optimality gap; within 5 hours, the optimality gap in this model was 10%, and the corresponding solution was quite reasonable. Since the Simple model in this case obtained nearly identical results in terms of DVH plots and LOT statistics in less time than the Detailed model, for the 51 projection case with 20 millisecond minimum LOT, the Simple model appears sufficient.

2.5.2.2 Increasing Minimum and Average LOTs

As discussed in the previous section, the average LOT attained by solutions to both the Simple and the Detailed models already exceed the lower bound of 170 milliseconds imposed as a constraint. Some lower bound constraints on the individual LOTs were active, as the histograms in Figure 2.10 indicate. However, if higher minimum LOT is desired, the value of T^M can still be increased with negligible impact on the dose, as shown in the DVH plot in Figure 2.11, where we compare treatment plans with a minimum LOT T^M of 20 milliseconds and 40 milliseconds, and $T^A = 170$ milliseconds.

Despite the similarities in DVH plots and average LOTs (which were 385 milliseconds and 330 milliseconds in the solutions to instances with $T^M = 20$ milliseconds and $T^M = 40$ milliseconds, respectively), there is a noticeable impact on the distribution of LOTs, as shown in Figure 2.12. The treatment plan obtained by setting $T^M = 40$ milliseconds included leaf opening sequences that were up to 4 projections long. This “chaining” of projections would not have been possible in the Simple model, since it allows merging of only up to two projections.

The DVH plots start showing small differences when we compare the solutions to the Simple and Detailed models using 153 projections per rotation; see Figure 2.13. The LOT histograms of these two plans are quite different, as shown in Figure 2.14 (note that the horizontal axes of the two histograms have different scales). Notice that the solution to the Simple model contains

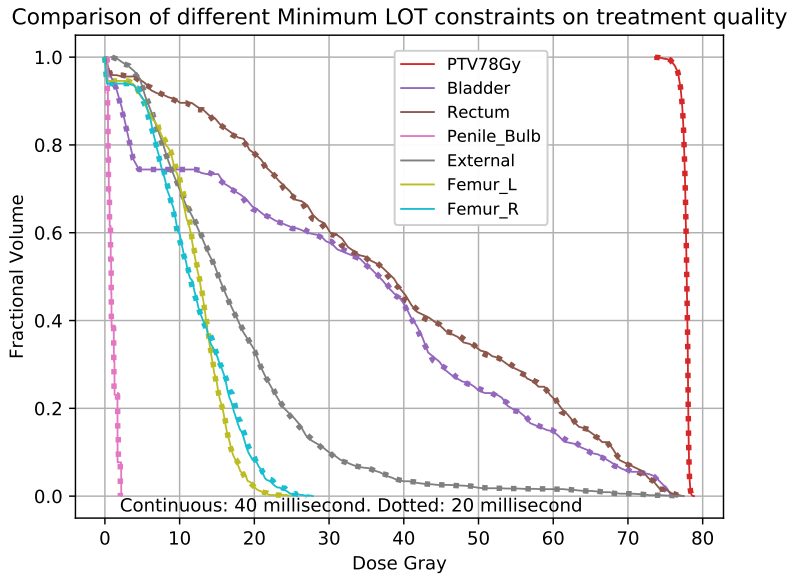


Figure 2.11: DVH comparison of prostate case with 51 projections per rotation, resulting from the Detailed model with $T^M = 20$ milliseconds (dotted curve) and $T^M = 40$ milliseconds (continuous curve). Both instances used $T^A = 170$ milliseconds and were solved to achieve a 0.01% optimality gap.

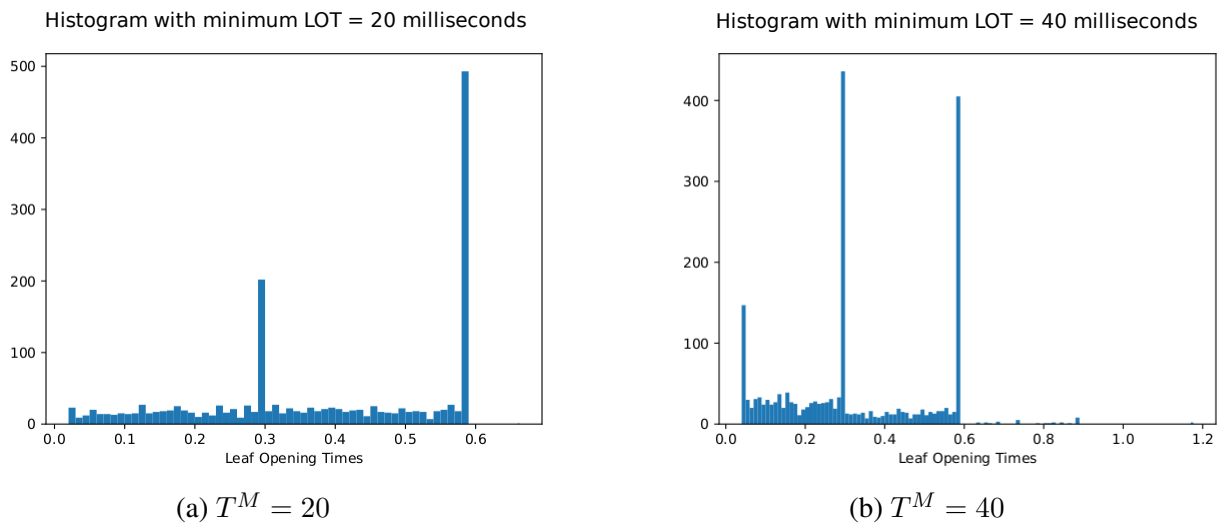


Figure 2.12: LOT histograms of prostate case plans with 51 projections per rotation, resulting from the Detailed model with $T^A = 170$ milliseconds and (a) $T^M = 20$ milliseconds and (b) $T^M = 40$ milliseconds.

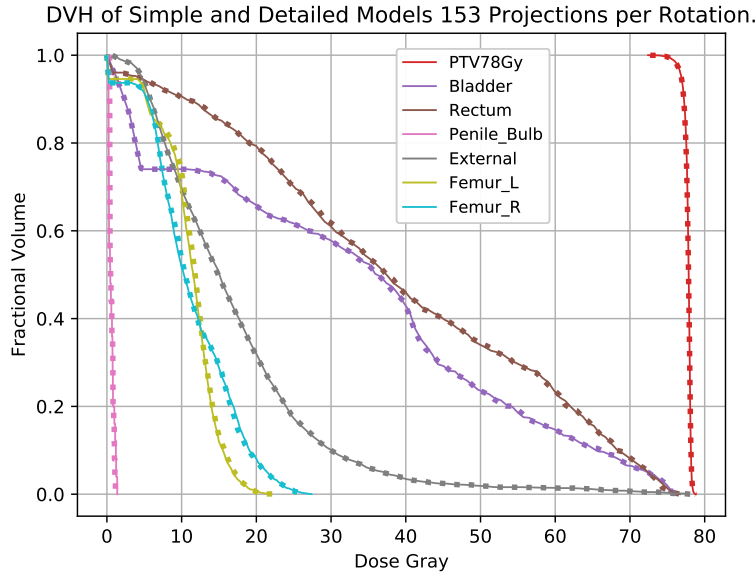


Figure 2.13: DVH comparison of prostate case with 153 projections per rotation, resulting from the Simple (continuous curves) and Detailed (dotted curves) models with $T^M = 20$ milliseconds and $T^A = 170$ milliseconds. Both instances were solved to achieve a 0.01% optimality gap. Some differences in the DVH plots are more prominent, especially for the rectum.

many leaf opening times that are equal to the maximum possible ($2\bar{t} = 200$ milliseconds with 153 projections per rotation), whereas the Detailed model does not have this limitation. Moreover, the average LOT constraint becomes binding in the solution to the Simple model, while the average LOT attained by the solution to the Detailed model is similar to the solutions obtained using 51 projections per rotation.

2.5.2.3 Discrepancies between Planned and Delivered Doses

At present, we don't have the tools to provide a precise estimate of the expected reduction in the discrepancy between the delivered and planned doses resulting from using our models for treatment planning. However, we hypothesize that the increase in the average LOT can be a predictor of significant improvement. Recall that our Simple model in the case with 51 projections per rotation achieved an increase of 188% in average LOT compared to the FMO model with $T^M = 0$, or 64%

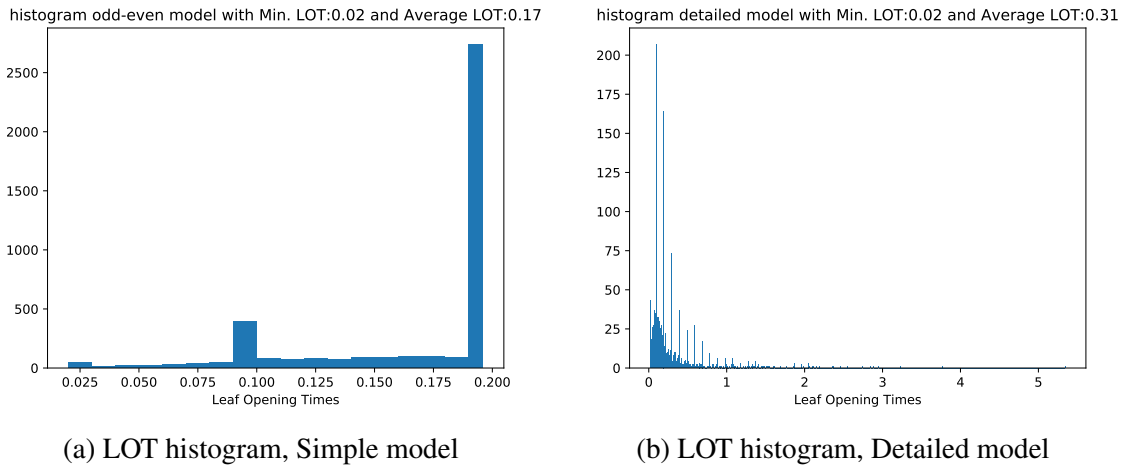


Figure 2.14: LOT histograms of prostate case plans with 153 projections per rotation, resulting from the Simple and Detailed models with $T^M = 20$ milliseconds and $T^A = 170$ milliseconds; (a) Simple model, (b) Detailed model.

compared to the FMO model with $T^M = 20$, without changing the pitch. For comparison, *Westerly et al. (2009)* achieved increases in average LOT ranging from 29.8% to 83.1% by increasing the pitch, which can have detrimental effects of its own. They conclude that the reduction of point dose discrepancies (points where the delivered dose deviated from the plan by more than 3%) was at least 68%.⁶ We conclude that our reduction in dose discrepancies is likely to be at least as high, if not higher than that.

Moreover, our proposed models, especially the Detailed model, allow for refinement in the discretization of the gantry trajectory into projections without affecting the delivery characteristics of the resulting treatment plans, which can further reduce dose discrepancies due to “blurring” (see Section 2.6.2 for further discussion).

⁶Table 4 in *Westerly et al. (2009)* presents a case where the proportion of point dose discrepancies at the 3% level decreased from 4.96% to 1.59%. For another case in the same table, the decrease is from 4.47% to 0.06%: a near complete elimination of dose discrepancies that will cause a treatment plan to fail QA.

Instances and parameters	root relaxation	branch-and-cut time to 10% gap	branch-and-cut time to 1% gap
Simple 51 v=8340 mLOT=20msec aLOT=170msec	2306	3064	3065.45
Simple 51 v=8340 mLOT=30msec aLOT=170msec	2260	3164	8691
Simple 51 v=16677 mLOT=20msec aLOT=170msec	7568	11353	28934
Detailed 51 v=8340 mLOT=20msec aLOT=170msec	18554	19302	37837
Detailed v=8340 51 mLOT=40msec aLOT=170msec	10431	37305	
Detailed 51 v=16677 mLOT=20msec aLOT=170msec	9753	18762	163194
Detailed 153 v=8344 mLOT=20msec aLOT=170ms	120857	173984	309523
Simple 153 v=16686 mLOT=20msec aLOT=170msec	320730	389211	
Detailed 153 v=16686 mLOT=20msec aLOT=170msec	320730	389211	

Table 2.3: Solution times of instances of Simple and Detailed models, in cumulative seconds since the beginning of the root relaxation phase. Here, $mLOT = T^M$, $aLOT = T^A$, v = number of voxels, and all times are shown in seconds

2.5.3 Running Times

Table 2.3 contains a representative summary of the solution times of several instances of Simple and Detailed model at different stages of the solution process. We include results for the down-sampled instances that were used to generate warm start solutions.

Whenever the Simple model and the Detailed model were compared, the former solved faster, which is to be expected due to its simplicity. The case with 51 projections per rotation can be solved at full voxel resolution within a 10% optimality gap in 11,353 seconds (3:10 hours). The cause of the “gridlock” for all instances is the time it takes to solve the root relaxation. We experimented with a reduction in the parameter *CutPasses* and *Cuts*, but got negligible improvement. It should be noted that the implementation of the solver for convex quadratic MIP in Gurobi currently provides

only one algorithm option for solving the root relaxation, which we have found to be significantly inferior to many other algorithms for solving convex quadratic programs (including the algorithm we used to solve the relaxation outside of the Gurobi’s Quadratic MIP solver to provide hints, as discussed in Section 2.4.6). An alternative solver with greater flexibility at this step would likely be able to achieve significantly faster solution times of instances of both Simple and Detailed models.

The sinogram in Figure 2.15 provides some insights into the quality of warm start solutions obtained by the procedure described in Section 2.4.7. We considered an instance of the Simple model with 51 projections per rotation and $T^M = 20$ milliseconds and $T^A = 170$ milliseconds. Recall that we used the solution to the instance with low voxel resolution (1385 voxels) as a warm start for solving the instance with high voxel resolution (16677 voxels). In the sinogram, the leaf opening times in the solution for the low-resolution instance are shown in red, high-resolution instance — in blue, and common leaf opening times are shown in purple. The presense of purple segments suggests the usefulness of a warm start; however, there are still significant differences between the warm start and ultimate solutions.

2.6 Discussion

2.6.1 Modulation Factor

The tomotherapy delivery LINAC uses constant rotational periods fixed ex-ante. Under the standard paradigm, *the longest* leaf opening time is essential because longer maximum leaf opening times demand proportionally slower rotational speeds. The standard paradigm defines the *Modulation Factor* as the ratio of maximum leaf opening time to the average of all non-zero leaf opening times. While higher values allow “greater variation in the leaf opening time” ([Sheng, 2017](#)), practical studies suggest that modulation factors below 2.5 are the recommended values for most cancer sites ([Westerly et al., 2009](#)). Plans with large modulation factors are less efficient than plans with

Sinograms of Low Resolution Model (red) vs. Full Resolution Model (blue)

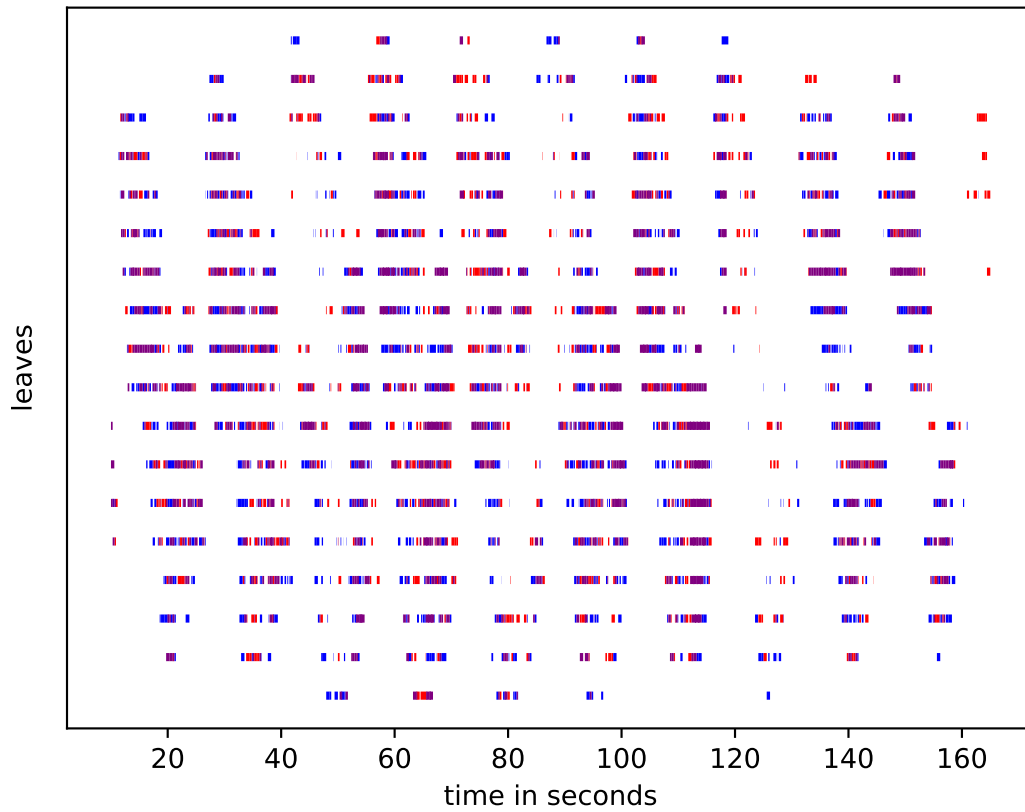


Figure 2.15: Leaf Control Sinogram comparison of solutions to instances of the Simple model with $T^M = 20$ milliseconds and $T^A = 170$ with low (red) and high (blue) voxel resolutions.

small modulation factors, because they require higher monitor units and produce more leakage and longer treatment times.

In our approach, the concept of the modulation factor becomes irrelevant for two reasons. First, we assume an a priori determined speed of the gantry that is not going to be influenced by a modulation factor. Second, the modulation factor loses its meaning when leaves are allowed to stay open longer than a single projection. Since a leaf-opening event can span more than one projection, high modulation factor in the new paradigm is not an indication that the machine should slow down its rotational speed. Furthermore, we are allowed to preserve greater variation in the leaf opening time mentioned above. All these considerations make the modulation factor obsolete under the new paradigm.

2.6.2 Resolution Increase Capabilities

The standard number of projections per gantry rotation is 51. Each of these projections spans roughly 7.06° , and while not a consequence of any particular physical limitations of the tomotherapy machine per se (see *Zhao et al., 2008*, who write, “the helical tomotherapy unit may deliver radiation continuously over a full rotation”), this number of projections is now standard in all research and practice since the publication of *Olivera et al. (1998)*. The Hi-Art TPS software uses a sinogram input with 51 projections as default which increases the popularity of this choice even further, although earlier models used to divide the arc into a different number of projections, some of which may be more intuitive, such as the 32 and 64 projection inverse-planning algorithms of *Holmes et al. (1995)*, or the 72 5° -projections used by *Kapatoes et al. (1999)*.

Some research articles have suggested that 51 projections per gantry rotation may be too coarse for some cancer sites where this discretization has been used. *Yang et al. (2012)* finds a 2° separation between control points to be the only acceptable level for QA evaluations in VMAT treatments. For tomotherapy, among others, *Hardcastle et al. (2012)* cite a 2011 field safety notice from To-

tomotherapy identifying a deficiency in the dose calculation for 51-projection treatment planning systems, especially in high-dose treatments with small, off-axis targets, due to what they called dose blurring resulting from the approximation of the 7.06° arcs by static gantry angles. As a remedy, they proposed “supersampling” of the gantry trajectory by virtue of tripling the number of projections. *Stambaugh et al. (2015)* also suggest further subdivision of each traditional projection into 2 or 3 projections to deal with this issue (102 or 153 projections per gantry rotation, respectively). *Tudor and Thomas (2013)* suggested that only a subdivision of 5 subprojections is acceptable for critical cases. They suggested a non-homogeneous distribution of gantry projections: 51 projections per rotation in non-critical regions, and 255 projections per rotation in critical segments of the gantry trip.

However, due to the requirement of closing every leaf whenever it opens during the same projection under the conventional treatment planning approach, increasing the projection resolution would have the unpleasant effect of shortening leaf opening times. In contrast, the new treatment planning and delivery paradigms we propose, especially the Detailed model, can increase the projection resolution while maintaining desired minimal and average LOTs. For the users, the only changes required are a different form sinogram specification, and several types of sinogram can be reported as output (*Van Dyk, 1999*).

2.7 Conclusions

The current tomotherapy standard for radiation therapy treatment assumes instantaneous leaf openings and closings. The helical Tomotherapy planning system partially solves this problem by modeling the pneumatic movement of the leaf with a constant speed. The LINAC MLC mechanism is rather precise in this endeavor.

Unfortunately, this approximation does not entirely solve the problem. The leaf movement is not exactly linear but has an acceleration/deceleration profile that is very difficult to model;

moreover, different leaves in a typical MLC have different motion patterns. These nonlinearities and distinctions accumulate and can create discrepancies in the dose delivery to the patient.

Moreover, the treatment delivery paradigm, and the corresponding treatment planning models, require excessive leaf pulsations, not only increasing the number of leaf transition events and associated errors, but also damaging the machinery. Other problems encountered in the conventional delivery include the insufficient resolution of the trajectory of the gantry. Paradoxically, we cannot refine this resolution without creating even more leaf opening and closing events.

Our new proposed paradigms address these problems by giving us the freedom to keep the leaves open for more than one projection, by cutting down on short leaf-opening times, and by opening the door for a smooth transition into high-resolution models. We have introduced two treatment planning optimization models under these paradigms. The first one achieves satisfactory practical results with faster computational times. The second model is a detailed version that achieves slightly better results with a more substantial computational burden. The first model can be implemented today and can have immediate impacts on Tomotherapy planning dose inaccuracy reduction. The second model looks forward into the near future and can only be deployed today on the simplest cases. The ultimate benefits to the patients will be the increased reliability in the radiation therapy treatments and control of side effects due to dose miscalculations.

To conclude, our work:

- Improves the delivery characteristics of treatment plans, namely,
 - increases leaf opening times by explicitly constraining minimum and average LOT, and
 - reduces the number of leaf pulsations, increasing the life expectancy of the multileaf collimator (*Mackie et al., 1993*),

which can reduce discrepancies between the planned dose and the dose delivered to the patient;

- Allows unlimited refinement of projections, this will reduce the resolution-related dose discrepancies (*Tudor and Thomas, 2013*).

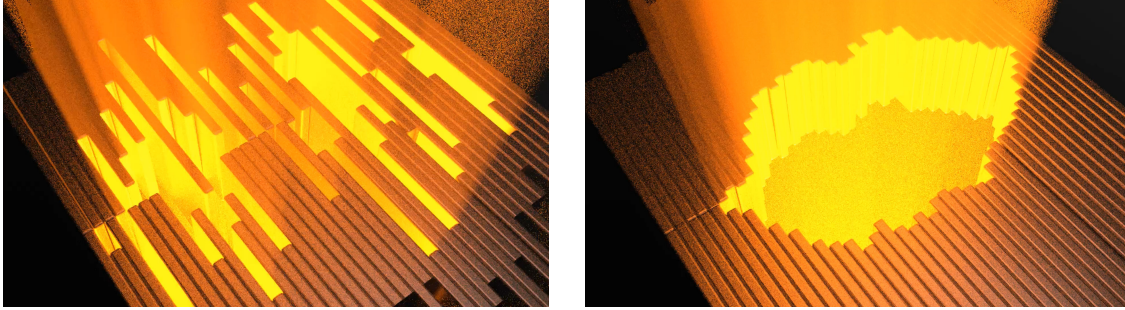
CHAPTER 3

VMAT with Aperture Control

3.1 Introduction

The utilization of Volumetric Modulated Arc Therapy (VMAT) has substantially increased since it was introduced into clinical use over a decade ago. Compared to the more traditional IMRT treatments using a limited number of fixed beam angles, VMAT can frequently deliver shorter treatments using fewer Monitor Units (MUs) without sacrificing conformity and other measures of treatment quality.

Since VMAT involves aperture shaping during continuous motion of the gantry along the treatment arc(s), it requires more complex approaches to treatment planning, many aspects of which still stand to be improved. For example, current approaches to treatment planning do not take into account the impact of the shapes of the apertures. Irregularly-shaped apertures such as those in Figure 3.1a add to the complexity of the treatment, making it more challenging to calculate dose deposition coefficients. This possible complexity of aperture shapes, the speed and acceleration of MLC leaves, and the tongue and groove effect, among other factors, can potentially produce errors in the calculation of the dose that is actually delivered to the patient (*Ezzell et al., 2003; Fredh et al., 2013; Heilemann et al., 2013; Hwang et al., 2014; Park et al., 2015a,b*). Our goal is to develop improved optimization-based treatment planning methods that reduce the discrepancy between the dose that is *planned* to be delivered and the dose that is *actually* delivered to the pa-



(a) Example of an irregular aperture shape.

(b) A desirable rounded and large shape.

Figure 3.1: Comparison of an irregular and a rounded aperture shape.

tient, by directly incorporating a metric related to complexity of aperture shapes into the objective function.

Excessive aperture shape complexity results in decreased dosimetric accuracy (*Bush et al., 2010*) and requires higher MLC positioning precision (*Das et al., 2008; Oliver et al., 2010*). Small aperture sizes correlate to irregular aperture shapes; *Fog et al. (2011)* reports that for such apertures, maximum dose and the overall width of the penumbra were underestimated by wide margins. Moreover, the problem is more pronounced in VMAT than it is in IMRT (*Du et al., 2014*), making complexity in VMAT treatments even more relevant.

Several complexity measures attempt to predict individual aperture dose accuracy: The Modulation Complexity Score (MCS) proposed by *McNiven et al. (2010)*, and the MCS applied to VMAT (MCSv) proposed by *Masi et al. (2013)*, consolidate Treatment Planning System (TPS) information such as leaf positions, aperture weight, field irregularity and area into a single score; unfortunately these metrics “perform poorly” in some particular sites (*McGarry et al., 2011*). *Du et al. (2014)* proposed an aperture irregularity (AI) metric calculated based on aperture area (AA) and aperture perimeter (AP): $AI = \frac{AP^2}{4\pi AA}$, which is similar to the Edge Metric (EM) proposed by *Younge et al. (2012)*, and to the “circumference/area” metric proposed by *Götstedt Julia Karlsson Hauer (2015)*. *Carlsson (2008)* proposes a different metric based on the ratio of differences of leaf overlaps.

In order to gauge the impact of aperture shape complexity on the final treatment quality, some of these metrics have already been tallied with actual QA results (*Agnew et al., 2014*). *McGarry et al. (2011)* and *Crowe et al. (2015)* compared their measures of complexity for anatomically different treatment sites with a corresponding QA measure. *Younge et al. (2012)* shows, via dosimetric validation of their EM measure, that it is a good predictor of dose calculation inaccuracies. We use this metric as our benchmark because their analysis suggests that “the majority of the error is concentrated on the edges of the apertures defined by the MLC leaves.”

The current claim is that: “Influencing the optimizer by integrating complexity metrics into the cost function has been little explored and requires more investigations” (*Chiavassa et al., 2019*). The few investigations of the matter so far include *Carlsson (2008)*, which used a “step-and-shoot” optimization that included their leaf overlap measure to generate round apertures that correlate well with the general circumference/area measures, and *Younge et al. (2012)*, which uses a local search to incrementally modify the aperture shape. A more straightforward approach includes limiting the beam intensity in order to force the creation of larger and rounder apertures (*Broderick et al., 2009*), but this presents a significant limitation to the treatment.

We propose an approach for integrating an aperture complexity metric into the objective function of an optimization-based treatment planning method by extending the VMAT treatment planning approach of *Peng et al. (2012)*. We will first define a modification of the edge metric and the resulting penalty function, and then explicitly incorporate this aperture shape penalty into the cost function. We will also adapt the column-generation heuristic of *Peng et al. (2012)* to our model, and finally, we will report some computational experiments and compare our edge metric to the one proposed by *Younge et al. (2012)*.

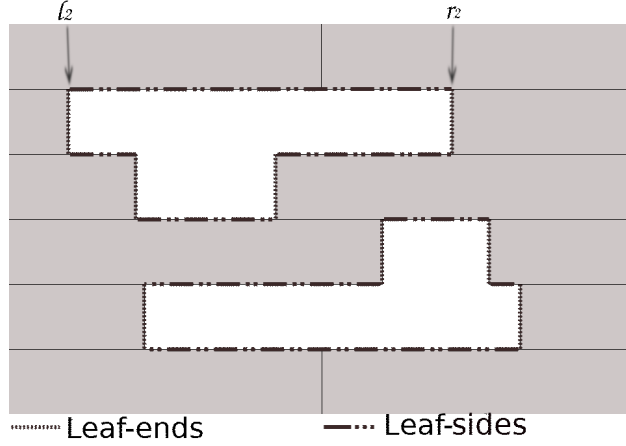


Figure 3.2: Example of an aperture illustrating vertical (stipples) and horizontal (dash-dot-dot) components of leaf edges contributing to the perimeter of the aperture.

3.2 Methods

3.2.1 A New Aperture-Edge Penalty

We will extend upon previous work on aperture penalization ([Younge et al., 2012](#)) and incorporate a comparable metric into our treatment planning approach. This aperture-edge metric will penalize excessive perimeter of the aperture relative to its area. Our goal is to include this aperture-edge metric in our optimization objective to explicitly control the aperture shape as part of treatment planning.

An aperture A in the multi-leaf collimator is determined by the configuration of the leaves. Let us denote the number of leaf pairs, or *rows*, in the MLC by M . We denote the position of the left and right leaves in row m by l_m and r_m , respectively, as shown in Figure 3.2, where $l_m, r_m \in [0, N]$, and N is the number of beamlets in each MLC row. The area of the aperture can be computed as

$$\text{area}(A) = b \sum_{m=1}^M (r_m - l_m), \quad (3.1)$$

where b is the width of each leaf. The aperture perimeter can be computed as the sum of the vertical

(μ) and the horizontal (λ) leaf-edge components, as illustrated in Figure 3.2. In particular, we can calculate $\mu(A)$ as

$$\mu(A) = \sum_{m=1}^M 2b\tau_m, \quad (3.2)$$

where τ_m is a binary indicator, signifying whether the leaf pair in row m is open or closed. In turn, $\lambda(A)$ can be calculated based on the location of endpoints of left and right leaves within the aperture as the sum of all the row-by-row horizontal perimeter components:

$$\lambda(A) = (r_1 - l_1) + \sum_{m=2}^M (|l_m - l_{m-1}| + |r_m - r_{m-1}| - 2(l_{m-1} - r_m)^+ - 2(l_m - r_{m-1})^+) + (r_M - l_M), \quad (3.3)$$

where $a^+ = \max(0, a)$. This formula takes into account any possible interdigitation, i.e., situations where $r_m < l_{m+1}$ so that the right leaf in row m protrudes over the left leaf in row $m + 1$, or situations where $l_m > r_{m+1}$. Figure 3.2 shows an example of interdigitation between the third and fourth rows.

Younge et al. (2012) introduced their Edge Metric, or EM, to quantify, in their words, “the amount of ‘edge’ in the aperture,” namely,

$$P_Y(A) = \frac{\tilde{C}_1\mu(A) + \tilde{C}_2\lambda(A)}{\text{area}(A)}, \quad (3.4)$$

where \tilde{C}_1 and \tilde{C}_2 are nonnegative parameters separating the contribution of leaf ends and sides into two individual terms, which, according to the authors, allows the user to “tailor the penalty depending on where dose calculation errors are observed for individual apertures,” and a penalty term

$$C \sum_{k=1}^K W_k P_Y(A_k), \quad (3.5)$$

with appropriate weights W_k that remove the bias to the regularization of apertures with the lowest

monitor units, and overall weight C , which can be added to the overall optimization objective for treatment planning (here, $k = 1, \dots, K$ are the control points used in the treatment). The penalty (3.5) was then added to the dose-related cost function, and an in-house treatment planning system described in [Fraass et al. \(2012\)](#) was used to improve the combined objective function by performing local search on the positions of endpoints of individual leaves. For an appropriate selection of the scaling parameter C , the resulting plans for a representative paraspinal Stereotactic Body Radiation Therapy (SBRT) case showed an improvement in aperture shapes, along with a reduction in MUs, without significant decrease in the plan quality.

However, the mathematical structure of the edge metric P_Y , which is not additive by MLC row, precludes its use within a wider class of VMAT treatment planning algorithms. To address this, we propose a *modified edge metric*, defined as

$$P(A) = C_1\mu(A) + C_2\lambda(A) - C_3\text{area}(A), \quad (3.6)$$

where C_1 , C_2 , and C_3 are non-negative parameters that reflect the relative importance of the contributions of the corresponding terms. Similarly to $P_Y(A)$, $P(A)$ is an increasing function of the perimeter and a decreasing function of the area of the aperture (note, however, that it can take on positive or negative values). A penalty term defined similarly to equation (3.5) is additive by MLC row, and can be added to the treatment planning objective function as discussed in the following sections.

3.2.2 Treatment Planning Problem Formulation

For simplicity of presentation, we will focus the discussion in this section on single-arc treatments, and discuss modifications for a multi-arc treatment separately. We discretize the trajectory of the gantry into K control points; it is sometimes convenient for notational purposes to also consider a “dummy” control point $K + 1$ at the end of the trajectory. Let δ_k denote the angular distance

between consecutive control points k and $k + 1$ (in degrees, or degs). At each control point $k = 1, \dots, K$, we will specify an aperture, A_k , and a fluence rate y_k (in MU deg⁻¹). A VMAT treatment plan can be constructed by interpolating the leaf positions in the MLC between the control points so that they match the specified apertures when the gantry reaches each control point. As discussed in Chapter 1, we represent the dose distribution delivered to the patient by the vector z , where each component z_v represents the dose to voxel $v \in \mathcal{V}$.

We formulate the VMAT treatment planning problem with aperture shape penalties as the following optimization problem, referred to as the Master Problem:

$$\text{(MP) } \underset{\substack{y_k, A_k: k \in \{1, \dots, K\} \\ z_v: v \in \mathcal{V}}}{\text{minimize}} \quad F(z) + C \sum_{k=1}^K P(A_k) \delta_k y_k \quad (3.7a)$$

$$\text{subject to} \quad z_v = \sum_{k=1}^K D_{kv}(A_k) \delta_k y_k \quad v \in \mathcal{V} \quad (3.7b)$$

$$y_k \in [0, Y] \quad k = 1, \dots, K \quad (3.7c)$$

$$S \leq S_{k,k+1}^U(A_k, A_{k+1}) \quad k = 1, \dots, K \quad (3.7d)$$

$$A_k \in \mathcal{A} \quad k = 1, \dots, K. \quad (3.7e)$$

The objective function (3.7a) is a weighted sum of the function $F(z)$, which evaluates treatment quality based on the dose distribution z , and a penalty term associated with the shapes of apertures used at each control point. Similarly to equation (3.5), we multiply each individual term $P(A_k)$ by the weight equal to the fluence $y_k \delta_k$ (in MU) specified at control point k , because any delivery errors due to an irregular shape of the aperture will be exacerbated by high monitor units associated with the control point. The parameter $C \geq 0$ is a scaling parameter that balances the relative importance of the two penalty components in the overall objective function. We will refer to the function

$$\sum_{k=1}^K P(A_k) \delta_k y_k \quad (3.8)$$

as the *modified edge metric penalty*, and to its counterpart based on the metric $P_Y(\cdot)$, namely,

$$\sum_{k=1}^K P_Y(A_k) \delta_k y_k, \quad (3.9)$$

as the *edge metric penalty*.

The addition of the aperture shape penalty to objective function is the main distinction between the problem (3.7a)–(3.7e) and the master problem in [Peng et al. \(2012\)](#); [Peng \(2013\)](#), and it necessitates modifications in the solution approach explored in the rest of this section.

The “dose deposition coefficient” $D_{kv}(A_k)$ in constraints (3.7b) denotes the dose received by voxel v from aperture A_k at control point k at unit fluence. To calculate doses z_v , $v \in \mathcal{V}$, in constraints (3.7b) we rely on a commonly-used approximation in which it is assumed that the aperture A_k , the fluence rate y_k , and the coefficient $D_{kv}(A_k)$ remain constant as the gantry travels between control points k and $k + 1$. As long as the discretization of the gantry trajectory into control points is such that the angular distances δ_k are small, this approximation is sufficiently accurate ([Otto, 2008](#); [Zwan et al., 2016](#)).

Constraints (3.7c) indicate that the fluence rate at each control point is nonnegative and bounded above by Y ; the upper bound is a reflection of the delivery machine’s maximum dose rate and gantry rotation speed S (in $\text{deg } s^{-1}$).

Constraints (3.7e) indicate that each selected aperture belongs to the set \mathcal{A} of *deliverable* apertures in the MLC system. In particular, we will assume that any aperture with non-overlapping leaf positions in each row is deliverable, and, unless stated otherwise, there are no restrictions on interdigitation in the adjacent rows.

Finally, constraints (3.7d) enforce *compatibility* of apertures at adjacent control points. This concept is illustrated in Figure 3.3. The MLC system imposes an upper bound on the leaf travel speed within the collimator. Thus, for the leaves to be able to move from their positions in A_k at control point k to their positions in A_{k+1} at control point $k + 1$, the gantry travel time between

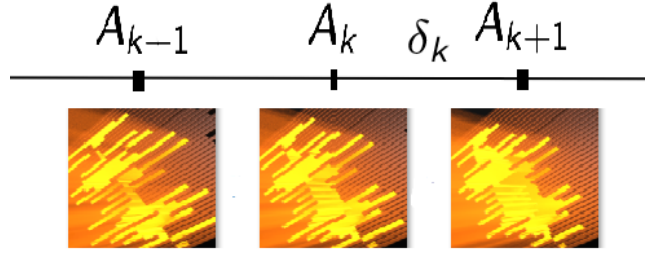


Figure 3.3: The compatibility constraint (3.7d) ensures that neighboring apertures remain “reachable” given the gantry rotation speed and constraints on the speed of movement of the MLC leaves.

control points k and $k + 1$ needs to be sufficiently long to allow the leaves to complete the required movement; equivalently, the gantry needs to travel sufficiently slowly. We follow [Peng et al. \(2012\)](#) and denote by $S_{k,k+1}^U(A_k, A_{k+1})$ the maximum gantry speed that would allow enough time for the leaves to complete the required motion; the apertures A_k at k and A_{k+1} at $k + 1$ are compatible if this speed is bounded below by the actual gantry rotation speed, S .

Note that, although in the formulation of (MP) we have assumed that the gantry travels at a constant speed, this assumption can be relaxed after the treatment plan — i.e., values of A_k and y_k for $k = 1, \dots, K$ — has been specified. In particular, it may be possible to increase gantry travel speed between some control points, and thus reduce treatment duration, as long as the source dose rate is adjusted accordingly to maintain the same fluence rate and the chosen speed is still bounded above by $S_{k,k+1}^U(A_k, A_{k+1})$ for each k . (See [Peng et al., 2012](#) for a detailed discussion.)

The problem (MP) is not convex, this means it is impossible to solve by conventional direct methods. The impact of the decision variables A_k , $k = 1, \dots, K$ on the objective and constraint functions is difficult to characterize. Therefore, in the remainder of this section we propose a *column-generation-based heuristic* to solve the problem approximately. Two optimization problems serve as the building blocks of this iterative heuristic: the restricted master problem, and the pricing problem. We motivate and formulate these problems, and describe methods for solving them, in the following two subsections, before discussing our main algorithm. While the moti-

vation behind both problems is similar to that of *Peng et al. (2012)*, we present their derivations for completeness, and note that the details of the pricing problem are more complex due to the presence of the shape penalty term in the master problem (3.7).

3.2.3 The Restricted Master Problem

Given a set of control points $\mathcal{C} \subseteq \{1, \dots, K\}$ and corresponding apertures \bar{A}_k , $k \in \mathcal{C}$, the restricted master problem ($\text{RMP}^{(\mathcal{C})}$) is designed to determine the optimal fluence rates y_k , $k \in \mathcal{C}$, associated with those apertures, assuming that $y_k = 0$ for $k \notin \mathcal{C}$:

$$\text{(RMP}^{(\mathcal{C})}) \quad \underset{\substack{y_k, k \in \mathcal{C} \\ z_v, v \in \mathcal{V}}}{\text{minimize}} \quad F(z) + C \sum_{k \in \mathcal{C}} P(\bar{A}_k) \delta_k y_k \quad (3.10a)$$

$$\text{subject to} \quad z_v = \sum_{k \in \mathcal{C}} D_{kv}(\bar{A}_k) \delta_k y_k \quad v \in \mathcal{V} \quad (3.10b)$$

$$y_k \in [0, Y] \quad k \in \mathcal{C}. \quad (3.10c)$$

Note that, since the apertures \bar{A}_k , $k \in \mathcal{C}$ are given as an input to ($\text{RMP}^{(\mathcal{C})}$), it is a continuous optimization problem in the variables (y, z) , and is a convex optimization problem as long as the function $F(z)$ is convex. Moreover, as long as the apertures \bar{A}_k , $k \in \mathcal{C}$ are deliverable and compatible with each other, any feasible solution of ($\text{RMP}^{(\mathcal{C})}$) corresponds to a deliverable VMAT treatment plan, where the apertures at each control point $k \notin \mathcal{C}$ can be defined by interpolating leaf positions at appropriate control points in \mathcal{C} and setting fluence rates to 0.

3.2.4 The Pricing Problem

Suppose $(\bar{y}_k, k \in \mathcal{C}; \bar{z})$ is an optimal solution to ($\text{RMP}^{(\mathcal{C})}$) with $\mathcal{C} \subsetneq \{1, \dots, K\}$ and \bar{A}_k , $k \in \mathcal{C}$. The goal of the pricing problem at the control point $c \notin \mathcal{C}$ is to determine whether the treatment can be improved by the addition of some aperture with a positive fluence at c .

3.2.4.1 Formulation of the Pricing Problem

To derive the objective function of the pricing problem we first consider a conceptual “intermediate” version of the master problem (cf. [Peng, 2013](#)):

$$\begin{aligned}
 \text{(MI)} \quad & \underset{z, y}{\text{minimize}} && F(z) + C \sum_{k \in \mathcal{C}} P(\bar{A}_k) \delta_k y_k + C \sum_{k \notin \mathcal{C}} \sum_{A \in \mathcal{A}} P(A) \delta_k y_{kA} \\
 & \text{subject to} && z_v = \sum_{k \in \mathcal{C}} D_{kv}(\bar{A}_k) \delta_k y_k + \sum_{k \notin \mathcal{C}} \sum_{A \in \mathcal{A}} D_{kv}(A) \delta_k y_{kA} && v \in \mathcal{V} && (\pi_v) \\
 & && y_k \geq 0 && k \in \mathcal{C} && (\rho_k) \\
 & && y_k \leq Y && k \in \mathcal{C} && (\gamma_k) \\
 & && \sum_{A \in \mathcal{A}} y_{kA} \leq Y && k \notin \mathcal{C} && (\gamma_k) \\
 & && y_{kA} \geq 0 && k \notin \mathcal{C}, A \in \mathcal{A}. && (\beta_k(A))
 \end{aligned} \tag{3.11}$$

In (MI), the apertures at control points $k \in \mathcal{C}$ are given and fixed, and for each control points $k \notin \mathcal{C}$, all possible apertures, and their associated fluence rates, are included in the problem, subject to an upper bound of Y on the total fluence rate at each control point. The optimal solution of (MI) is not intended to correspond to a deliverable VMAT treatment plan — in fact, this problem is not meant to be solved, but instead, the analysis of its optimality conditions will help us motivate the pricing problem.

Associating multipliers indicated in equation (3.11) in parentheses with constraints of (MI), we

can write the Lagrangian for this problem as

$$\begin{aligned}
\mathcal{L}(y_k, \gamma_k, \rho_k, y_{kA}, \gamma_k, \beta_k(A), z_v) = & F(z) + C \sum_{k \in \mathcal{C}} P(\bar{A}_k) \delta_k y_k + C \sum_{k \notin \mathcal{C}} \sum_{A \in \mathcal{A}} P(A) \delta_k y_{kA} + \\
& + \sum_{v \in \mathcal{V}} \pi_v \left(z_v - \sum_{k \in \mathcal{C}} D_{kv}(\bar{A}_k) \delta_k y_k - \sum_{k \notin \mathcal{C}} \sum_{A \in \mathcal{A}} D_{kv}(A) \delta_k y_{kA} \right) - \\
& - \sum_{k \in \mathcal{C}} \rho_k y_k + \sum_{k \in \mathcal{C}} \gamma_k (y_k - Y) + \sum_{k \notin \mathcal{C}} \gamma_k \left(\sum_{A \in \mathcal{A}} y_{kA} - Y \right) - \\
& - \sum_{\substack{k \notin \mathcal{C} \\ A \in \mathcal{A}}} \beta_k(A) y_{kA}.
\end{aligned} \tag{3.12}$$

First-order KKT conditions are necessary and, if $F(z)$ is convex, sufficient for optimality for (MI) (see, for instance, [Bazaraa et al., 2006](#)). These conditions are:

$$\frac{\partial \mathcal{L}}{\partial z_v} = \frac{\partial F(z)}{\partial z_v} + \pi_v = 0 \quad v \in \mathcal{V} \tag{3.13a}$$

$$\frac{\partial \mathcal{L}}{\partial y_k} = C \delta_k P(\bar{A}_k) - \sum_{v \in \mathcal{V}} \pi_v D_{kv}(\bar{A}_k) \delta_k - \rho_k + \gamma_k = 0 \quad k \in \mathcal{C} \tag{3.13b}$$

$$\gamma_k (y_k - Y^U) = 0 \quad k \in \mathcal{C} \tag{3.13c}$$

$$\rho_k y_k = 0 \quad k \in \mathcal{C} \tag{3.13d}$$

$$\gamma_k \geq 0 \quad k \in \mathcal{C} \tag{3.13e}$$

$$\rho_k \geq 0 \quad k \in \mathcal{C} \tag{3.13f}$$

$$\frac{\partial \mathcal{L}}{\partial y_{kA}} = C \delta_k P(A) - \sum_{v \in \mathcal{V}} \pi_v D_{kv}(A) \delta_k - \beta_k(A) + \gamma_k = 0 \quad k \notin \mathcal{C}, A \in \mathcal{A} \tag{3.13g}$$

$$\gamma_k \left(\sum_{A \in \mathcal{A}} y_{kA} - Y \right) = 0 \quad k \notin \mathcal{C} \tag{3.13h}$$

$$\beta_k(A) y_{kA} = 0 \quad k \notin \mathcal{C}, A \in \mathcal{A} \tag{3.13i}$$

$$\gamma_k \geq 0 \quad k \notin \mathcal{C} \tag{3.13j}$$

$$\beta_k(A) \geq 0 \quad k \notin \mathcal{C}, A \in \mathcal{A}, \quad (3.13k)$$

along with feasibility conditions for (MI).

Suppose $(\bar{y}_k, k \in \mathcal{C}; \bar{z})$ is an optimal solution to $(\text{RMP}^{(C)})$, along with associated KKT multipliers $\bar{\pi}$ and $\bar{\rho}_k, \bar{\gamma}_k, k \in \mathcal{C}$ (in particular, $\bar{\pi} = -\nabla F(\bar{z})$). Together, they automatically satisfy conditions (3.13a)–(3.13f). We can naturally extend this solution to a feasible solution to (MI) by setting $y_{kA} = 0$ for $k \notin \mathcal{C}, A \in \mathcal{A}$, and set $\gamma_k = 0$ for $k \notin \mathcal{C}$ in order to satisfy equations (3.13h) and (3.13j); note that constraint (3.13i) is satisfied automatically. It remains to verify whether conditions (3.13k) are satisfied for $\beta_k(A)$'s defined based on constraint (3.13g). If this is the case, the current solution is optimal for (MP). If not, there exists a control point $c \notin \mathcal{C}$ and an aperture $A \in \mathcal{A}$ such that

$$\beta_c(A) = C\delta_c P(A) - \sum_{v \in \mathcal{V}} \bar{\pi}_v D_{cv}(A) \delta_c \quad (3.14)$$

is negative. $\beta_c(A)$ can be interpreted as a *price*, or *marginal value*, of aperture $A \in \mathcal{A}$ at control point $c \notin \mathcal{C}$, and its negative value suggests that the objective value of (MI), and hence (MP), can be improved if the fluence rate associated with this aperture is increased from the current value of 0.

Motivated by the above discussion, the goal of the pricing problem at $c \notin \mathcal{C}$ is to minimize $\beta_c(A)$ given by equation (3.14) over all apertures $A \in \mathcal{A}$ that are compatible with the already-specified apertures $\bar{A}_k, k \in \mathcal{C}$. Following [Peng et al. \(2012\)](#), we ensure aperture compatibility by imposing the following constraints in the specification of the pricing problem (PP_c) at control point c :

$$(\text{PP}_c) \quad \beta_c^* = \min \quad \beta_c(A) \quad (3.15a)$$

$$\text{subject to} \quad A \in \mathcal{A} \quad (3.15b)$$

$$S \leq S_{c,c^+}^U(A, A_{c^+}) \quad (3.15c)$$

$$S \leq S_{c^-,c}^U(A_{c^-}, A). \quad (3.15d)$$

Here, c^- is the predecessor of c in \mathcal{C} , i.e., the control point in \mathcal{C} with the largest index smaller than c , and c^+ is the successor of c in \mathcal{C} , defined similarly. (If c has no predecessor and/or successor in \mathcal{C} , the corresponding constraint can be dropped from the formulation.)

To derive an explicit mathematical representation of the constraints of (3.15b), recall from Section 3.2.1 that any aperture A can be represented by vectors \vec{l} and \vec{r} , which are, respectively, vectors of positions of left and right leaves in the collimator forming this aperture. Assuming that interdigitation is allowed, constraint (3.15b) is equivalent to the linear inequalities

$$0 \leq l_m \leq r_m \leq N, \quad m = 1, \dots, M. \quad (3.16)$$

Constraints (3.15c) and (3.15d), which ensure compatibility of the aperture with the rest of the plan, can be interpreted as constraints on leaf positions, dictated by the gantry travel speed S and upper bound on the leaf travel speed v (in beamlets $\times s^{-1}$):

$$|l_{c^+m} - l_m| \leq \frac{v\delta_{cc^+}}{S}, \quad |r_{c^+m} - r_m| \leq \frac{v\delta_{cc^+}}{S}, \quad m = 1, \dots, M, \quad (3.17)$$

where $\delta_{cc^+} = \sum_{k=c}^{c^+-1} \delta_k$ is the angular distance between control points c and c^+ and $(\vec{l}_{c^+}, \vec{r}_{c^+})$ are the leaf positions in the aperture \bar{A}_{c^+} , and

$$|l_{c^-m} - l_m| \leq \frac{v\delta_{c^-c}}{S}, \quad |r_{c^-m} - r_m| \leq \frac{v\delta_{c^-c}}{S}, \quad m = 1, \dots, M, \quad (3.18)$$

where $\delta_{c^-c} = \sum_{k=c^-}^{c-1} \delta_k$ is the angular distance between control points c^- and c and $(\vec{l}_{c^-}, \vec{r}_{c^-})$ are the leaf positions in the aperture \bar{A}_{c^-} . (If c^+ or c^- are not defined for c , the corresponding bounds on l_m and r_m should be omitted.)

Constraints (3.16)–(3.18) can be combined to derive lower and upper bounds $(\underline{l}_m, \bar{l}_m)$ and $(\underline{r}_m, \bar{r}_m)$ on the positions of, respectively, the left and the right leaf in row m of the aperture at control point

c. We can therefore restate the pricing problem at control point c as

$$(PP_c) \quad \beta_c^* = \min \quad \beta_c(\vec{l}, \vec{r}) \quad (3.19a)$$

$$\text{subject to} \quad l_m \leq r_m \quad m = 1, \dots, M \quad (3.19b)$$

$$l_m \leq \tilde{l}_m \leq \bar{l}_m \quad m = 1, \dots, M \quad (3.19c)$$

$$\underline{r}_m \leq r_m \leq \bar{r}_m \quad m = 1, \dots, M. \quad (3.19d)$$

Note that, as long as apertures \bar{A}_{c^-} and \bar{A}_{c^+} at control points c^- and c^+ , respectively, are deliverable apertures and are compatible with each other, (PP_c) specified in problem (3.19) is feasible. Indeed, consider $(l_m, r_m) = (\tilde{l}_m, \tilde{r}_m)$ obtained by interpolating positions of corresponding leaves between control points c^- and c^+ , assuming that each leaf moves at a constant speed as the gantry travels between these control points:

$$\tilde{l}_m = l_{c^-m} + \frac{l_{c^+m} - l_{c^-m}}{\delta_{c^-c^+}} \cdot \delta_{c^-c}, \quad \tilde{r}_m = r_{c^-m} + \frac{r_{c^+m} - r_{c^-m}}{\delta_{c^-c^+}} \cdot \delta_{c^-c}, \quad (3.20)$$

where $\delta_{c^-c^+} = \sum_{k=c^-}^{c^+-1} \delta_k$ is the angular distance between control points c^- and c^+ . Then, if apertures at control points c^- and c^+ were chosen to be deliverable and compatible with each other, then leaf positions \tilde{l}_m and \tilde{r}_m satisfy appropriate constraints (3.19).

3.2.5 Solving the Pricing Problem

The pricing problem (PP_c) can be solved approximately using an approach that builds on the one first described in [Romeijn et al. \(2005\)](#) (however, the method described here is more involved since it needs to account for the added complexity due to the presence of the edge penalty $P(A) = P(\vec{l}, \vec{r})$ in the objective function). We will describe an approach for approximately solving (PP_c) by providing an approximate reformulation as a shortest path problem in a directed acyclic network, which can then be solved using, e.g., a Dynamic Programming algorithm. The refor-

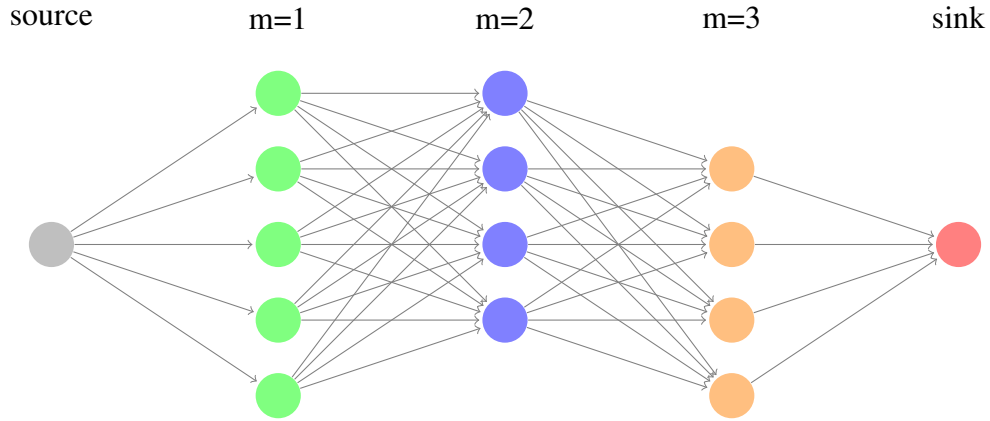


Figure 3.4: Illustration of a network constructed to solve the pricing problem for the case $M = 3$ rows in the MLC.

mulation takes advantage of the common approximation of the dose deposition coefficients by constant values within their beamlet domains, and the consequent piece-wise linear characteristics of the pricing problem’s objective function, and explores *most of* the critical values of the (PP_c) objective, corresponding to the breakpoints of such functions, while limiting the size of the underlying network to achieve a reasonable tradeoff between solution quality and computational effort required to achieve it.

3.2.5.1 Network Representation

The network consists of $M + 2$ layers. The first and the last layer contain one node each, called the source and the sink, respectively. The intermediate layers correspond to rows $m = 1, \dots, M$ of the MLC, and each node in layer m corresponds to a possible position of the left and right leaves, l_m and r_m , in that row; thus, we index each node by the triple (m, l_m, r_m) . (We will address the issue of which nodes, i.e., which combinations of leaf positions, should be included in each layer later on.) Each node is connected by a directed arc to each node in the subsequent layer. Figure 3.4 illustrates this construction for the case $M = 3$.

Note that every directed path from the source to the sink in this network passes through a se-

quence of nodes of the form $(1, l_1, r_1), (2, l_2, r_2), \dots, (M, l_M, r_M)$, and thus can be interpreted as an aperture.

We can associate costs with the arcs of this network so that the total cost of the arcs in each path from the source to the sink is equal to the value of $\beta_c(\cdot)$ associated with the corresponding aperture. Substituting equation (3.6) into (3.14), we can write

$$\beta_c(A) = \beta_c(\vec{l}, \vec{r}) = \delta_c \left[C(C_1\mu(\vec{l}, \vec{r}) + C_2\lambda(\vec{l}, \vec{r}) - C_3\text{area}(\vec{l}, \vec{r})) - \sum_{v \in \mathcal{V}} \bar{\pi}_v D_{cv}(\vec{l}, \vec{r}) \right],$$

where $\text{area}(\vec{l}, \vec{r})$, $\mu(\vec{l}, \vec{r})$, and $\lambda(\vec{l}, \vec{r})$ are defined by equations (3.1), (3.2), and (3.3), respectively.

Moreover, if we use the common representation

$$D_{cv}(\vec{l}, \vec{r}) = \sum_{m=1}^M D_{cmv}(l_m, r_m), \quad (3.21)$$

where $D_{cmv}(l_m, r_m)$ is the dose deposition coefficient for voxel v and control point c associated with “row aperture” (l_m, r_m) in row m (i.e., an aperture where all rows except for m are closed, and the left and right leaves in row m are positioned at l_m and r_m , respectively), we can represent the last term in the square brackets above as

$$- \sum_{v \in \mathcal{V}} \bar{\pi}_v D_{cv}(\vec{l}, \vec{r}) = - \sum_{m=1}^M \sum_{v \in \mathcal{V}} \bar{\pi}_v D_{cmv}(l_m, r_m). \quad (3.22)$$

Thus, we will define the costs of the arcs in the network as follows (we omit the multiplicative constant δ_c they all have in common):

- arcs between the source and nodes $(1, l_1, r_1)$ in layer 1:

$$C(2C_1b\tau_1 + C_2(r_1 - l_1) - C_3b(r_1 - l_1)) - \sum_{v \in \mathcal{V}} \bar{\pi}_v D_{c1v}(l_1, r_1) \quad (3.23)$$

- arcs between nodes $(m - 1, l_{m-1}, r_{m-1})$ and (m, l_m, r_m) , for $m = 2, \dots, M$:

$$C \left(C_2 (|l_m - l_{m-1}| + |r_m - r_{m-1}| - 2(l_{m-1} - r_m)^+ - 2(l_m - r_{m-1})^+) \right. \\ \left. + C(2C_1 b \tau_m - C_3 b (r_m - l_m)) - \sum_{v \in \mathcal{V}} \bar{\pi}_v D_{cmv}(l_m, r_m) \right) \quad (3.24)$$

- arcs between nodes (M, l_M, r_M) and the sink:

$$C(C_2(r_M - l_M)), \quad (3.25)$$

where $\tau_m = 1$ when $r_m > l_m$, and zero otherwise.

Finally, to calculate dose deposition coefficients associated with row apertures we used the common approximation (cf. [Peng et al., 2012](#))

$$D_{cmv}(l, r) = \int_l^r \phi_{cmv}(x) dx, \text{ where } \phi_{cmv}(x) = D_{cmnv}, \quad n - 1 < x \leq n; \quad n = 1, \dots, N, \quad (3.26)$$

where the beamlet dose deposition coefficients D_{cmnv} for every control point c , voxel v , MLC row m , and beamlet n were calculated using Varian's pencil-beam convolution-superposition algorithm (Varian Medical Systems, Inc., Palo Alto, CA, USA).

3.2.5.2 Selecting Nodes to Include in the Network

It remains to specify which nodes should be included in each layer of the network.

Consider the nodes in the layer corresponding to row m of the MLC. For every node (m, l_m, r_m) included in this layer, l_m and r_m must satisfy equations (3.19b)–(3.19d). The layer typically includes a node corresponding to each combination of l_m and r_m where l_m is equal to an integer in the interval $[\underline{l}_m, \bar{l}_m]$, and $r_m \geq l_m$ is equal to either l_m (if $r_m = l_m$ satisfies equation (3.19d)), or an integer in the interval $[\underline{r}_m, \bar{r}_m]$, i.e., it consists of nodes (m, l_m, r_m) for all combinations of values

l_m and r_m in the set

$$\begin{aligned} \{(l_m, r_m) : l_m \leq r_m, l_m \in \{\lceil \underline{l}_m \rceil, \lceil \underline{l}_m \rceil + 1, \dots, \lfloor \bar{l}_m \rfloor - 1, \lfloor \bar{l}_m \rfloor\}, \\ r_m \in \mathcal{L}_m \cup \{\lceil \underline{r}_m \rceil, \lceil \underline{r}_m \rceil + 1, \dots, \lfloor \bar{r}_m \rfloor - 1, \lfloor \bar{r}_m \rfloor\}\}, \end{aligned} \quad (3.27)$$

where $\mathcal{L}_m = \{l_m\}$ if $\underline{r}_m \leq l_m \leq \bar{r}_m$ and $\mathcal{L}_m = \emptyset$ otherwise. In other words, we solve the pricing problem approximately by considering only integer values of leaf positions, i.e., positions corresponding to the beamlet endpoints within each MLC row. In our computational experiments we occasionally had to make an exception to the above rule when one or both intervals specified in the set (3.27) contained no integer values. In this situation, we instead used interpolated values \tilde{l}_m and/or \tilde{r}_m defined in equations (3.20).

Unlike the pricing problem in [Peng et al. \(2012\)](#), even without interdigitation constraints, our pricing problem cannot be decomposed by row due to the dependence of the penalty metric $P(A)$ on the relative positions of the leaves in adjacent rows. Furthermore, the pricing problem could be solved exactly by considering a subset of feasible values of the variables, namely, those associated with (i) boundaries of constraints (3.19b)–(3.19c) and (ii) the breakpoints of the objective function, which is piece-wise linear when equation (3.26) is used to define functions $D_{cmv}(l_m, r_m)$. While in [Peng et al. \(2012\)](#) the breakpoints of the objective function of the pricing problem were limited to integer values of the variables, the breakpoints of our objective function are much more numerous, including all the combinations of values corresponding to breakpoints of the piece-wise linear function λ in equation (3.3). Therefore, we are effectively limited to approximate solutions of the pricing problems, which were based on integer variable values in our implementation.

Due to the layered structure of the network, we can solve the shortest path problem by forward dynamic programming, or *forward induction*. In a typical iteration, each node in layer m is labeled with the cost of the shortest path from the source to this node, $c(m, l_m, r_m)$, and the cost of the

shortest path from the source to each node in layer $m + 1$ is calculated as

$$c(m+1, l_{m+1}, r_{m+1}) = \min_{(m, l_m, r_m) \text{ in layer } m} \{c(m, l_m, r_m) + \text{arc cost}((m, l_m, r_m), (m+1, l_{m+1}, r_{m+1}))\},$$

with the arc costs defined in Section 3.2.5.1. (The *predecessor* nodes are also recorded, to enable the reconstruction of the shortest path.) Since each row may contain up to $\mathcal{O}(N^2)$ nodes, updating the cost labels requires $\mathcal{O}(N^4)$ operations per row. Thus, we can compute the shortest path in $\mathcal{O}(MN^4)$ operations. Note, however, that the above analysis does not account for the time required to calculate the arc costs given by the expressions in (3.23) and (3.24), which depends on the specifics of how the summations (in the last term) over voxels as well as calculations of dose deposition coefficients for row apertures are carried out.

3.2.5.3 Extensions and Generalizations

We conclude our discussion of the pricing problem by considering two extensions.

Control point-specific limits on leaf positions It may be desirable to replace bounds (3.16) with

$$L_m \leq l_m \leq r_m \leq R_m, \quad m = 1, \dots, M, \quad (3.28)$$

where $0 \leq L_m \leq R_m \leq N$ for $m = 1, \dots, M$, and replace constraints (3.19b) and (3.19c) with

$$\max\{\underline{l}_m, L_m\} \leq l_m \leq \bar{l}_m, \quad \underline{r}_m \leq r_m \leq \min\{\bar{r}_m, R_m\}, \quad m = 1, \dots, M. \quad (3.29)$$

For example, the dosimetrist may determine that making the aperture (at a particular control point) wider than the specified limits \vec{L} on the left and \vec{R} on the right is undesirable because the additional dose delivered to the targets is insufficient to justify the additional dose delivered to the healthy tissue. Aperture-edge penalties, by design, create a preference for large and round apertures over small and irregular ones. While the function $F(z)$ evaluating treatment quality aims to achieve an

appropriate tradeoff between treating the targets and sparing the OARs and other healthy tissues, balancing these tradeoffs with the aperture-edge penalty can present undue challenges for calibration of the objective function (3.7a). Instead, imposing explicit bounds on leaf positions (which are typically informed by the beam’s projections onto the targets) provides more direct controls on dose spillage to adjacent tissues. Furthermore, the reduction of the number of beamlets reduces computation time of solving the pricing problem.

While replacing bounds (3.19c) and (3.19d) with (3.29) may seem like a simple change to the pricing problem, if the values of (\vec{L}, \vec{R}) are control point-dependent, doing so may render the pricing problem infeasible. For example, it is possible that at control point c , for some m , the value of, say, L_m is too large for $l_m = L_m$ to satisfy the m th inequality in (3.17) and/or (3.18); i.e., the left leaf in row m cannot travel sufficiently fast to reach position L_m at control point c starting from position l_{c-m} at control point c^- and/or to reach position l_{c+m} at control point c^+ starting from position L_m at control point c ; a similar phenomenon occurs if R_m is too small. If this is the case, (PP_c) will be infeasible, and layer m in the network constructed based on bounds in (3.29) according to the rules similar to those summarized in the set (3.27) will be empty.

In our implementation, we attempt to respect the control point-specific limits on leaf positions to the extent possible, violating them only when faced with an infeasible pricing problem. In particular, when $\bar{l}_m < L_m$, we relax the limit and include nodes with $l_m = \bar{l}_m$ in layer m of the network. Similarly, when $\underline{r}_m < R_m$, we include nodes with $r_m = \underline{r}_m$. These allowances allowed us to prevent infeasibility of all the pricing problem instances encountered in our experiments.

MLC with no interdigitation If the MLC used in the VMAT delivery system does not allow interdigitation, the pricing problem can be modified to enforce this restriction. In particular, in the network representation of the pricing problem, the arcs between nodes $(m-1, l_{m-1}, r_{m-1})$ and (m, l_m, r_m) should be removed whenever $r_{m-1} < l_m$ or $l_{m-1} > r_m$.

Note that, without interdigitation, the “positive part” terms in the definition (3.3) of λ are always equal to zero, thus simplifying the structure of the cost function of the pricing problem.

3.2.6 Aperture Selection/Refinement Heuristic for (MP)

In this section we describe a heuristic approach for solving (MP). The heuristic is inspired by the algorithms described in *Peng et al. (2012)* and *Peng et al. (2015)*, with an enhanced approach to aperture refinement in its second phase.

3.2.6.1 Initial Aperture Selection via Column Generation

The first phase of the algorithm produces a deliverable VMAT treatment plan by using a column generation-like procedure to perform initial selection of apertures; this phase of the algorithm follows the framework of *Peng et al. (2012)*. In particular, we start with an “empty” treatment plan with $\mathcal{C} = \emptyset$ and $\bar{z} = 0$. At each iteration that follows, we first formulate and solve an instance of the pricing problem (PP_c) at each control point $c \notin \mathcal{C}$, given the current value of \bar{z} and the already-specified apertures at control points in \mathcal{C} . If the minimum of the objective values $\beta_c(A)$ of solutions found for all (PP_c)’s is negative, the aperture with the smallest objective value is added at the corresponding control point, the control point is added to \mathcal{C} , the restricted master problem (RMP^(C)) is re-solved with the updated set of apertures to obtain \bar{z} , and the algorithm proceeds to the next iteration. This process continues until $|\mathcal{C}| = K$, i.e., we have obtained a plan with an aperture specified at every control point, or the objective values at solutions of all (PP_c)’s are nonnegative, i.e., we are unable to identify an improving aperture to add to the current plan.

The formal statement of this procedure is provided in Algorithm 1.

3.2.6.2 Aperture Refinement

The process for initial aperture selection described in Section 3.2.6.1 is a greedy heuristic: at each iteration, it chooses to add an aperture that stands to improve the quality of the treatment plan specified in the previous iteration the most, based on local derivative information reflected by the KKT conditions. Due to its greedy nature, the apertures selected by the procedure in its early

Algorithm 1 Initial aperture selection

```
1: procedure INITIAL APERTURE SELECTION
2:   Set  $\mathcal{C} = \emptyset$  and  $\bar{z} = 0$ 
3:   while  $|\mathcal{C}| < K$  do
4:     For each  $c \notin \mathcal{C}$ , formulate and solve (PPc) with  $\bar{z}$  and  $\bar{A}_k$ ,  $k \in \mathcal{C}$ , to find  $\beta_c^*$ 
5:     Find  $\beta^* \leftarrow \min_{c \notin \mathcal{C}} \beta_c^*$ 
6:     if  $\beta^* \geq 0$  then
7:       Exit while loop
8:     else
9:       Let  $\bar{c} \leftarrow \operatorname{argmin}_{c \notin \mathcal{C}} \beta_c^*$  and let  $\bar{A}_{\bar{c}}$  be the corresponding solution to (PP $\bar{c}$ )
10:       $\mathcal{C} \leftarrow \mathcal{C} \cup \bar{c}$ 
11:    end if
12:    Solve (RMP( $\mathcal{C}$ )) with  $\mathcal{C}$  and  $\bar{A}_k$ ,  $k \in \mathcal{C}$ ; let  $(\bar{y}, \bar{z})$  be the optimal solution found and  $W_0$ 
    — the optimal objective value
13:  end while
14:  If necessary, complete the treatment plan by identifying feasible apertures at  $k \notin \mathcal{C}$ , augmenting  $\mathcal{C}$ , and setting  $\bar{y}_k = 0$ 
15:  Return  $\mathcal{C}$ ,  $\bar{A}_k$ ,  $k \in \mathcal{C}$ ,  $\bar{y}$ , and  $W_0$ 
16: end procedure
```

iterations may no longer be beneficial for the ultimate treatment plan. Indeed, in the instances of (RMP^(\mathcal{C})) solved in later iterations of the algorithm, the optimal values of the fluence rate variables at control points populated early in the process are frequently equal to, or close to, 0.

In the second phase of our algorithm, we perform *aperture refinement*. In particular, in this phase we revisit control points on the gantry trajectory in a specified order, and consider whether the aperture at the control point can be replaced with a different one, which is better in the context of the current treatment plan. This process is also accomplished by solving instances of restricted master problems (RMP^(\mathcal{C})) and pricing problems (PP_c).

The aperture refinement procedure is formally presented in Algorithm 2. The input for the procedure consists of the output of the initial aperture selection procedure in Algorithm 1, namely, \mathcal{C} — the set of control points where apertures have been specified, \bar{A}_k , $k \in \mathcal{C}$ — the selected apertures, and W_0 — the optimal value of the corresponding restricted master problem instance.

We will use n as the iteration counter for the refinement procedure, and use W_n to keep track of

Algorithm 2 Aperture Refinement

```
1: procedure APERTURE REFINEMENT( $\mathcal{C}$ ;  $\bar{A}_k, k \in \mathcal{C}$ ;  $W_0$ )
2:    $n \leftarrow 0$ 
3:   repeat
4:      $n \leftarrow n + 1$ 
5:      $W_n \leftarrow W_{n-1}$ 
6:      $c \leftarrow 1$ 
7:     repeat
8:       Formulate and solve (RMP $^{(\mathcal{C} \setminus \{c\})}$ ) to find  $\bar{z}$ 
9:       Formulate and solve ( $PP_c$ ); store  $\beta_c^*$ 
10:       $c \leftarrow c + 1$ 
11:     until  $c \geq K + 1$ 
12:     Let  $\mathcal{C}'$  be the set of  $[1, \dots, K]$  sorted in increasing order by  $\beta_c^*$ 
13:      $i \leftarrow 0$ 
14:     for each control point  $c$  in the ordered set  $\mathcal{C}'$  do
15:       Formulate and solve (RMP $^{(\mathcal{C} \setminus \{c\})}$ ) to find  $\bar{z}$  and  $\tilde{W}$  — the optimal value
16:       Formulate and solve ( $PP_c$ ) to find  $\beta_c^*$  and  $A_c^*$ 
17:       if  $\beta_c^* \geq 0$  then
18:          $i \leftarrow i + 1$ 
19:         if  $\tilde{W} < W_n$  then
20:            $\mathcal{C} \leftarrow \mathcal{C} \setminus c$ 
21:            $W_n \leftarrow \tilde{W}$ 
22:         end if
23:       else
24:          $i \leftarrow 0$ 
25:         Formulate and solve (RMP $^{(\mathcal{C})}$ ) with  $A_c^*$  at  $c$ ; let  $W^*$  be its optimal value
26:         if  $W^* < W_n$  then
27:            $\mathcal{C} \leftarrow \mathcal{C} \cup c, \bar{A}_c \leftarrow A_c^*$ 
28:            $W_n \leftarrow W^*$ 
29:         end if
30:       end if
31:       if  $i = 5$  then
32:          $n \leftarrow n + 1$ , exit the for loop
33:       end if
34:     end for
35:     until  $\left| \frac{W_n - W_{n-1}}{W_{n-1}} \right| \leq \epsilon$ 
36:     If necessary, complete the treatment plan by identifying feasible apertures at  $k \notin \mathcal{C}$ , augmenting  $\mathcal{C}$ , and setting  $\bar{y}_k = 0$ 
37:     Return  $\bar{A}_k, k \in \mathcal{C}, \bar{y}$ , and  $W_n$ 
38: end procedure
```

the objective function value of the best plan found in this iteration.

In every iteration of the aperture refinement procedure, we first determine an ordering of control points $\{1, \dots, K\}$, and then proceed to examine the aperture at each control point, in the specified order, and consider replacing it with a different one.

Given a treatment plan specified by \bar{A}_c , $c \in \mathcal{C}$, we order the control points based on a prediction of potential improvement of the treatment plan from replacing their current apertures. For each $c = 1, \dots, K$, we first remove the aperture specified at this control point (if any) and solve the corresponding instance of the restricted master problem, namely, $(\text{RMP}^{(\mathcal{C} \setminus \{c\})})$. Then, using the optimal solution of $(\text{RMP}^{(\mathcal{C} \setminus \{c\})})$, we formulate and solve the corresponding instance of (PP_c) , and record the objective value found, β_c^* . We sort the control points in the increasing order of β_c^* 's. The ordering process is described in lines 2.7–2.12 of Algorithm 2, and produces \mathcal{C}' — an ordering of $\{1, \dots, K\}$.

Next, we move on to the refinement process described in lines 2.14–2.34 of the algorithm. We proceed in the order specified by \mathcal{C}' and consider whether we can add, remove, or replace the aperture at each control point in a way that improves the overall treatment plan. The process for doing this is somewhat similar to the one used in the ordering step above. Conceptually, if c is the control point currently under consideration, we remove the aperture \bar{A}_c (assuming one is specified) from the plan, solve the resulting instance of $(\text{RMP}^{(\mathcal{C} \setminus \{c\})})$, and use its optimal solution to formulate and solve an instance of (PP_c) . If the optimal value of (PP_c) is nonnegative, i.e., no beneficial aperture could be found, we remove the control point from \mathcal{C} . If the optimal value of (PP_c) is negative, we place the discovered aperture at control point c (note that this aperture may, in fact, be equal to \bar{A}_c , which was removed to formulate $(\text{RMP}^{(\mathcal{C} \setminus \{c\})})$), confirming that its presence is beneficial to the overall plan).

In practice, however, we proceed with more caution, due to the following considerations: (i) the objective value of an aperture in (PP_c) is a prediction of its contribution to the master problem based on local derivative information, but the actual contribution needs to be verified by solving

the restricted master problem, and (ii) our solution to (PP_c) is based on a discretization of its feasible region, and the “original” aperture \bar{A}_c , while feasible, might not be representable in the current discretization, and thus might not be re-discoverable. Below, we provide a more detailed description of the procedure.

Let the optimal value of $(RMP^{(C \setminus \{c\})})$ be \tilde{W} , and let A_c^* be the aperture found by solving (PP_c) , with objective value β_c^* .

Case 1: $\beta_c^* \geq 0$. In this case, the pricing problem fails to find a beneficial aperture at c . We use a counter i to keep track of the number of such control points encountered consecutively within the current iteration, i.e., the number of consecutive control points where we failed to discover an improving aperture by solving the pricing problem. In this case, we increment i . If the current treatment plan includes some aperture \bar{A}_c at this control point, this would suggest that the aperture should be removed since, apparently, it does not have a negative value of $\beta_c(\cdot)$. However, since the pricing problem was only solved approximately, before removing \bar{A}_c , we explicitly check whether doing so would improve the plan, i.e., whether $\tilde{W} < W_n$. If this is the case, we update \mathcal{C} by removing control point c and set $y_c = 0$ and $W_n = \tilde{W}$. This case is described in lines 2.17–2.22 of the algorithm.

Case 2: $\beta_c^* < 0$. In this case, the pricing problem indicates that the aperture A_c^* should be used at control point c , replacing the current aperture if one exists. We formulate and solve a new instance of the restricted master problem $(RMP^{(C)})$, which uses A_c^* at c , and denote its optimal objective value by W^* . If this new W^* is better (lower) than W_n , we update the aperture at this control point to A_c^* , and update the value of W_n . Otherwise, the aperture found via the pricing problem was ultimately not beneficial, and we keep the original aperture in the plan. This case is described in lines 2.23–2.29 of the algorithm.

Once the value of counter i reaches a threshold (we use 5 in our implementation), we conclude that the likelihood of finding an improving aperture at the remaining control points in the cur-

rent ordering sequence is small. Therefore, in this case we terminate the current iteration of the refinement process.

We terminate the algorithm when the relative improvement in the objective function value of the restricted master problem in the latest iteration fails to exceed a pre-specified threshold (line 2.35).

3.3 Experiments and Results

3.3.1 Test Cases and Implementation Details

We performed the experiments discussed in this section on a dataset consisting of clinical cases at body sites that are suitable for VMAT treatment and tend to generate small irregular apertures during treatment: Head and Neck, Brain, Lung, and Spine. The Head and Neck case was a single-arc coplanar version of the case originating from the CORT dataset (*Craft et al., 2014*), and the remaining cases were provided by our collaborators at the University of Michigan Hospital System Radiation Oncology department, who assisted us with case data acquisition, including specification of the gantry trajectory and the number and location of control points for each case. Beamlet dose deposition coefficients were calculated using Varian’s Analytical Anisotropic Algorithm (AAA) with a calculation model that is based on a 3D pencil-beam convolution-superposition algorithm that accounts for tissue heterogeneities (*Ulmer et al., 2005*).

We chose these cases and their discretization settings in part with the goal of testing the computational performance of our treatment planning algorithm, studying how it scales for problems with varying numbers of control points, beamlets, and voxels. These features are summarized in Table 3.1. All single-arc cases use a total of $K = 180$ control points equally spaced around the patient with $\delta_k = 2^\circ$, while the multi-arc version of the brain case includes a second semi-circular arc in a different plane, with 90 additional control points. The machine settings were assumed to follow the parameters from *Varian Medical Systems (2011)* TrueBeam specifications.

	Brain	Lung	Spine	Brain multi-arc	Head and Neck
# Voxels	8,094	1,500	30,815	8,094	251,893
# Beamlets	69,120	36,000	20,160	91,800	202,640
Beamlet size	5mm × 5mm	1cm × 1cm	5mm × 5mm	5mm × 5mm	5mm × 5mm
# Control points	180	180	180	270	180
# OARs	16	5	8	16	25
# Targets	2	1	2	2	6

Table 3.1: Summary of case sizes in VMAT experiments.

Control point-specific limits \vec{L}_k and \vec{R}_k on leaf positions (see Section 3.2.5.3) were generated for each of the cases. For the head and neck case, their values could be deduced from the information included as part of the CORT data set. For the other cases, we determined the limits as follows. For each control point k , we located all the beamlets that were deemed *necessary* according to the following calculation. First, we calculated the maximum dose deposition coefficient D_{kmnv} from any beamlet (m, n) at this control point to any voxel v in the targets. We characterized a beamlet as *necessary* if it delivers at least a certain (control point-specific) percentage $X_k\%$ of that maximum dose to any target. For each row m , we set the limits L_m and R_m at that control point to be the left-most and the right-most necessary beamlets in that row. (Thus, some non-necessary beamlets will also be included because they lie between two necessary beamlets.) The reason for making the cutoffs X_k in this calculation specific to each control point is that the lesions are rarely located in the center of the body, and the ranges of values of dose deposition coefficients to voxels in the targets are different at different beam angles. To determine control point-specific values of X_k , we used manual iteration and visual inspection of the resulting \vec{L}_k and \vec{R}_k , and corresponding limits on possible apertures, overlaid on top of a beam’s eye view of the targets.

Recall from Section 3.2.5.3 that we may occasionally have to violate the control point-specific limits on leaf positions to ensure feasibility of the pricing problem. As a final step, we computed the case-specific envelope $\vec{L} = \min_k \vec{L}_k$ and $\vec{R} = \max_k \vec{R}_k$. While the apertures generated in the process of solving the pricing problems may occasionally violate the control point-specific limits

in some rows, they will always be contained in this envelope. Thus, we can reduce the memory requirements of the algorithm by only retaining dose deposition coefficients for beamlets contained in this envelope.

In our computational experiments, we did occasionally encounter apertures that were outside of the control point-specific limits, but it happened very infrequently.

We used an objective function $F(z)$ with the same convex smooth piece-wise quadratic structure as described in Section 2.4.1, assigning structure-dependent coefficients in each component. These coefficients were determined using manual trial and error on a VMAT case with $C = 0$ (i.e., no aperture shape control), until a satisfactory plan was obtained.

We implemented our algorithms using Python v.3.7. In particular, the restricted master problem was solved using the SciPy package; its results were corroborated using Gurobi's quadratic programming solver. There was an opportunity to solve the instances of the pricing problem at different control points in parallel; these parallelizations were carried out under the *multiprocessing* library. The full implementation is available as the GitHub repository <https://github.com/wilmerhenao/VMATwPenCode>.

We ran most cases on a desktop powered by an Intel Core™ i7 processor at 3.50GHz, with four cores and 32 GB of RAM. Where parallelization was possible, our implementation took advantage of all eight threads.

3.3.2 Calibration of Parameters in Metric P

In [Younge et al. \(2012\)](#), the authors showed that the number of open MLC rows in an aperture did not have a strong correlation with aperture irregularity, and therefore used $\tilde{C}_1 = 0$ and $\tilde{C}_2 = 1$ in their edge metric $P_Y(A)$ of equation (3.4). Similarly, in our modified edge metric $P(A)$ of equation (3.6) we used $C_1 = 0$, which, after scaling, allowed us to re-write the expression for

$P(A)$ simply as

$$P(A) = \lambda(A) - \xi \text{area}(A). \quad (3.30)$$

The next step was to calibrate the value of the parameter ξ so that the resulting values of $P(A)$ are in reasonable agreement with values of the edge metric $P_Y(A)$ (with $\tilde{C}_1 = 0$ and $\tilde{C}_2 = 0$).

We used the lung case to perform the experiments that informed our ultimate choice of ξ . First, we set $\xi = 1$ and applied the first phase of our heuristic algorithm (i.e., the initial aperture selection phase, but not the aperture refinement phase) to the resulting instances of the Master Problem (3.7) for several values of the scaling parameter $C \geq 0$. The apertures in each generated treatment plan were used to produce a data set to fit the linear regression model

$$P_Y(A) = \lambda(A) - \xi \text{area}(A) + \text{error}, \quad (3.31)$$

minimizing the sum of squared errors over each set of apertures.

Using different values of C in the Master Problem resulted in different treatment plans with different apertures, which in turn led to different values of ξ attaining the optimal fit. However, we observed that all values of ξ in the range $[0.4, 0.8]$ provided a reasonably good choice for most sets of apertures. This result is in agreement with the results of *Götstedt Julia Karlsson Hauer (2015)*, who studied Pearson's correlations between several aperture metrics and found that most of the edge-related metrics positively correlate with each other. Based on these observations, we used $\xi = 0.75$ in our subsequent experiments. As an illustration, we applied the initial aperture selection phase of the algorithm to the problem instance with $\xi = 0.75$ and $C = 0.0001$, generating an aperture at each of the 180 control points. The values of the two metrics for each of these apertures are depicted in Figure 3.5. We observe that the relationship between the two metrics is heteroscedastic, that is, greater dispersion of the values of the edge metric P_Y is observed for apertures with larger values of the modified edge metric P . However, overall the relationship is nearly monotone, and a linear approximation provides a good fit: the red line in the figure

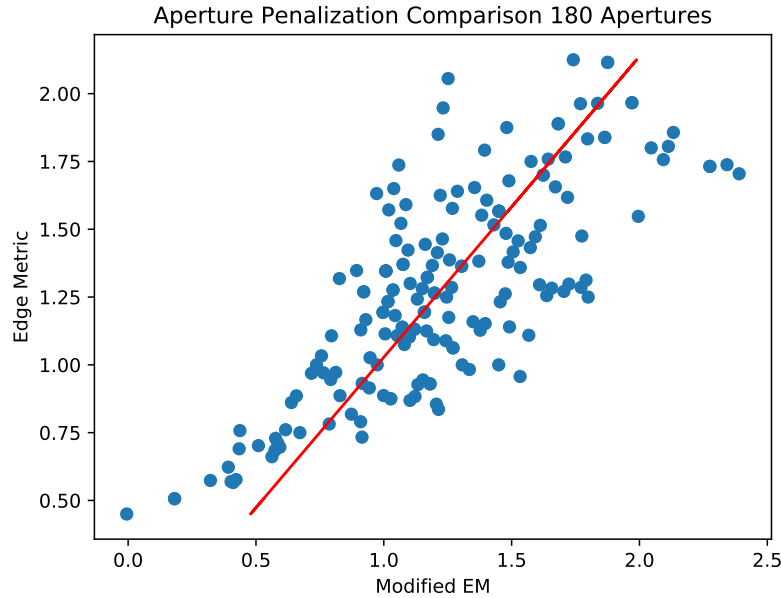


Figure 3.5: Comparison of values of $P(A)$ (horizontal axis) and $P_Y(A)$ (vertical axis) for 180 apertures in the plan obtained by applying the initial aperture selection phase of the algorithm to the lung case with $\xi = 0.75$ and $C = 0.0001$. Each dot associated with aperture A has coordinates $(P(A), P_Y(A))$. The red line represent the best linear fit for this data set, achieving $R^2 = 0.77$.

represents the best linear fit for this set of apertures, achieving $R^2 = 0.77$.

We performed additional experiments to confirm that, with the above choice of ξ , the values of the modified edge metric $P(A)$ provide a reasonable proxy for the values of the edge metric $P_Y(A)$. We calculated the values of both metrics for all apertures in plans produced for different values of C for the brain, spine, and lung cases after the initial aperture selection phase of the algorithm. The results are presented in Figure 3.6. The apertures from plans generated using different values of C are depicted using dots of different shapes and colors. Not surprisingly, plans generated using higher values of C tend to contain more regularly-shaped apertures, which tend to form clusters at the lower-left of the scatter plots, while apertures generated using lower values of C tend to cluster at the upper-right of the scatter plots. Combining the results for all the plans produced for each case in a single plot allows us to conclude that, although the relationship between the metrics

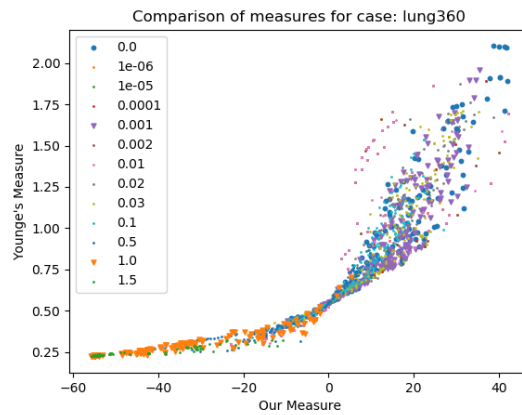
is not linear over the entire range of values of interest, the value of $P(A)$ provides a reasonable proxy for the value of $P_Y(A)$ for many different types of apertures, and the apertures with lower values of $P(A)$ should have better delivery characteristics, similarly to apertures with lower values of $P_Y(A)$, as studied in [Younge et al. \(2012\)](#).

With this hypothesis confirmed, we used $\xi = 0.75$ in the subsequent experiments, which included adding the aperture refinements phase to the algorithm. Although it is possible that a better fit can be obtained by selecting a case- or site-specific value of ξ , we kept its value constant for the purposes of this project, to focus our experiments on the impact of the scaling parameter C on the tradeoff between the two components of the objective function of the Master Problem.

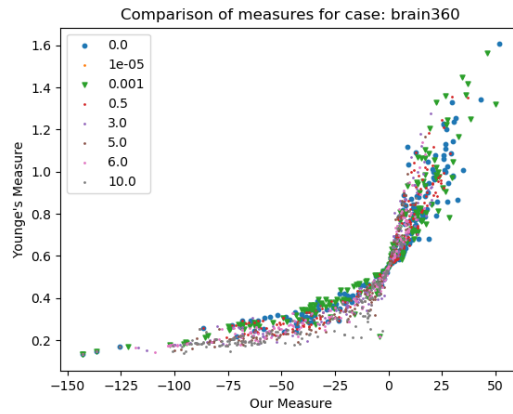
3.3.3 Aperture Refinement

We will use the spine case to illustrate the importance of the aperture refinement phase of the algorithm described in Section 3.2.6.2 by comparing three sets of DVH plots in Figure 3.7. (Once again, it is important to emphasize that all the DVH plots presented in this chapter are created using *planned* dose distributions of the proposed treatments, whereas *delivered* dose distributions will deviate from the planned ones. Moreover, since the treatments resulting from different planning and delivery paradigms will have different delivery characteristics, these dose discrepancies will be different as well; if our proposed approaches indeed produce plans that have better delivery characteristics, their delivered doses should adhere to the planned ones more closely.)

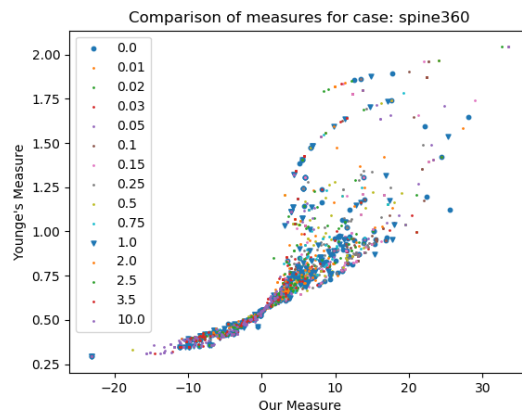
To create an illustrative “proof-of-concept” example to which the heuristic can be applied efficiently, we modified the spine case described in Table 3.1 by considering a subset of 36 control points uniformly distributed around the arc with $\delta_k = 10^\circ$ for all k : $\{0^\circ, 10^\circ, 20^\circ, \dots, 350^\circ\}$. The first DVH plot in Figure 3.7 corresponds to a hypothetical IMRT treatment obtained by solving an FMO model (where the intensity of each beamlet is determined individually) with 36 beam angles corresponding to the aforementioned VMAT control points, and objective function $F(z)$. Clearly,



(a) Lung case



(b) Brain case



(c) Spine case

Figure 3.6: Comparison of values of $P(A)$ (horizontal axis) and $P_Y(A)$ (vertical axis). Dots of different shapes and colors correspond to apertures in plans obtained with different values of C , as shown in the legends, after the initial aperture selection phase of the algorithm.

such a treatment cannot be delivered, but it provides us with a useful, if unattainable, benchmark, since the optimal objective value of this FMO model provides a lower bound on the value of $F(z)$ of any VMAT treatment with the same set of control points. While we cannot solve the optimization problem (3.7) exactly, we can compare the quality of dose distributions of its approximate solutions by how closely they approach the dose distribution of this IMRT benchmark. The second DVH plot in Figure 3.7 corresponds to the plan obtained after the initial aperture selection phase (Algorithm 1) with $C = 0.0001$, and the third plot corresponds to the plan obtained after several subsequent iterations of the aperture refinement phase (Algorithm 2). While the last plot is not quite as good as the FMO benchmark, it represents a significant improvement over the second one.

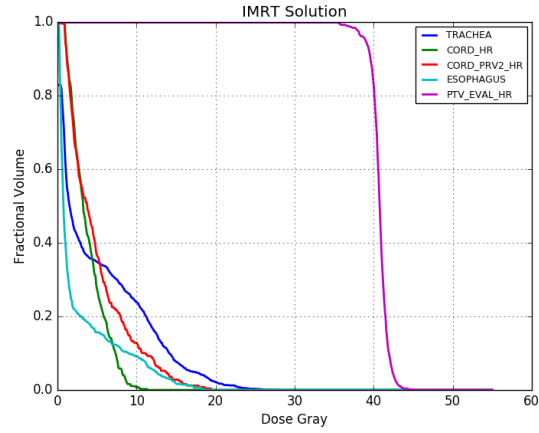
(It should be noted that the parameters of the function $F(z)$ used in this experiment are different than those in the forthcoming study of the spine case in Section 3.3.5, and thus the DVH plots in that section should not be directly compared to Figure 3.7.)

3.3.4 The Lung Case

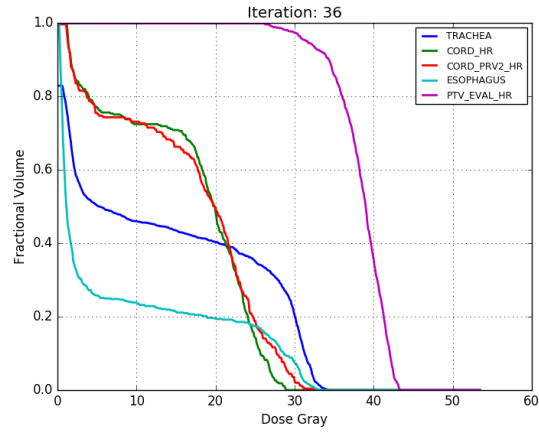
The lung case is the smallest (voxel-wise) case we considered, containing only 1,500 voxels. It is the most straightforward case when it comes to computational times.

In this case, we used a simpler version of the refinement algorithm. We did not control for strict improvement in the objective function value in each aperture refinement. Having said that, objective value improvement was attained after the application of the aperture refinement phase.

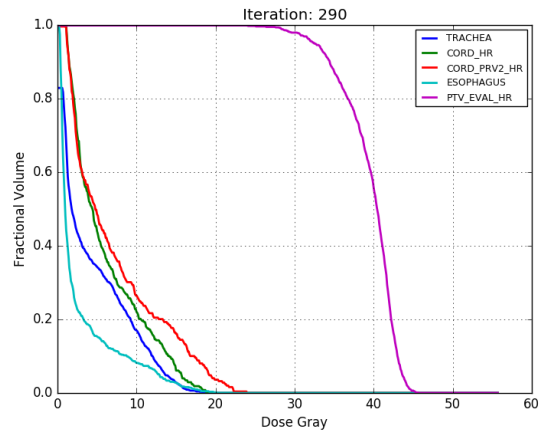
Our clinical goals for this case, defined according to the RTOG protocols, are outlined in Table 3.2. We can satisfy most of these objectives using our heuristic with $C = 0$ and $C = 1$. The DVH results shown in Figure 3.8 confirm the satisfaction of most of these goals with the only exception of the target (PTV) hot spot. The main limiting organ in this case is the spinal cord, but a steep penalization near the 45Gy threshold in the $F(z)$ function achieves the desired maximum dose. Esophagus, lung, and heart achieve doses below their goals.



(a) Benchmark FMO plan



(b) VMAT treatment without aperture refinement



(c) VMAT treatment after aperture refinement

Figure 3.7: DVH plots of plans for the spine case with 36 equispaced control points around the circular arc. (a) Benchmark FMO plan; (b) Output of the initial aperture selection phase with scaling parameter $C = 0.0001$; (c) Output of several subsequent iterations of the aperture refinement phase.

Goal	Delivered
Target Dose	60 Gy
Target Hot Spot	99% below 66 Gy
Target Cold Spot	99% above 55 Gy
Cord	100% below 45 Gy
Heart	50% below 20 Gy
Lungs	50% below 20 Gy
Esophagus	99% below 60 Gy
Esophagus	50% below 34 Gy

Table 3.2: Treatment goals for the lung case.

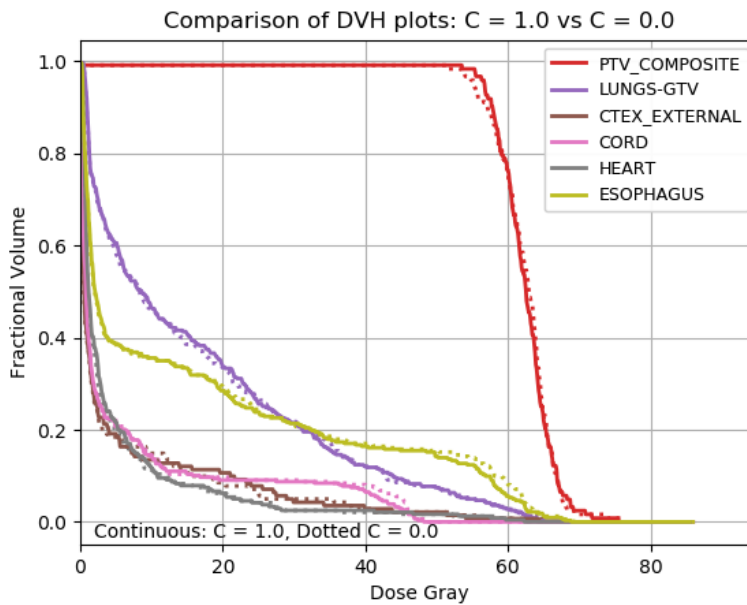


Figure 3.8: Comparison of Dose Volume Histograms (DVHs) for the lung case plans obtained with scaling parameter $C = 0$, i.e., without aperture shape penalty (dotted lines) and with $C = 1$ (solid lines). The structures shown correspond to the target and the most important OARs.

C = 0.0	C = 1.0
2581.50	1508.72
5949.03	4296.96
7828.04	4737.23
5701.68	4801.40
5585.59	4820.94
4995.61	4900.32

Table 3.3: The times (in seconds) spent in the initial aperture generation phase (first row) and in each “pass” through the aperture refinement loop (subsequent rows) for the lung case.

The plans for the lung case can be obtained in a few hours of computation, with most of this time spent on the aperture refinements. The times spent in the initial aperture generation phase and each “pass” through the refinement loop are reported in Table 3.3. The first phase runs faster than the subsequent aperture refinement passes; we attribute this to the greater parallelization opportunities in the first phase. Curiously, in this particular case the algorithm runs faster when $C = 1$; this is not indicative of a general pattern (in fact, the pricing problem solution subroutine can be implemented much more efficiently if we were only planning to consider instances with $C = 0$).

By increasing the value of the scaling parameter C , we can reduce the edge metric penalty $\sum_{k=1}^K P_Y(A_k)\delta_k y_k$ for the lung case by as much as 15%. The top plot of Figure 3.9 shows the values of this penalty for plans obtained by applying the heuristic algorithm to instances with different values of the scaling parameter C . The plots are scaled by the edge metric penalty corresponding to the plan obtained with $C = 0$. The edge metric penalty does not strictly decrease as the scaling parameter C increases, but the correlation between this penalty and the scaling parameter is negative. Figure 3.10 presents a plot of the values of the modified edge metric penalty $\sum_{k=1}^K P(A_k)\delta_k y_k$ versus the scaling parameter C . There is a stronger negative correlation between this penalty and the scaling parameter C , which is to be expected since this penalty is explicitly included in the objective function of the optimization problem (3.7). The dependence is not strict, which can be attributed to the fact that the plans are obtained by a heuristic algorithm, which is not guaranteed

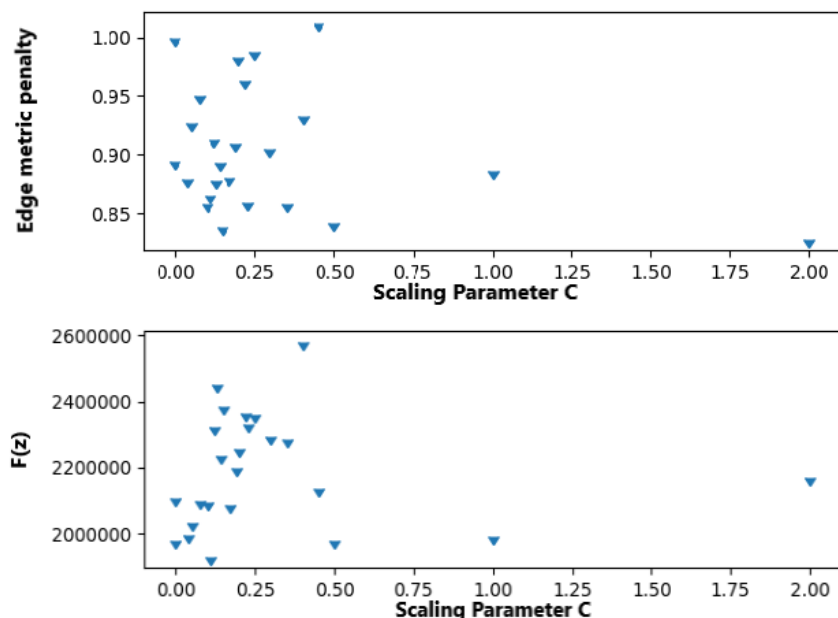


Figure 3.9: The edge metric penalty $\sum_{k=1}^K P_Y(A_k)\delta_k y_k$ (top) and the quality function $F(z)$ (bottom) of treatment plans for the lung case obtained for different values of the scaling parameter C . In the top graph, the values are scaled by the edge metric penalty corresponding to the plan obtained with $C = 0$.

to achieve global optimality.

Increasing the value of C shifts the focus of optimization away from the treatment quality penalty function $F(z)$. As shown in the bottom plot of Figure 3.9, the values of this function tend to increase as we increase the scaling parameter C (again, the relation is not monotonic, but the correlation is positive).

In Figure 3.11, we compare representative apertures from plans obtained with values $C = 0.0$ and $C = 1.0$ in the left and the right column, respectively. We show these results at 5 of the control points (one in each row). The heat maps in these pictures show beamlets at each control point represented by squares of different colors. The color of each beamlet represents the maximum dose deposition coefficient from this beamlet to any voxel in the targets using a dark-blue to bright-yellow spectrum: the darkest blue beamlets don't deliver dose to any target, and the brightest

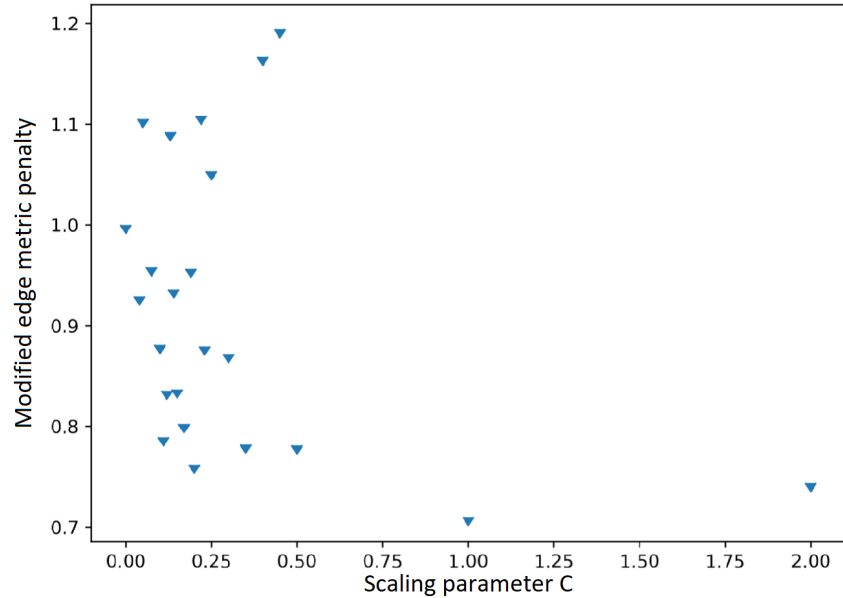


Figure 3.10: The modified edge metric penalty $\sum_{k=1}^K P(A_k)\delta_k y_k$ of treatment plans for the lung case obtained for different values of C . The values are scaled by the modified edge metric penalty corresponding to the plan obtained with $C = 0$.

yellow beamlets have the largest dose deposition coefficient to any target voxel from that control point. The aperture contours are shown by the green outlines. As demonstrated in the figure, a higher value of the scaling parameter C leads to a plan with apertures that have visually more rounded shapes (the five apertures on the right).

3.3.5 The Spine Case

The spine case has a total of 180 control points, comprising 20,160 beamlets. There are 30,815 voxels grouped into ten structures. The significant limiting factor for spine cases is the cord, which is surrounded by the lesion. Treatments for these types of lesions are characterized by small aperture fields and the tendency to produce irregular aperture shapes. The clinical goals for this case are shown in Table 3.4.

Figure 3.12 shows the values of the edge metric penalty (on the top) and the modified edge

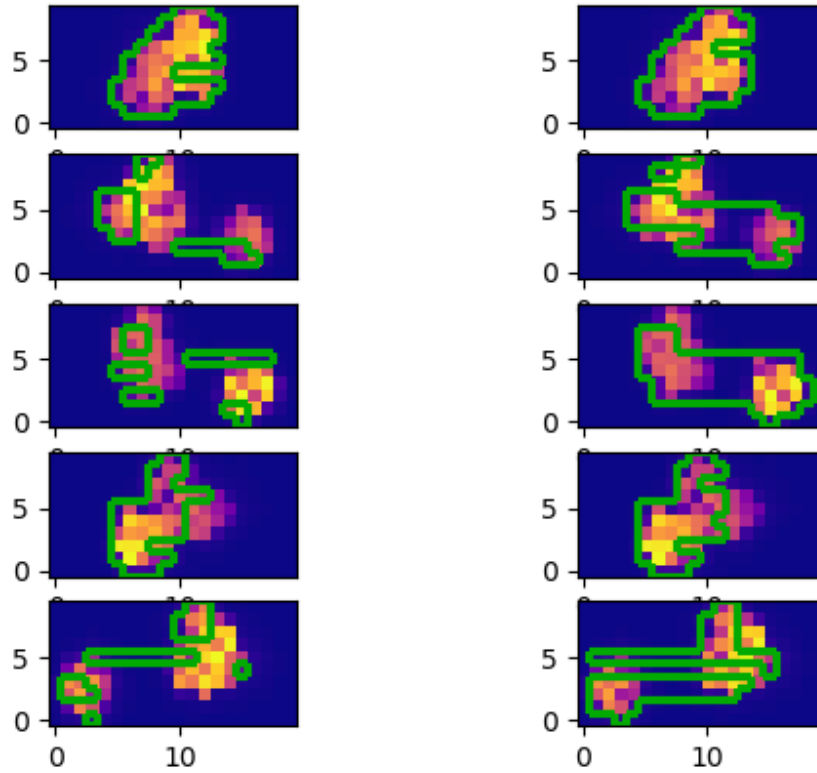


Figure 3.11: Comparison of aperture shapes in lung case plans obtained using values $C = 0.0$ (left column) and $C = 1.0$ (right column) at 5 of the control points (one in each row). Beamlets are shown as squares of different colors representing the maximum dose deposition coefficient from this beamlet to any voxel in the targets using a dark-blue to bright-yellow spectrum. The darkest blue beamlets don't deliver dose to any target, and the brightest yellow beamlets have the largest dose deposition coefficient to any target voxel from that control point. The aperture contours are shown by green outlines.

Goal	Delivered
Target Dose	44 Gy
Target Hot Spot	99% below 46 Gy
Target Cold Spot	99% above 37 Gy
Cord	100% below 30 Gy
Esophagus	95% below 40 Gy

Table 3.4: Treatment goals for the spine case.

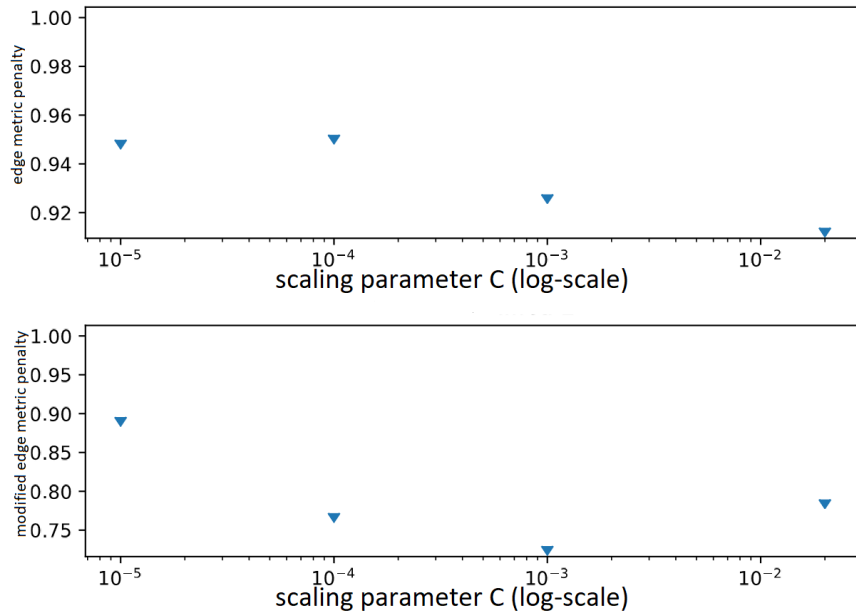
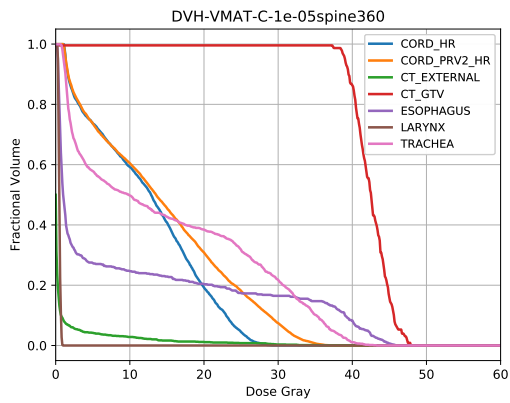


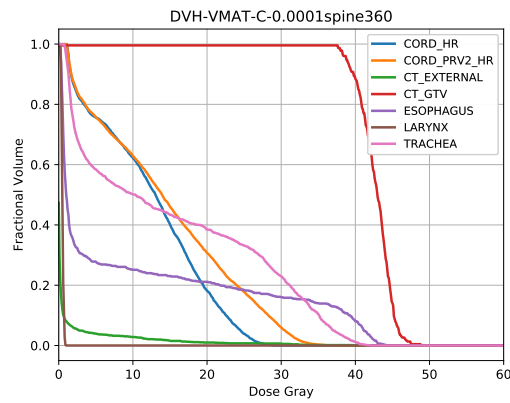
Figure 3.12: The edge metric penalty (top) and the modified edge metric penalty (bottom) of treatment plans for the spine case obtained for different values of the scaling parameter C . The values in each graph are scaled by the value of the corresponding penalty associated with the plan obtained with $C = 0$. The horizontal axis of both plots uses a logarithmic scale.

metric penalty (on the bottom) for the spine case obtained by applying the heuristic algorithm to instances with different values of the scaling parameter C . (Note that both plots use logarithmic scale for the horizontal axis.) As expected, these plots show a downward trend in the penalty values. The DVH plots of the resulting plans are presented in Figure 3.13. Increasing the value of C slightly deteriorates the quality of the treatment plans; however for the chosen values of C , the plans still satisfy the treatment goals.

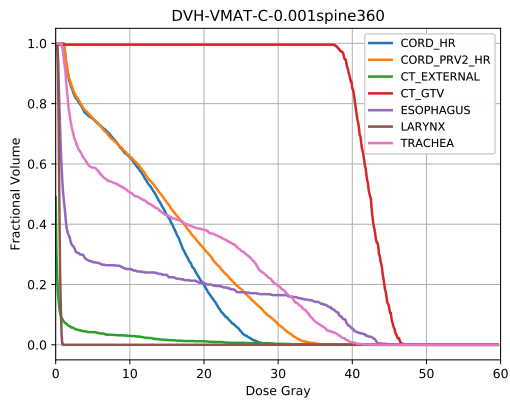
Several values of the scaling parameter C proved effective for reducing edge metric penalties in this case. For example, $C = 0.02$ may be a good choice to achieve a satisfying tradeoff: it reduces the edge metric penalty by more than 8% (while reducing the modified edge metric penalty by 21.1%), while the resulting treatment plan satisfies the goals with only small violations. In fact, some aspects of the treatment plan obtained with $C = 0.02$ improved over the plan obtained with



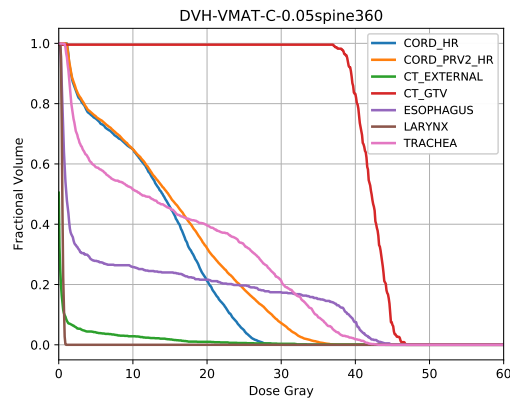
(a) $C = 0.00001$



(b) $C = 0.0001$



(c) $C = 0.001$



(d) $C = 0.05$

Figure 3.13: Comparison of DVH plots for the spine case plans obtained for different values of the scaling parameter C .

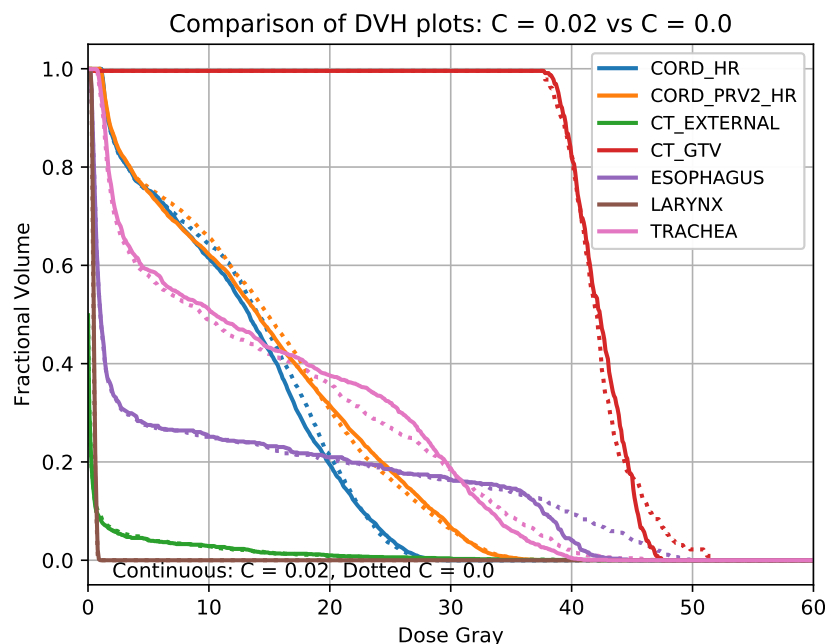


Figure 3.14: Comparison of DVH plots for the spine case plans obtained with scaling parameter $C = 0$, i.e., without aperture shape penalty (dotted lines) and with $C = 0.02$ (solid lines). The structures shown correspond to the target and the most important OARs.

$C = 0$ (notice, for example, the reduction in maximum dose to the esophagus in Figure 3.14); however, this improvement may be attributable to the heuristic nature of the algorithm used to generate these plans.

3.3.6 The Head and Neck Case

The Head and Neck case was taken from the CORT dataset, which is an open dataset available to researchers for developing and comparing radiation therapy planning algorithms. The case contains a total of 1983 control points on a sphere around the target, to allow creation of non-coplanar treatment plans. In our study, we limited our attention to the 180 control points on a single coplanar VMAT arc orthogonal to the couch length. The case contains 251,893 voxels, separated into 25 structures, with 25,388 of the voxels corresponding to one of the 6 targets. We

Goal	Delivered
Target Dose	78 Gy
Spinal Cord	100% below 48 Gy
Parotid Glands (Left)	50% below 40 Gy
Parotid Glands (Right)	50% below 40 Gy
Left Optic Nerve	100% below 30 Gy
Right Optic Nerve	100% below 30 Gy

Table 3.5: Treatment goals for the head and neck case.

illustrate 6 of the OARs and 2 of the targets in Figure 3.15. These are the structures highlighted on Figure 3 of *Craft et al. (2014)*, and they are the critical structures for this particular case.

In this case, as in the lung case, we used a simpler version of the refinement algorithm, where we did not control for strict improvement in the objective function value in each aperture refinement. As in the lung case, objective value improvement was attained after the application of the aperture refinement phase.

The treatment goals for this case are outlined on Table 3.5. By calibrating the parameters in function $F(z)$, we were able to satisfy all of the goals, and by setting $C = 0.5$, we can reduce the edge metric penalty $\sum_{k=1}^K P_Y(A_k)\delta_k y_k$ by 19%. Figure 3.15 presents a comparisons of the DVH plots for plans obtained with $C = 0$ and $C = 0.5$. Some deterioration of treatment quality can be seen in the latter plan, although all treatment goals remain satisfied, with the exception of a slight violation of the goal for the right parotid. In a case such as this, the treatment planner would consider whether the plan obtained for the chosen value of C achieves the desired tradeoff between the planned treatment quality and edge metric penalty reduction, and re-run the algorithm on instances with different values of C to explore the range of the tradeoffs available. A recalibration of the function $F(z)$ (in particular, parameters associated with the terms corresponding to the right parotid voxels) can also be used to attempt to maintain the goal corresponding to this structure.

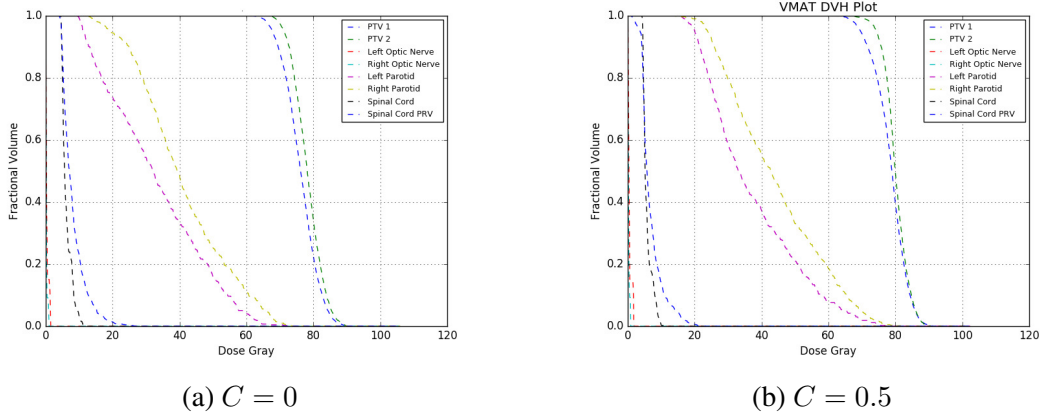


Figure 3.15: Comparison of DVH plots for the head and neck case plans obtained with scaling parameter $C = 0$, i.e., without aperture shape penalty (left) and with $C = 0.5$ (right). The structures shown correspond to the targets and the most important OARs.

3.4 Multi-Arc VMAT

It is possible to extend the VMAT treatment planning model and solution heuristic presented in Section 3.2 to a multi-arc setting by making the following modifications, and keeping track of the case data according to the following recommendations:

- Concatenate the arcs, but do not impose compatibility constraints (3.7d) between the last control point of one arc and the first control point of the next arc. Additionally, when solving the pricing problem, do not include bounds (3.17) or (3.18) if control points c^+ or c^- , respectively, belong to a different arc than c . Lastly, aperture envelope bounds (\vec{L}, \vec{R}) should also be determined separately for each arc.
- Make sure that the data generated for different arcs maintains consistency of indexing of voxels and structures.
- Make sure that in the data generated for different arcs, control points' and beamlets' indices are unique across all arcs.
- MLC sizes and discretizations may be different for different arcs, i.e., each arc will have

Goal	Delivered
Target Dose	60 Gy
Brainstem	100% below 60 Gy
Left Optic Nerve	100% below 55 Gy
Right Optic Nerve	100% below 55 Gy
Optic Chiasm	100% below 54 Gy
Left Eye	100% below 40 Gy
Right Eye	100% below 40 Gy
Left Lens	100% below 10 Gy
Right Lens	100% below 10 Gy

Table 3.6: Treatment goals for the brain case.

different values of parameters M and N associated with it. In this case, make sure that the upper bound on the leaf travel speed v (in beamlets $\times s^{-1}$) is calculated for each arc accordingly.

- If the arcs are independent (i.e., there are no voxels that have positive dose deposition coefficients for beamlets corresponding to control points in different arcs), Algorithm 1 can be applied to each arc separately and in parallel. Algorithm 2 can also be applied to each arc separately, or the arcs can be merged in this phase to speed up the process (namely, to enable combined updates of the counter i across different arcs).

We tested the multi-arc case extension on a brain case discussed in Subsection 3.4.1.

3.4.1 The Brain Case

This case is particularly challenging because of the proximity of one of the lesions to the right eye. The lesion surrounds the optic nerve, which severely complicates the treatment. The clinical goals for the brain case are shown on Table 3.6.

Due to the anatomical features of this case, it is not possible to achieve a good treatment with only one arc. Figure 3.16 shows a DVH plot of the single-arc plan obtained for $C = 0$; i.e., without

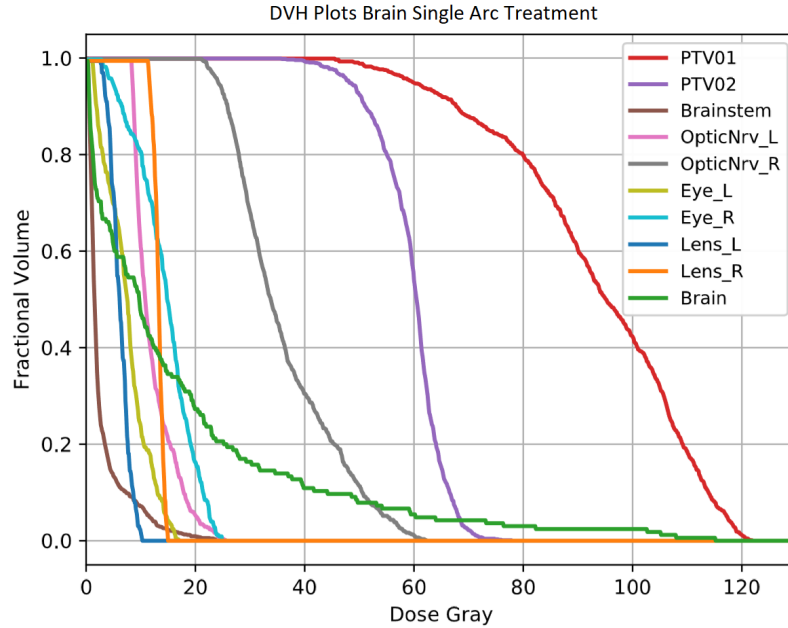


Figure 3.16: A DVH plot for the single-arc brain case plan obtained for $C = 0$, i.e., without aperture shape penalty. The structures shown correspond to the target and the most important OARs.

aperture shape penalty. In this experiment, the parameters in the function $F(z)$ were selected in an attempt to satisfy the upper bounds on all OARs to the extent possible while avoiding cold spots in both PTVs and hot spots in the second PTV. As a result, the first PTV has a significant hot spot. Many more combination of parameters of $F(z)$ were tried, all of them with different unsatisfactory results, with goals for at least one of the structures being significantly violated.

In order to achieve better results, a second arc was added. While the first arc has 180 control points equally spaced around a full circle orthogonal to the couch length, the second arc is a semi-circle perpendicular to the first arc, with 90 control points. With the addition of the second arc, significantly better treatment plans could be obtained. After re-calibrating the function $F(z)$, our algorithm (applied to the instance with $C = 0$) was able to find a treatment plan that satisfies most of the goals connected to the OARs associated with the left eye (lens, nerve, and the eye itself); see the dotted DVH curves in Figure 3.17. However, we were not able to satisfy the goals related to

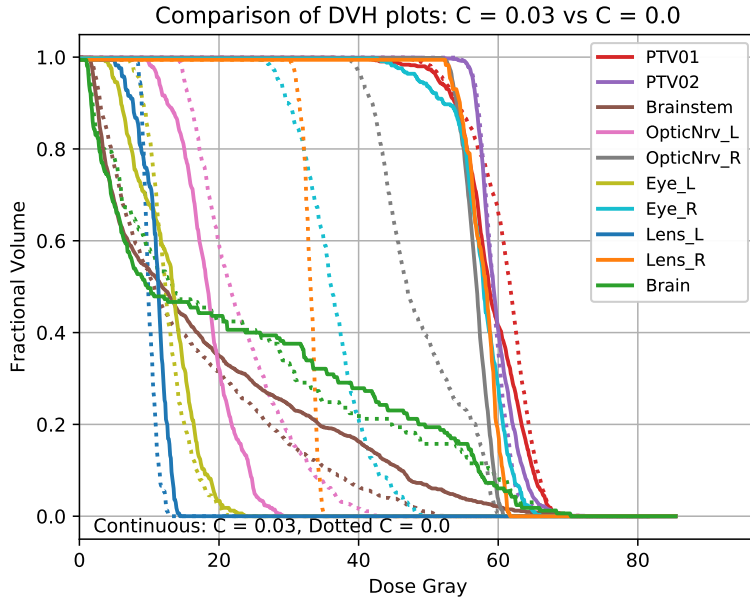


Figure 3.17: Comparison of Dose Volume Histograms (DVHs) for the brain case plans obtained with scaling parameter $C = 0$, i.e., without aperture shape penalty (dotted lines) and with $C = 0.03$ (solid lines). The structures shown correspond to the target and the most important OARs.

the right eye OARs. From the clinical point of view, in a case such as this, with the right eye being very close to one of the PTVs, the physician may consider sacrificing the function of the right eye for the sake of achieving adequate coverage of cancerous regions.

We tested the algorithm on an instance with a relatively large value of the scaling parameter $C = 0.03$. With this value of C , the aperture shape penalty plays a significant role, and the solid DVH curves in Figure 3.17 for the resulting plan look quite different. However, it should be noted that due to our choice of parameters in the function $F(z)$, the most significant deterioration is still limited to the OARs corresponding to the right eye, while the left eye and the brainstem remain fairly well protected and the coverage of the PTVs is not significantly affected. Therefore, the physician may in fact choose this plan, which sacrifices the right eye to meet the goals for the other structures, while simultaneously reducing the edge metric penalties and thus reducing the discrepancy between planned and delivered doses.

It should be noted that the increase in the number of control points resulted in significant increase in the computational demands of the heuristic. It took 8,861 seconds to execute the first phase, and 62,740 seconds to execute the first refinement loop (after which the algorithm was stopped). To enable clinical use of our algorithm and its variation, more aggressive use of parallelization will be required.

3.5 Translation to the Clinical System

Translation to the clinical system involves the transfer of the plan specified using the values of the variables of the optimization model of the treatment planning problem into the specifications used in a clinical setting. In our particular case, it involves the passing of the machine parameters corresponding to control point orientations, intensities, and positions to the LINAC Treatment Planning System. The “Beam Meterset” attribute that assigns intensity for each successive control point needs to be specified in “Monitor Units or minutes as defined by the Primary Dosimeter Unit (The measurement unit of machine dosimeter)” (as defined in the DICOM standard Information Object Definitions). Leaf positions are also passed to the LINAC at each discrete control point.

The information is transferred using a Digital Image and Communications in Medicine (DICOM) file (see <https://www.dicomlibrary.com>). The files can store treatment information, namely, machine parameters derived from the optimization, high-quality medical images (CT, MRI, and ultrasound results), and provides integration with multiple medical devices. One of these devices is the Varian LINAC. Varian’s Treatment Planning System (TPS) allows the calculation of dose effects, enables the visualization of the dose distribution, and allows for the visualization of the aperture field along the gantry trajectory. The TPS extrapolates a smooth trajectory from the provided discretization. In particular, the variables of the optimization model (3.7) can be converted into the input types of VMAT treatments in Varian’s TPS via a simple bijective transformation.

3.6 Expected Reduction of Dosimetric Discrepancies

It is difficult to provide a definitive forecast of dose discrepancies based on our results. Our experiments show that a significant reduction in the modified edge metric penalty is medically feasible for the cases we considered. However, at this stage in our research, we have to rely on the physical measurements performed by *Younge et al. (2012)* to predict the effect of greater regularity of aperture shapes on reductions in dosimetric discrepancies. As we have pointed out throughout our analysis, our modified edge metric provides a reasonable proxy for the edge metric of *Younge et al. (2012)*, but the relationship is not linear or fully monotone.

To recap some of the experimental results presented in *Younge et al. (2012)*, they achieved an aggressive reduction of 55% in edge metric penalization in their spine case. The next step was the experimental measurement of the resulting reduction in dose discrepancies using gamma analysis (a measure that combines dosimetric and spatial discrepancy regularly used in QA). For gamma analysis with 3% dose and 1mm distance, the percentage of voxels passing the test went from 79.5% to 95.4% when the aperture penalization was implemented, which is a significant improvement.

In our experiments, we didn't achieve such a significant reduction in the edge metric penalty, but we can still expect an increase in the percentage of voxels passing the QA tests. Moreover, due to the different mathematical structure of the original and modified edge metrics, the apertures obtained by using a penalty based on the latter may have subtly different geometric features and have a different effect on actual dose discrepancy and QA pass rates. Further computational experiments should be performed to more fully explore the clinically feasible tradeoffs potentially leading to greater reductions in the penalties, and the modified edge metric penalty should be evaluated dosimetrically in order to measure its actual impact.

3.7 Conclusions

We have developed a new optimization model for VMAT treatment planning that explicitly penalizes aperture shape irregularities, and can be used to study the tradeoffs between the quality of planned treatments and a new edge metric penalty. We developed a significant extension of a heuristic solution algorithm for VMAT treatment planning applicable to the new model. The edge metric penalty used in our model is shown to be a good proxy for the metric used in prior work by *Younge et al. (2012)*, which in turn was shown to correlate with measured discrepancies between planned and delivered doses; based on this analysis we predict that the plans obtained using our model and algorithm will have less dosimetric discrepancies between the planned and delivered doses.

Future research in this area should focus on the development of improved solution methods for the proposed treatment planning model, both in terms of algorithmic approaches leading to better heuristics, and implementation techniques to speed up the computation. Moreover, dosimetric studies should be used to verify the reduction in dosimetric errors and discrepancies achieved by using the modified edge metric. We do not rule out the possibility that our modified metric may convey somewhat different information about the aperture shape than the original metric of *Younge et al. (2012)*, and thus using a penalty based in this metric in the treatment planning optimization may have a qualitatively different impact on dosimetry.

We also believe that our modeling approach can be a useful tool in other applications where geometric properties, specifically, “excessive edges,” are an important aspect of design decisions. Gerrymandering in districting decisions is an excellent example of a field that could benefit from edge penalties.

CHAPTER 4

Conclusions and Future Research Suggestions

The rapid advances in the field of radiation therapy planning and the constant development of new machinery and technologies demand the constant evolution of models and methods for treatment planning to squeeze the most benefit out of the potential offered by the new technology. The newest machines require models that take into account the complex physical properties inherent to the treatment modality. Together, physicists, physicians, and operations research professionals work in tandem in order to devise the tools that can take the field of radiation therapy planning to the next level.

We attempted to devise optimization models and solution methods that allow the creation of treatment plans that are less prone to discrepancies between the dose distributions that are planned and that are actually delivered to the patients. If these methods prove to be clinically feasible, they can be used to design more precise treatment strategies, leading to treatments that are more effective and safer for the patients.

In closing, we mention a few promising future research directions.

4.1 Tomotherapy

We think tomotherapy delivery can be made better by increasing the projection resolution. Our proposed modeling paradigms avoid the limitations of current treatment planning methods in this

regard, paving the way for more precise representation of tomotherapy treatments. However, this flexibility comes at the cost of complexity of the associated MIP problems. We anticipate that future work should focus on improvement of the solution approaches for these MIP problem, possibly including design of specialized algorithms.

Moreover, the hypothesized improvements in dosimetric accuracy due to finer discretization and increased Leaf Opening Times (LOTs) should be tested at other treatment sites, both computationally and dosimetrically. While our experiments suggest that longer LOTs can be achieved without negative impact on (planned) treatment quality for typical prostate cases, for other treatment sites shorter LOTs may be unavoidable. In the latter case, increased precision in dose delivery modeling due to improved resolution may still have a positive effect on treatment planning.

4.2 VMAT

The first line of future research would investigate whether the proposed modified edge metric is actually a good predictor of dose discrepancies. Two potential complementary approaches are to perform experiments based on dosimetric measurements, similar to those done by *Younge et al. (2012)* and other researchers, as well as attempt to develop analytical approaches for estimating (e.g., providing lower and upper bounds) the dose discrepancies based on aperture shapes.

Another line of research to be explored is whether it is possible to either avoid high computational cost of solving the pricing problem, or leverage parallelization opportunities to a greater extent, to bring computational demands of our algorithms, and other similar solution methods, down to clinically reasonable levels. As more complex multi-arc VMAT treatments with greater numbers of control points are being considered, while computer systems with large multi-core processors and distributed capabilities become increasingly available, parallelization opportunity will enable these two developments to proceed hand in hand.

Lastly, outside of the field of radiation oncology, there are many other geometric design prob-

lems where “excessive edges” are views as undesirable, and there is a need to balance some type of edge metric penalty with other measures of quality; gerrymandering in political redistricting is a particularly prominent exaple. Optimization models similar to the one developed in this thesis, and appropriate solution methods, can be a promising direction to address such problems.

BIBLIOGRAPHY

- Agnew, C. E., D. M. Irvine, and C. K. McGarry (2014), Correlation of phantom-based and log file patient-specific QA with complexity scores for VMAT, *Journal of Applied Clinical Medical Physics*, 15(6), 204–216, doi:10.1120/jacmp.v15i6.4994.
- Bazaraa, M. S., H. D. Sherali, and C. M. Shetty (2006), *Nonlinear Programming: Theory and Algorithms*, 3rd ed., John Wiley & Sons, Inc.
- Biegala, M., and A. Hydzik (2016), Analysis of dose distribution in organs at risk in patients with prostate cancer treated with the intensity-modulated radiation therapy and arc technique, *Journal of Medical Physics*, 41(3), 198–204, doi:10.4103/0971-6203.189490.
- Broderick, M., M. Leech, and M. Coffey (2009), Direct aperture optimization as a means of reducing the complexity of intensity modulated radiation therapy plans, *Radiation Oncology*, 4(1), 8, doi:10.1186/1748-717X-4-8.
- Bush, K., S. Zavgorodni, I. Gagne, R. Townson, W. Ansbacher, and W. Beckham (2010), Monte Carlo evaluation of RapidArc oropharynx treatment planning strategies for sparing of midline structures, *Phys Med Biol*, 55(16), 4465–4479, doi:10.1088/0031-9155/55/16/s03.
- Cao, D., M. K. N. Afghan, J. Ye, F. Chen, and D. M. Shepard (2009), A generalized inverse planning tool for volumetric-modulated arc therapy, *Physics in Medicine and Biology*, 54(21), 6725, doi:10.1088/0031-9155/54/21/018.
- Carlsson, F. (2008), Combining segment generation with direct step-and-shoot optimization in intensity-modulated radiation therapy, *Medical physics*, 35(9), 3828–3838.
- Chen, M., Y. Chen, Q. Chen, and W. Lu (2011), Theoretical analysis of the thread effect in helical TomoTherapy, *Medical Physics*, 38(11), 5945–5960, doi:10.1118/1.3644842.
- Chiavassa, S., I. Bessieres, M. Edouard, M. Mathot, and A. Moignier (2019), Complexity metrics for IMRT and VMAT plans: a review of current literature and applications, *The British Journal of Radiology*, 92, 20190,270, doi:10.1259/bjr.20190270.
- Craft, D., D. McQuaid, J. Wala, W. Chen, E. Salari, and T. Bortfeld (2012), Multicriteria VMAT optimization, *Medical Physics*, 39(2), 686–696, doi:10.1118/1.3675601.

- Craft, D., M. Bangert, T. Long, D. Papp, and J. Unkelbach (2014), Shared data for intensity modulated radiation therapy (IMRT) optimization research: the CORT dataset, *GigaScience*, 3(1), doi:10.1186/2047-217X-3-37.
- Crowe, S. B., T. Kairn, N. Middlebrook, B. Sutherland, B. Hill, J. Kenny, C. M. Langton, and J. V. Trapp (2015), Examination of the properties of IMRT and VMAT beams and evaluation against pre-treatment quality assurance results, *Physics in Medicine and Biology*, 60(6), 2587–2601, doi:10.1088/0031-9155/60/6/2587.
- Das, I. J., G. X. Ding, and A. Ahnesjö (2008), Small fields: Nonequilibrium radiation dosimetry, *Medical Physics*, 35(1), 206–215, doi:10.1118/1.2815356.
- De Meerleer, G. O., G. M. Villeirs, L. Vakaet, L. J. Delrue, and W. J. De Neve (2004), The Incidence of Inclusion of the Sigmoid Colon and Small Bowel in the Planning Target Volume in Radiotherapy for Prostate Cancer, *Strahlentherapie und Onkologie*, 180(9), 573–581, doi:10.1007/s00066-004-1267-5.
- Du, W., S. H. Cho, X. Zhang, K. E. Hoffman, and R. J. Kudchadker (2014), Quantification of beam complexity in intensity-modulated radiation therapy treatment plans, *Medical physics*, 41(2), 21,716, doi:10.1118/1.4861821.
- Ezzell, G. A., et al. (2003), Guidance document on delivery, treatment planning, and clinical implementation of IMRT: Report of the IMRT subcommittee of the AAPM radiation therapy committee, *Medical Physics*, 30(8), 2089–2115, doi:10.1118/1.1591194.
- Fog, L. S., J. F. B. Rasmussen, M. Aznar, F. Kjær-Kristoffersen, I. R. Vogelius, S. A. Engelholm, and J. P. Bangsgaard (2011), A closer look at RapidArc® radiosurgery plans using very small fields, *Physics in medicine and biology*, 56(6), 1853–1863, doi:10.1088/0031-9155/56/6/020.
- Fraass, B. A., J. M. Steers, M. M. Matuszak, and D. L. McShan (2012), Inverse-optimized 3D conformal planning: Minimizing complexity while achieving equivalence with beamlet IMRT in multiple clinical sites, *Medical Physics*, 39, 3361–3374, doi:10.1118/1.4709604.
- Fredh, A., J. B. Scherman, L. S. Fog, and P. Munck (2013), Patient QA systems for rotational radiation therapy : A comparative experimental study with intentional errors, *Medical Physics*, 031716, doi:10.1118/1.4788645.
- Gibbons, J. P., K. Smith, D. Cheek, and I. Rosen (2009), Independent Calculation of Dose from a Helical TomoTherapy, *Journal of Applied Clinical Medical Physics*, 10(1), 103–119.
- Götstedt Julia Karlsson Hauer, B. A. (2015), Development and evaluation of aperture-based complexity metrics using film and EPID measurements of static MLC openings, *Medical physics*, 42(7), 3911–3921.
- Grigorov, G., T. Kron, E. Wong, J. Chen, J. Sollazzo, and G. Rodrigues (2003), Optimization of helical tomotherapy treatment plans for prostate cancer, *Physics in Medicine and Biology*, 48(13), 1933.

- Gurobi Optimization (2018), *Gurobi optimizer reference manual*, Gurobi Optimization, LLC.
- Gutiérrez, A. N., D. C. Westerly, W. A. Tomé, H. A. Jaradat, T. R. Mackie, S. M. Bentzen, D. Khuntia, and M. P. Mehta (2007), Whole Brain Radiotherapy With Hippocampal Avoidance and Simultaneously Integrated Brain Metastases Boost: A Planning Study, *International Journal of Radiation Oncology Biology Physics*, 69(2), 589–597, doi:10.1016/j.ijrobp.2007.05.038.
- Hardcastle, N., A. Bayliss, J. H. D. Wong, A. B. Rosenfeld, and W. A. Tom (2012), Improvements in dose calculation accuracy for small off-axis targets in high dose per fraction tomotherapy, *Medical Physics*, 39(8), 4788–4794, doi:10.1118/1.4736811.
- Hawkins, R. B. (1994), A Statistical Theory of Cell Killing by Radiation of Varying Linear Energy Transfer, *Radiation Research*, 140(3), 366, doi:10.2307/3579114.
- Heilemann, G., B. Poppe, and W. Laub (2013), On the sensitivity of common gamma-index evaluation methods to MLC misalignments in Rapidarc quality assurance, *Medical Physics*, 40(3), doi:10.1118/1.4789580.
- Holmes, T. W., T. R. Mackie, and P. Reckwerdt (1995), An iterative filtered backprojection inverse treatment planning algorithm for tomotherapy, *International Journal of Radiation Oncology*Biological*Physics*, 32(4), 1215–1225, doi:https://doi.org/10.1016/0360-3016(94)00465-W.
- Hwang, S., S. Ye, J. Park, and J. Kim (2014), Study on the Sensitivity of Gamma-Index Evaluation Methods to Positioning Errors of High-Definition MLC of True Beam STx in VMAT Quality Assurance for SBRT, *International Journal of Radiation Oncology*Biological*Physics*, 90(1), S858—S859, doi:10.1016/j.ijrobp.2014.05.2456.
- Jeraj, R., T. R. Mackie, J. Balog, G. Olivera, D. Pearson, J. Kapatoes, K. Ruchala, and P. Reckwerdt (2004), Radiation characteristics of helical tomotherapy, *Medical Physics Med. Phys*, 31(30), doi:10.1118/1.1639148.
- Kampfer, S., S. Schell, M. N. Duma, J. J. Wilkens, and P. Kneschaurek (2011), Measurements to predict the time of target replacement of a helical tomotherapy, *Journal of applied clinical medical physics / American College of Medical Physics*, 12(4), 3596.
- Kapatoes, J. M., G. H. Olivera, P. J. Reckwerdt, E. E. Fitchard, E. A. Schloesser, and T. R. Mackie (1999), Delivery verification in sequential and helical tomotherapy, *Physics in Medicine and Biology*, 44(7), 1815–1841, doi:10.1088/0031-9155/44/7/318.
- Kissick, M. W., J. Fenwick, J. A. James, R. Jeraj, J. M. Kapatoes, H. Keller, T. R. Mackie, G. Olivera, and E. T. Soisson (2005), The helical tomotherapy thread effect, *Medical physics*, 32(5), 1414–1423, doi:10.1118/1.1896453.
- Langen, K. M., et al. (2010), QA for helical tomotherapy: report of the AAPM Task Group 148, *Medical physics*, 37(9), 4817–4853, doi:10.1118/1.3462971.

- Long, T. C. (2015), Optimization Problems in Radiation Therapy Treatment Planning, Ph.D. thesis, The University of Michigan.
- Lu, W. (2010), A non-voxel-based broad-beam (NVBB) framework for IMRT treatment planning, *Physics in Medicine and Biology*, 55(23), 7175–7210, doi:10.1088/0031-9155/55/23/002.
- Mackie, T. R., J. W. Scrimger, and J. J. Battista (1985), A convolution method of calculating dose for 15-MV x rays, *Medical physics*, 12(2), 188–196.
- Mackie, T. R., T. Holmes, S. Swerdloff, P. Reckwerdt, J. O. Deasy, J. Yang, B. Paliwal, and T. Kinsella (1993), Tomotherapy: A new concept for the delivery of dynamic conformal radiotherapy, *Medical Physics*, 20(6), 1709–1719, doi:10.1118/1.596958.
- Mahnam, M., M. Gendreau, N. Lahrichi, and L.-M. Rousseau (2017), Simultaneous delivery time and aperture shape optimization for the volumetric-modulated arc therapy (VMAT) treatment planning problem, *Physics in Medicine & Biology*, 62(14), 5589–5611, doi:10.1088/1361-6560/aa7447.
- Marias, K., et al. (2011), Clinically driven design of multi-scale cancer models : the ContraCancrum project paradigm, *Interface Focus*.
- Masi, L., R. Doro, V. Favuzza, S. Cipressi, and L. Livi (2013), Impact of plan parameters on the dosimetric accuracy of volumetric modulated arc therapy, *Medical Physics*, 40(7), 71,718, doi:10.1118/1.4810969.
- McGarry, C. K., et al. (2011), Assessing software upgrades, plan properties and patient geometry using intensity modulated radiation therapy (IMRT) complexity metrics, *Medical physics*, 38(4), 2027–2034, doi:10.1118/1.3562897.
- McNiven, A. L., M. B. Sharpe, and T. G. Purdie (2010), A new metric for assessing IMRT modulation complexity and plan deliverability, *Medical physics*, 37(2), 505–515, doi:10.1118/1.3276775.
- Men, C., H. E. Romeijn, X. Jia, and S. B. Jiang (2010), Ultrafast treatment plan optimization for volumetric modulated arc therapy (VMAT), *Medical physics*, 37(11), 5787–5791, doi:10.1118/1.3491675.
- Oliver, M., I. Gagne, K. Bush, S. Zavgorodni, W. Ansbacher, and W. Beckham (2010), Clinical significance of multi-leaf collimator positional errors for volumetric modulated arc therapy, *Radiotherapy and Oncology*, 97(3), 554–560, doi:10.1016/j.radonc.2010.06.013.
- Olivera, G. H., D. M. Shepard, P. J. Reckwerdt, K. Ruchala, J. Zachman, E. E. Fitchard, and T. R. Mackie (1998), Maximum likelihood as a common computational framework in tomotherapy, *Physics in Medicine and Biology*, 43(11), 3277–3294, doi:10.1088/0031-9155/43/11/008.
- Otto, K. (2008), Volumetric modulated arc therapy: IMRT in a single gantry arc, *Medical Physics*, 35(1), 310–317, doi:10.1118/1.2818738.

- Papp, D., and J. Unkelbach (2014), Direct leaf trajectory optimization for volumetric modulated arc therapy planning with sliding window delivery, *Med Phys*, 41(1), 11,701, doi:10.1118/1.4835435.
- Park, J. M., S. Y. Park, and H. Kim (2015a), Modulation index for VMAT considering both mechanical and dose calculation uncertainties, *Physics in Medicine and Biology*, 60(18), 7101–7125, doi:10.1088/0031-9155/60/18/7101.
- Park, J. M., H. G. Wu, J. H. Kim, J. N. K. Carlson, and K. Kim (2015b), The effect of MLC speed and acceleration on the plan delivery accuracy of VMAT, *British Journal of Radiology*, 88(1049), 16–24, doi:10.1259/bjr.20140698.
- Peng, F. (2013), Optimization Methods for Volumetric Modulated Arc Therapy and Radiation Therapy Under Uncertainty, Ph.D. thesis, University of Michigan.
- Peng, F., X. Jia, X. Gu, M. A. Epelman, H. E. Romeijn, and S. B. Jiang (2012), A new column-generation-based algorithm for VMAT treatment plan optimization, *Physics in Medicine and Biology*, 57(14), 4569–4588.
- Peng, F., S. B. Jiang, H. E. Romeijn, and M. A. Epelman (2015), VMATc: VMAT with constant gantry speed and dose rate, *Physics in Medicine and Biology*, 60(7), 2955–2979.
- Romeijn, H. E., and J. F. Dempsey (2008), Intensity modulated radiation therapy treatment plan optimization, *TOP*, 16(2), 215–243, doi:10.1007/s11750-008-0064-1.
- Romeijn, H. E., R. K. Ahuja, J. F. Dempsey, and A. Kumar (2005), A column generation approach to radiation therapy treatment planning using aperture modulation, *SIAM Journal on Optimization*, 15(3), 838–862.
- Salari, E., J. Wala, and D. Craft (2012), Exploring trade-offs between VMAT dose quality and delivery efficiency using a network optimization approach, *Physics in medicine and biology*, 57(17), 5587–5600, doi:10.1088/0031-9155/57/17/5587.
- Seco, J., and F. Verhaegen (2013), *Monte Carlo techniques in radiation therapy*, CRC press.
- Sheng, K. (2017), TomoTherapy, in *Principles and Practice of Image-Guided Radiation Therapy of Lung Cancer*, edited by J. Cai, J. Y. Chang, and F.-F. Yin, chap. 8, pp. 141–161, CRC Press.
- Shepard, D. M., G. H. Olivera, P. J. Reckwerdt, and T. R. Mackie (2000), Iterative approaches to dose optimization in tomotherapy, *Physics in medicine and biology*, 45(1), 69–90, doi:10.1088/0031-9155/45/1/306.
- Stambaugh, C., B. Nelms, T. Wolf, R. Mueller, M. Geurts, D. Opp, G. Zhang, E. Moros, and V. Feygelman (2015), Measurement-guided volumetric dose reconstruction for helical tomotherapy, *Journal of Applied Clinical Medical Physics*, 16(2), 302–321, doi:10.1120/jacmp.v16i2.5298.

- Sterpin, E., F. Salvat, R. Cravens, K. Ruchala, G. H. Olivera, and S. Vynckier (2008), Monte Carlo simulation of helical tomotherapy with PENELOPE, *Physics in medicine and biology*, 53(8), 2161–2180, doi:10.1088/0031-9155/53/8/011.
- Tudor, G. S. J., and S. J. Thomas (2013), Impact of the fixed gantry angle approximation on dosimetric accuracy for helical tomotherapy plans, *Medical Physics*, 40(1), 11,711, doi:10.1118/1.4769120.
- Ulmer, W., J. Pyry, and W. Kaissl (2005), A 3D photon superposition/convolution algorithm and its foundation on results of Monte Carlo calculations, *Physics in Medicine and Biology*, 50(8), 1767–1790, doi:10.1088/0031-9155/50/8/010.
- Van Dyk, J. (1999), *The modern technology of radiation oncology: a compendium for medical physicists and radiation oncologists*, 521–587 pp., Medical Physics Pub., Madison, Wis.
- Varian Medical Systems (2011), *TrueBeam STx System Specifications*, Varian Medical Systems, Inc., Palo Alto, CA, USA.
- Wala, J., E. Salari, W. Chen, and D. Craft (2012), Optimal partial-arcs in VMAT treatment planning, *Physics in medicine and biology*, 57(18), 5861–5874, doi:10.1088/0031-9155/57/18/5861.
- Webb, S. (2001), *Intensity-Modulated Radiation Therapy*, 64–74 pp., Taylor and Francis Group.
- Westerly, D. C., E. Soisson, Q. Chen, K. Woch, L. Schubert, G. Olivera, and T. R. Mackie (2009), Treatment Planning to Improve Delivery Accuracy and Patient Throughput in Helical Tomotherapy, *International Journal of Radiation Oncology Biology Physics*, 74(4), 1290–1297, doi:10.1016/j.ijrobp.2009.02.004.
- Yang, K., D. Yan, and N. Tyagi (2012), Sensitivity analysis of physics and planning SmartArc parameters for single and partial arc VMAT planning, *Journal of Applied Clinical Medical Physics*, 13(6), 34–45, doi:10.1120/jacmp.v13i6.3760.
- Younge, K. C., M. M. Matuszak, J. M. Moran, D. L. McShan, B. A. Fraass, and D. A. Roberts (2012), Penalization of aperture complexity in inversely planned volumetric modulated arc therapy, *Medical Physics*, 39(11), 7160–7170.
- Yu, C. X. (1995), Intensity-modulated arc therapy with dynamic multileaf collimation: an alternative to tomotherapy, *Physics in medicine and biology*, 40(9), 1435–1449, doi:10.1088/0031-9155/40/9/004.
- Zhao, Y.-L., M. Mackenzie, C. Kirkby, and B. G. Fallone (2008), Monte Carlo calculation of helical tomotherapy dose delivery, *Medical Physics*, 35(8), 3491–3500, doi:10.1118/1.2948409.
- Zwan, B. J., J. Hindmarsh, E. Seymour, K. Kandasamy, K. Sloan, R. David, and C. Lee (2016), The dosimetric impact of control point spacing for sliding gap MLC fields, *Journal of Applied Clinical Medical Physics*, 17(6), 204–216, doi:10.1120/jacmp.v17i6.6345.