

UCSF

UC San Francisco Previously Published Works

Title

A GPU-based multi-criteria optimization algorithm for HDR brachytherapy

Permalink

<https://escholarship.org/uc/item/8q60216z>

Journal

Physics in Medicine and Biology, 64(10)

ISSN

0031-9155

Authors

Bélanger, Cédric
Cui, Songye
Ma, Yunzhi
[et al.](#)

Publication Date

2019-05-01

DOI

10.1088/1361-6560/ab1817

Peer reviewed

A GPU-based multi-criteria optimization algorithm for HDR brachytherapy

Cédric Bélanger^{1,2‡}, Songye Cui^{1,2‡}, Yunzhi Ma², Philippe Després^{1,2}, J. Adam M. Cunha³, Luc Beaulieu^{1,2}

¹Department of Physics, Engineering Physics and Optics and Cancer Research Center, Université Laval, Quebec City, QC, G1V 0A6, Canada

²Department of Radiation Oncology and Research Center of CHU de Québec - Université Laval, Quebec City, QC, G1R 2J6, Canada

³Radiation Oncology, University of California, San Francisco, CA 94115, USA

[‡]Co-first authorship

E-mail: Luc.Beaulieu@phy.ulaval.ca

Abstract. Currently in HDR brachytherapy planning, a manual fine-tuning of an objective function is necessary to obtain case-specific valid plans. This study intends to facilitate this process by proposing a patient-specific inverse planning algorithm for HDR prostate brachytherapy: GPU-based multi-criteria optimization (gMCO).

Two GPU-based optimization engines including simulated annealing (gSA) and a quasi-Newton optimizer (gL-BFGS) were implemented to compute multiple plans in parallel. After evaluating the equivalence and the computation performance of these two optimization engines, one preferred optimization engine was selected for the gMCO algorithm. Five hundred sixty-two previously treated prostate HDR cases were divided into validation set (100) and test set (462). In the validation set, the number of Pareto optimal plans to achieve the best plan quality was determined for the gMCO algorithm. In the test set, gMCO plans were compared with the physician-approved clinical plans.

Our results indicated that the optimization process is equivalent between gL-BFGS and gSA, and that the computational performance of gL-BFGS is up to 67 times faster than gSA. Over 462 cases, the number of clinically valid plans was 428 (92.6%) for clinical plans and 461 (99.8%) for gMCO plans. The number of valid plans with target V_{100} coverage greater than 95% was 288 (62.3%) for clinical plans and 414 (89.6%) for gMCO plans. The mean planning time was 9.4 s for the gMCO algorithm to generate 1000 Pareto optimal plans.

In conclusion, gL-BFGS is able to compute thousands of SA equivalent treatment plans within a short time frame. Powered by gL-BFGS, an ultra-fast and robust multi-criteria optimization algorithm was implemented for HDR prostate brachytherapy. Plan pools with various trade-offs can be created with this algorithm. A large-scale comparison against physician approved clinical plans showed that treatment plan quality could be improved and planning time could be significantly reduced with the proposed gMCO algorithm.

Keywords: brachytherapy, prostate cancer, patient-specific, treatment planning, optimization, GPU

1 Introduction

About 52.3% of non-skin cancer patients receive radiation therapy during the course of their illness (Citrin 2017, DeVita *et al* 2015, Delaney *et al* 2005). The most common radiation therapy treatment particle type used is the photon, which can be delivered either externally from a medical linear accelerator (External Beam Radiation Therapy - EBRT) or internally from an inserted small radioactive source (brachytherapy, high dose rate (HDR) or low dose rate (LDR)).

Dose prescriptions in modern radiation treatment planning contain both tumor and healthy organ objectives. These objectives are often conflicting and can be generalized as: treating the tumor with high radiation dose and sparing the healthy organs with low radiation dose. Computerized treatment planning systems were used to formulate clinical prescriptions into a mathematical optimization problem, and to find treatment plans that well presented these prescriptions with treatment facilities.

However, most available algorithms are not inherently patient-specific in a sense that manual re-plannings are usually inevitable to find a clinically acceptable plan for each patient. As a result, the planning procedure can be time consuming and the planning output is planner dependant (Moore *et al* 2011, Nelms *et al* 2012, Wu *et al* 2009).

Several patient-specific inverse planning algorithms such as knowledge-based planning (KBP), auto-planning (AP) and multi-criteria optimization (MCO) have been proposed in EBRT. In KBP, one plan is created for a new case by searching in a prior physician-approved plan dataset based on the geometric features (Moore *et al* 2011, Wu *et al* 2011, Wu *et al* 2009, Petit *et al* 2012). In AP, a clinical plan can be obtained by interactively and automatically adapting objectives, constraints and dose shaping contours (Hazell *et al* 2016). In MCO, a plan pool is constructed by generating plans with various trade-offs on Pareto surfaces (Craft *et al* 2006, Teichert *et al* 2011). Similar studies can also be found in brachytherapy (van der Meer *et al* 2018, Shen *et al* 2018, Zhou *et al* 2017, Cui *et al* 2018a, Cui *et al* 2018b).

Our prior studies (Cui *et al* 2018a, Cui *et al* 2018b) showed that a patient-specific treatment plan can be created without any user interventions in HDR prostate brachytherapy. However, the optimization engine of these studies was stochastic, and was implemented on CPU hardware (Cui *et al* 2018a, Cui *et al* 2018b). As a result, the algorithm inevitably involved an intensive computation (41 s), which may restrain its application in clinical practice, because the patient is under general anesthesia in the operating room waiting for the treatment to be delivered.

The capability of graphics processing unit (GPU) architecture in reducing calculation time in medical physics were reviewed in (Pratx and Xing 2011, Jia *et al* 2014, Després and Jia 2017). The purpose of this study is to propose an ultra-fast patient-specific inverse planning algorithm on GPU for HDR brachytherapy.

2 Methods and Materials

This section begins with a detailed description of experimental setups including patient selection, mathematical formulations and computational specifications. Next, two inverse planning optimization engines were implemented on GPU architecture to calculate multiple plans in parallel and to populate the Pareto surfaces. Powered by the preferred optimization engine, a GPU-based multi-criteria optimization algorithm (gMCO) which is able to automatically generate clinical plans was proposed to eliminate the re-planning problem in HDR brachytherapy. In the end, a comprehensive comparison, including dosimetric performance as well as planning time, between clinical plans and gMCO plans was made.

2.1 Experimental setup

2.1.1 Patient selection An anonymous dataset that contains 562 prostate cancer patients who received an HDR brachytherapy treatment as a boost to EBRT from April 2011 to July 2016 at our institution was studied. This dataset incorporates the cases studied in prior works (Edimo *et al* 2019, Cui *et al* 2018a, Cui *et al* 2018b). Among the dataset, 100 random cases (validation set) were used to determine the number of Pareto optimal plans with the gMCO algorithm, and 462 random cases (test set) were used in the performance evaluation of the gMCO generated plans.

After inserting 16-18 plastic catheters into the prostate under a transrectal ultrasound guidance, the anatomy of these patients was obtained from CT scans. Organ structures (prostate, urethra, bladder and rectum) were delineated and were imported into a commercial treatment planning system (Elekta Oncentra Brachy IPSA, Veenendaal, The Netherlands). The prescription was to deliver 15 Gy in a single fraction to the prostate. Plans were delivered using a Flexitron afterloader (Elekta Brachy, Veenendaal, The Netherlands) with an Ir-192 radioactive source.

The dwell positions were extracted from the DICOM-RT files of clinical plans, and the mean number of active dwell positions (N_{act}) used for the optimization was 171 (range:102-385). The mean number of dose calculation points (N_{pnt}) used for the optimization was 5913 (range:2753-15 998), and the mean number of dose calculation points used for the dose-volume histogram (DVH) computations was 31 039 (range:11 451-66 089).

2.1.2 Quadratic objective function formulation Inverse Planning Simulated Annealing (IPSA) (Lessard and Pouliot 2001) was used as a dose optimization engine in our prior studies (Cui *et al* 2018a, Cui *et al* 2018b). In IPSA, piecewise linear objective functions were solved with simulated annealing (Lessard and Pouliot 2001), a stochastic optimizer. These objective functions were constructed with a population based planning template called a class solution (Cui *et al* 2018a, Cui *et al* 2018b).

In order to implement an efficient optimizer, one option is to replace the stochastic optimizer with a gradient-based optimizer. Therefore, it may be necessary to replace

the IPSA linear piecewise objective functions with piecewise quadratic objective functions, so that the first derivative (gradient) of the objective function is continuous. Quadratic objective functions are usually solved with gradient-based optimizers in radiation therapy (Milickovic *et al* 2002, Lahanas, Baltas and Giannouli 2003, Lahanas, Schreibmann and Baltas 2003, Men *et al* 2009).

The dose at the i^{th} dose calculation point in the j^{th} organ, denoted by d_{ij} , is described in equation (1)

$$d_{ij} = \sum_{l=1}^{N_{act}} \dot{d}_{ijl} t_l \quad (1)$$

where \dot{d}_{ijl} is the dose rate contribution of the l^{th} dwell position to the i^{th} dose calculation point in the j^{th} organ, and t_l is the dwell time of the l^{th} dwell position. In order to avoid negative dwell times, new decision variables called dwell weight ($x_l = t_l^{1/2}$) were introduced as in (Milickovic *et al* 2002, Lahanas, Baltas and Giannouli 2003). With this substitution, the dwell times are always non-negative ($t_l = x_l^2$).

The piecewise quadratic objective function f_{ij} at the i^{th} dose calculation point of the j^{th} organ is given in equation (2)

$$f_{ij}(d_{ij}) = \begin{cases} w_{\min} \cdot (D_{\min} - d_{ij})^2 & d_{ij} < D_{\min} \\ 0 & D_{\min} \leq d_{ij} \leq D_{\max} \\ w_{\max} \cdot (d_{ij} - D_{\max})^2 & d_{ij} > D_{\max} \end{cases} \quad (2)$$

Variables D_{\min} and D_{\max} are the underdose limit and the overdose limit respectively, and variables w_{\min} and w_{\max} are the corresponding weights. The corresponding gradient function g_{ij} of equation (2) is described in equation (3)

$$g_{ij}(x_l) = \frac{\partial f_{ij}}{\partial x_l} = \begin{cases} 4 \cdot \dot{d}_{ijl} \cdot x_l \cdot w_{\min} \cdot (d_{ij} - D_{\min}) & d_{ij} < D_{\min} \\ 0 & D_{\min} \leq d_{ij} \leq D_{\max} \\ 4 \cdot \dot{d}_{ijl} \cdot x_l \cdot w_{\max} \cdot (d_{ij} - D_{\max}) & d_{ij} > D_{\max} \end{cases} \quad (3)$$

The single joint MCO objective function to be minimized is defined as a weighted sum in equation (4)

$$F = \sum_{j=1}^{N_O} w_j \cdot \frac{1}{N_{\text{pnt},j}} \sum_{i=1}^{N_{\text{pnt},j}} f_{ij}(d_{ij}) \quad (4)$$

where N_O is the number of organs, $N_{\text{pnt},j}$ is the number of dose calculation points in the j^{th} organ. w_j is a hidden weight applied to the objectives (surface and volume) of the j^{th} organ to introduce trade-off in the solution space around the population-based starting point as in (Cui *et al* 2018a, Cui *et al* 2018b). The hidden weights are always non negative and their sum is one (because of the weighted sum method).

The original class solution designed for the piecewise linear objective functions (Cui *et al* 2018a, Cui *et al* 2018b) will no longer be appropriate to construct the new quadratic objective functions, and so a new one must be designed (table 1).

Table 1: The class solution to formulate quadratic objective functions (equation (2)) for 15 Gy prostate boost HDR treatment (Surface: surface dose calculation points, Volume: volume dose calculation points).

Organ	Surface				Volume			
	w_{\min}	$D_{\min}(\text{Gy})$	$D_{\max}(\text{Gy})$	w_{\max}	w_{\min}	$D_{\min}(\text{Gy})$	$D_{\max}(\text{Gy})$	w_{\max}
Target	200	15	22.5	80	200	15	22.5	1
Urethra	30	14	16	160	30	14	16	160
Bladder	0	0	7.5	60	0	0	7.5	60
Rectum	0	0	7.5	15	0	0	7.5	15

2.1.3 Computational specifications The CPU algorithm was written in *C++*, compiled with *g++* (7.3.0) and executed on a six-core Intel Xeon CPU (E5-2620 v3 @ 2.40 GHz). The GPU algorithms were written in CUDA C, compiled with *nvcc* (CUDA toolkit 10.0.130) and executed on an NVIDIA Titan X (Pascal) GPU.

2.2 GPU-based efficient optimization engines

Previous studies showed that it is feasible to find clinically acceptable treatment plans after exploring Pareto surfaces with MCO approaches (Craft *et al* 2006, Cui *et al* 2018a). However, constructing Pareto surfaces could be inefficient, if performed sequentially.

2.2.1 IPSA on GPU A traditional CPU-based inverse planning algorithm such as IPSA (or cSA) (Lessard and Pouliot 2001) can be divided into several serial computing steps (figure 1). In each step, the same operation is repeated over a large dataset. For example, the following five steps are essential in cSA:

- (i) Initialization and dose rate matrix calculation (repeated for: N_{pnt} dose calculation points \times N_{act} dwell positions),
- (ii) dwell time updates (repeated for: N_{act} dwell positions),
- (iii) dose calculations based on equation (1) (repeated for: N_{pnt} dose calculation points),
- (iv) objective function values calculation based on equation (2) (repeated for: N_{pnt} dose calculation points),
- (v) mean objective function evaluation based on equation (4) (repeated for: one accumulation over N_{pnt} dose calculation points).

To obtain an optimal solution or a treatment plan, steps (ii)-(v) are iteratively repeated in cSA. Furthermore, in order to explore Pareto surfaces by computing N_{plan} treatment plans, it is usually necessary to repeat the aforementioned steps N_{plan} times.

To increase the efficiency of MCO approaches, GPU-based IPSA (or gSA) was implemented on GPU architecture to compute treatment plans with various trade-offs in parallel. Two strategies were applied to achieve this purpose.

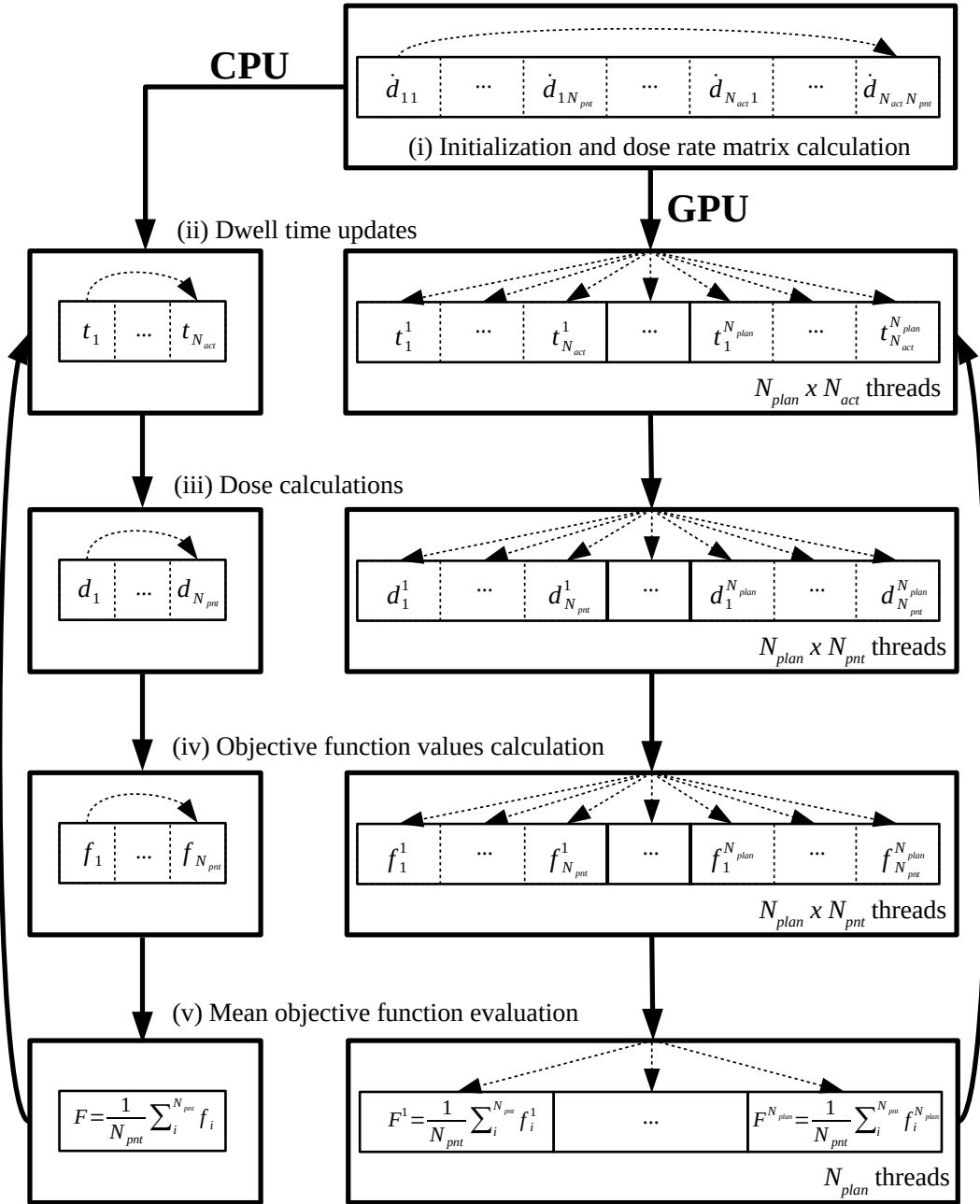


Figure 1: Illustration of the iterative procedure to optimize one treatment plan on CPU and N_{plan} plans on GPU. In each CPU or GPU iteration, the steps (ii)-(v) are executed sequentially. In each step on the CPU, the operations are executed sequentially in a loop. In each step on the GPU, the operations are executed in parallel on different threads for N_{plan} plans. (The superscript indicates the plan number on GPU).

First, the serial operations computed in each step in cSA were adapted to run in parallel on GPU, so the operations within each step can be executed simultaneously on different threads (figure 1). Note that in each step on GPU, the computational burden is N_{plan} times larger than in the CPU implementation (N_{plan} plans on GPU vs. one plan on CPU in figure 1). However, a performance gain can be achieved with the GPU implementation, as the huge burden of updating the values for all plans in each step is processed in parallel on different threads. To obtain N_{plan} optimal solutions or N_{plan} treatment plan with the proposed implementation, it is necessary to iteratively repeat steps (ii)-(v) in gSA.

Second, as frequent data transfers between CPU and GPU will slow down the computation, data transfer only occurs twice in gSA: once when preparing the data used for the optimization (CPU to GPU), once more when saving the dosimetric results onto the disk after the optimization (GPU to CPU).

2.2.2 Deterministic optimizer In section 2.2.1, a stochastic optimizer was implemented on CPU and on GPU. To further improve the computational performance, a deterministic optimizer (Limited-memory Broyden-Fletcher-Goldfarb-Shanno, L-BFGS) (Liu and Nocedal 1989, Wetzl and Taubmann 2013, Wetzl *et al* 2013) was introduced to replace the stochastic optimizer. There are two reasons to choose this quasi-Newton optimizer, (1) BFGS and its variants are widely studied in brachytherapy (Milickovic *et al* 2002, Lahanas, Baltas and Giannouli 2003), and (2) L-BFGS is widely used in clinic after being integrated in Hybrid Inverse Planning Optimization (HIPO) (Elekta Brachy, Veenendaal, The Netherlands) (Karabis *et al* 2005).

So far, four optimization engines were implemented: cSA, gSA (simulated annealing on CPU and on GPU), cL-BFGS and gL-BFGS (L-BFGS on CPU and on GPU). The description of L-BFGS implementation on CPU and on GPU is omitted in this study, due to the similarity with the context and figure 1 in section 2.2.1.

2.2.3 Equivalence between the four optimizers The equivalence between the four optimization engines was evaluated based on the same objective function (class solution in table 1) as tested over the validation set. For cSA, gSA, and cL-BFGS, one plan using uniform 5 s initial dwell times as a starting point was generated. For gL-BFGS, 1000 degenerated plans were calculated to evaluate the convergence of different starting points (randomly distributed between 0 and 10 s). The stopping criteria for cSA and gSA was specified by the number of iterations. The stopping criteria for cL-BFGS and gL-BFGS was specified by the parameter ϵ (based on the relative variation of the objective function (Men *et al* 2009)). To measure the equivalence between the four optimizers, 1 000 000 iterations and $\epsilon = 10^{-7}$ were used as the stopping criteria, because no significant improvements in the objective function were observed.

2.2.4 Pareto surfaces characterization with gSA and gL-BFGS Planning efficiency is a key factor when designing an inverse planning algorithm. For SA, a clinically useful

stopping criteria (50 000 iterations) can be used to reach Pareto surfaces (Cui *et al* 2018a). For gradient-based method, it is also desirable to find a stopping criteria that can well approximate the Pareto surfaces.

By computing solutions in parallel with various combinations of hidden weights, Pareto surfaces can be populated either with gSA and with gL-BFGS. Such solutions were Pareto optimal, or non-dominated, if no solution that improves any individual objective value without worsening at least one of the other individual objective values exists. A clinically useful stopping criteria was determined for gL-BFGS to approximate the Pareto surfaces, after examining the effect of different stopping criteria (ranging from $\epsilon = 10^{-7}$ to $\epsilon = 10^{-2}$) based on the fraction of non-dominated solutions and the speedup factor of the optimization time for all 100 validation cases.

2.2.5 Computational performance under clinically useful scenarios The benefits of the proposed GPU implementation over a traditional CPU implementation of inverse planning algorithms were explored. Based on the clinically useful stopping criteria, the computational performance of cSA, gSA, cL-BFGS and gL-BFGS were measured against the number of generated plans.

2.3 Patient-specific multi-criteria optimization algorithm

Usually, plans obtained with a population-based planning template are not always directly acceptable, and manual weights adjustments are required to obtain a patient-specific deliverable plan. After reviewing the definition of acceptable plans, a GPU-based multi-criteria optimization algorithm (gMCO) powered with gL-BFGS was proposed to eliminate the procedure of manual weights adjustments.

2.3.1 Plan evaluation The schedules of dose fractionation and the evaluation criteria of HDR prostate brachytherapy plans may vary between centers (Yamada *et al* 2012). According to the Radiation Therapy Oncology Group (RTOG) 0924 protocol (Radiation Therapy Oncology Group 2016), RTOG acceptable plans (or valid solutions) can be summarized as follows:

- Prostate/Target coverage constraint: $V_{100} \geq 90\%$ of the volume.
- Urethra constraint: $D_{10} < 118\%$ of the prescription dose.
- Bladder constraint: $V_{75} < 1$ cc.
- Rectum constraint: $V_{75} < 1$ cc.

Note:

(1) V_x refers to the absolute volume that receives $x\%$ of the prescription dose, and D_x refers to the percent of the prescription dose that covers $x\%$ of the volume.

(2) In this study, a more stringent set of criteria was introduced. It is designated by the RTOG+ symbol and is the same as the RTOG criteria set except that it specifies

a higher target coverage requirement of 95% for the V_{100} . This is usually attainable in the clinic without sacrificing the OAR protection.

2.3.2 gMCO algorithm Compared with our previous studies (Cui *et al* 2018a, Cui *et al* 2018b), there are three main differences in gMCO: (1) the trade-off between target and urethra is now explored, (2) the Pareto surfaces are widely explored with a large number of plans, as no prior knowledge of the RTOG+ valid solution space is involved, and (3) the validation cases were used to determine the number of parallel plans (from 1 to 10 000) needed to achieve high RTOG and RTOG+ acceptance rates with random hidden weights. In gMCO, the parallel plan computations were executed with gL-BFGS.

2.4 Comparison between clinical plans and gMCO plans

A plan pool was created with the gMCO algorithm. One plan was selected from the plan pool and was referred to as the gMCO plan.

The criteria used for plan selection are, in descending order of priority: RTOG+ valid plan, RTOG valid plan, RTOG invalid plan (violates at least one criteria). If multiple RTOG or RTOG+ valid plans existed, the one with a highest target V_{100} was selected. If multiple RTOG invalid plans existed, the one with the lowest bladder and rectum V_{75} (while not violating the criteria for target and urethra) was selected.

2.4.1 Dosimetric performance The dosimetric results of clinical plans were retrieved from Oncentra Brachy (Elekta Brachy, Veenendaal, The Netherlands). Dosimetric comparisons between clinical plans and gMCO plans were analyzed for 462 test cases. The overall result was examined based on RTOG and RTOG+ acceptance rates (the criteria of all organs were met). The acceptance rate (i.e. target V_{100} , urethra D_{10} , bladder V_{75} , and rectum V_{75}) for each organ was also reported.

2.4.2 Planning time The planning time consists of the time taken for dose calculation points creation, dose rate matrix calculation, optimization, and DVH calculation on GPU. The calculation time of each portion was recorded for gMCO plans. The total planning time was compared between clinical plans and gMCO plans.

3 Results

3.1 GPU-based optimization engines

3.1.1 Equivalence between the four optimizers The optimization processes of the four optimizers for one random validation case are illustrated in figure 2. From this figure, (1) gL-BFGS plans obtained with different initial dwell times converge to the SA objective function value, (2) no significant differences (within 0.02%) in objective function values resulted from the four optimizers were observed. Over all 100 validation cases, similar

results were observed, because the final objective function values of the four algorithms were in agreement within 0.2%.

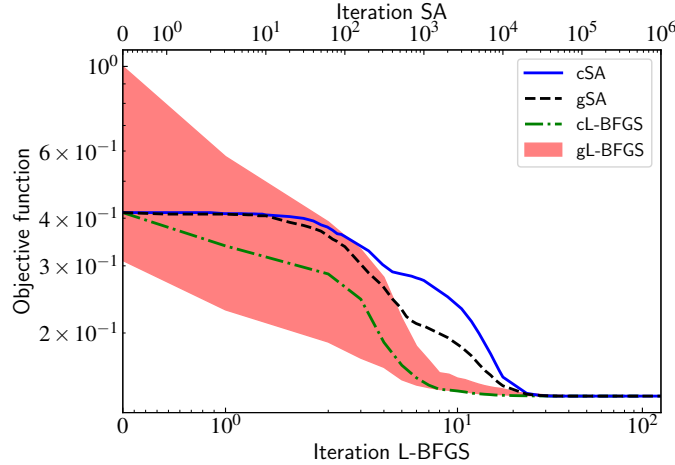


Figure 2: Illustration of cSA, gSA, cL-BFGS and gL-BFGS objective function values against the number of iterations for one random validation case. (The difference between CPU and GPU random number generators accounts for the different trajectories for cSA and gSA).

3.1.2 Pareto surfaces characterization with gSA and gL-BFGS To characterize Pareto surfaces, 100 000 different solutions were generated with gSA and gL-BFGS (1000 solutions/case for all 100 validation cases). For gSA, the mean fraction of non-dominated solutions was 99.6% under 50 000 iterations.

For gL-BFGS, the results in figure 3a indicate that the fraction of non-dominated solutions decreased (from 100% to 89.3%) as the stopping criteria increased (from $\epsilon = 10^{-7}$ to $\epsilon = 10^{-2}$). On the other hand, the speedup factor in the optimization time increased (from 1 to 10) as the stopping criteria increased (from $\epsilon = 10^{-7}$ to $\epsilon = 10^{-2}$). It should be noted that over 99.3% of the solutions obtained with a larger stopping criteria ($\epsilon = 10^{-3}$) are Pareto optimal solutions. Given that reaching optimality and a reasonable calculation time are important criteria for clinical applicability, the results in figure 3a suggest that there could be a time advantage in using a larger stopping criteria ($\epsilon = 10^{-3}$).

Furthermore, a single 2D Pareto surface characterization with gSA and gL-BFGS is shown in figure 3b. The results suggest that no significant difference in Pareto surfaces approximations is observed with GPU-based optimization engines under clinically useful stopping criteria and under more strict stopping criteria as specified in section 2.2.3. From these results, $\epsilon = 10^{-3}$ is used as the stopping criteria in gL-BFGS afterwards.

3.1.3 Computational performance under clinically useful scenarios Under clinically useful scenarios, the optimization time of cSA, gSA, cL-BFGS and gL-BFGS are shown

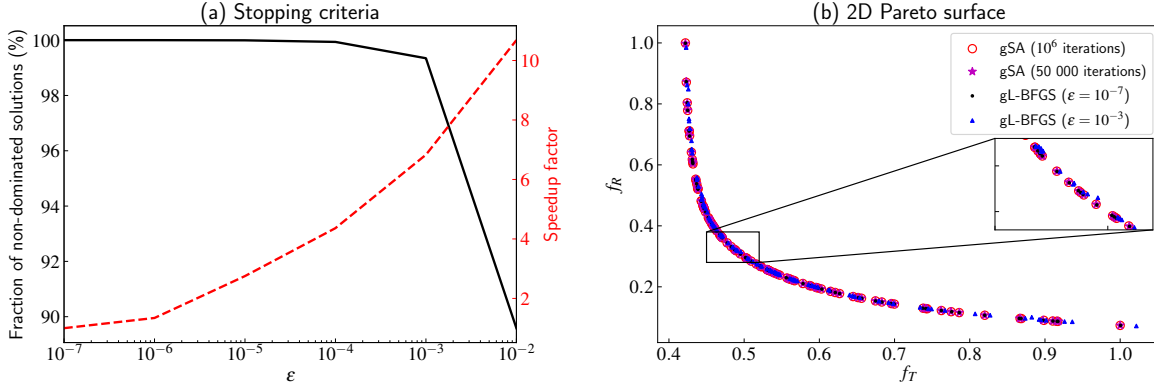


Figure 3: (a) Effect of the stopping criteria on the fraction of non-dominated solutions in the Pareto front characterized with gL-BFGS (black solid line) and the speedup factor of the optimization time (red dashed line). The speedup factor are normalized to the values obtained with a stopping criteria of $\epsilon = 10^{-7}$. (b) A comparison of 2D Pareto surface approximations with gSA and gL-BFGS optimization engines for a random case. (f_T is denoted for target individual objective function and f_R is denoted for rectum individual objective function).

in figure 4a. From the results, the time of all four engines increased as the number of plans increased. For 1000 plans, the mean optimization time was 9.2s/plan (cSA), 60ms/plan (gSA), 1s/plan (cL-BFGS), and 0.9ms/plan (gL-BFGS). In other words, compared with the cSA result, cL-BFGS can achieve a speedup factor up to 9, gSA can achieve a speedup factor of up to 176, and gL-BFGS can achieve a speedup factor of up to 10 990.

Figure 4b shows that the mean GPU memory usage increased with the number of plans for the GPU algorithms, and that the increase rate becomes significantly large when the number of plans reaches approximately 1000.

3.2 Patient-specific multi-criteria optimization algorithm

As the hidden weights were randomly generated in gMCO algorithm, the RTOG and RTOG+ acceptance rates were measured multiple times with different random hidden weight vectors in equation (4). In figure 5, the RTOG+ acceptance rate increases (from 17% to 85%) and the spread of the acceptance rate distributions decreases with the number of plans. However, a number of 1000 plans was selected as the best compromised between optimization time (which increases after 1000 plans, see figure 4a) and the RTOG+ acceptance rate (which does not increase significantly after 1000 plans) for gMCO algorithm.

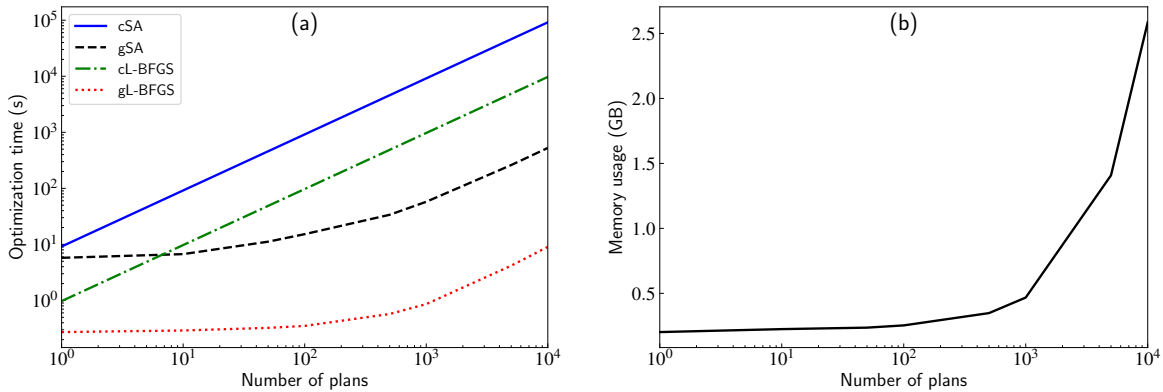


Figure 4: Computational performance against the number of plans for cSA, gSA, cL-BFGS and gL-BFGS under clinically useful scenarios (cSA and gSA: 1000 iterations, cL-BFGS and gL-BFGS: $\epsilon = 10^{-3}$): (a) the mean optimization time, (b) the mean GPU memory usage of gL-BFGS (the result of gSA was ignored, for its similarity to the gL-BFGS one).

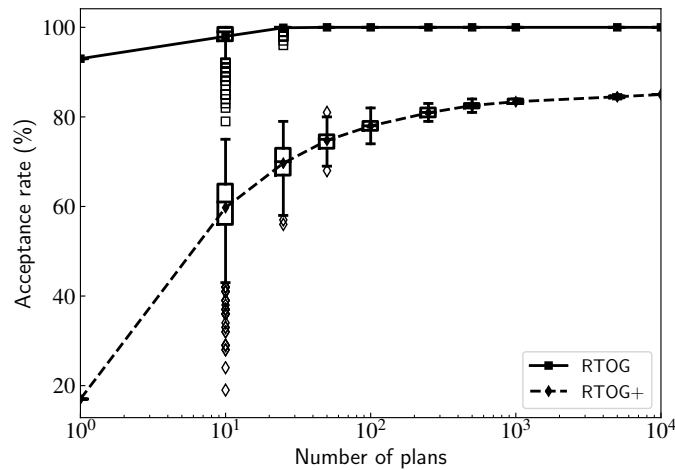


Figure 5: The effect of the number of plans on RTOG and RTOG+ acceptance rates for gMCO (including the spread of the distributions in the boxes).

3.3 Comparison between clinical plans and gMCO plans

3.3.1 Dosimetric performance The dosimetric comparison between clinical plans and gMCO plans is illustrated in figure 6. These results suggest that the mean target coverage was higher for gMCO plans (97.2%) than for clinical plans (95.3%). The mean urethra D_{10} was significantly higher for gMCO plans (115.7%) than for clinical plans (109.1%). The mean bladder V_{75} was 0.53 cc for clinical plans, and 0.78 cc for gMCO plans. For rectum sparing, the mean rectum V_{75} was 0.56 cc for clinical plans, and 0.52 cc for gMCO plans.

The acceptance rate results are summarized in table 2. For overall dosimetric

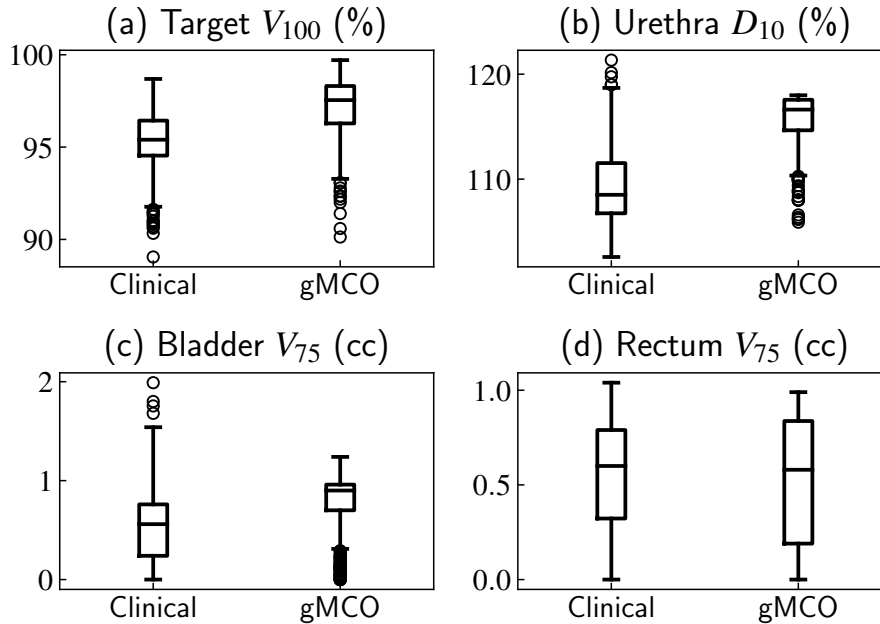


Figure 6: Dosimetric comparison between IPSA physician-approved plans and gMCO plans over the test cohort: (a) target V_{100} , (b) urethra D_{10} , (c) bladder V_{75} , and (d) rectum V_{75} .

performances, the number of RTOG valid plans was 428 (92.6%) for clinical plans, and 461 (99.8%) for gMCO plans. The number of RTOG+ valid plans was 288 (62.3%) for clinical plans, and 414 (89.6%) for gMCO plans.

Table 2: RTOG and RTOG+ acceptance rates (%) for clinically approved plans and gMCO plans over 462 test cases.

	RTOG					RTOG+		Time
	Target	Bladder	Rectum	Urethra	All	Target	All	
Clinical	99.8	95.2	98.7	98.5	92.6	64.1	62.3	mins
gMCO	100.0	99.8	100.0	100.0	99.8	89.6	89.6	9.4 s

The number of plans with a target coverage greater than 95% was 296 (64.1%) for clinical plans, and 414 (89.6%) for gMCO plans. The number of plans that exceeded the urethra sparing constraint was 7 for clinical plans, and 0 for gMCO plans. The number of plans that exceeded the bladder sparing constraint was 22 for clinical plans, and 1 for gMCO plans. The number of plans that exceeded the rectum sparing constraint was 6 for clinical plans, and 0 for gMCO plans. In addition, the mean number of RTOG valid plans was 617/1000 (61.7%), and the mean number of RTOG+ valid plans was 268/1000 (26.8%) for the gMCO plan pool.

As a supplement to the general comparisons described above, one example case was chosen to illustrate the advantage of gMCO in terms of the results of DVHs and isodose curves in figure 7.

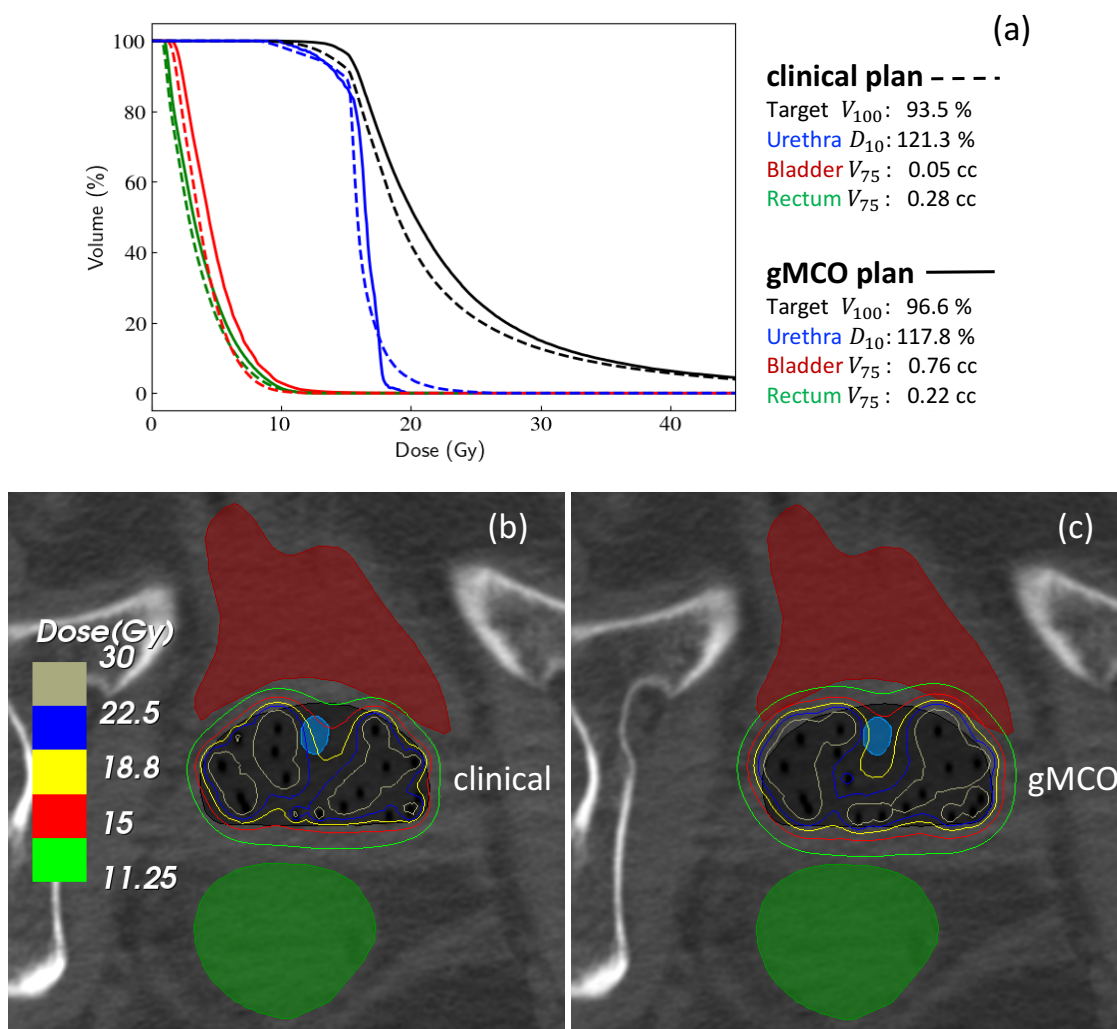


Figure 7: A comparison between the clinical plan and the gMCO plan for one example case: (a) DVHs and dosimetric parameters, (b) and (c) isodose curves.

3.3.2 Planning time The time to create a plan is of the order of a few minutes in clinic, including manual tweaking of the objective function and/or dwell times. On the other hand, the mean planning time was 9.4s for gMCO to generate 1000 optimal plans. Among these numbers, the mean dose calculation points creation time was 7.4s, which represents 79% of the mean planning time. The mean optimization time was 0.8s (8.5% of the mean total planning time). Dose rate matrix calculation and DVH calculation on GPU contribute to the rest of the mean planning time. In addition, automatically plan selection from the plan pools was performed in batch for 462 cases, and the corresponding time was negligible (4.2s for plan selection for 462 cases).

4 Discussion

Our recent studies (Cui *et al* 2018a, Cui *et al* 2018b) showed that it is possible to obtain a RTOG valid plan without any user interventions. In order to further increase the planning efficiency, four optimization engines were implemented and compared. Our results indicated that (1) gSA and gL-BFGS can speedup the optimization time by two or three orders of magnitude compared to their CPU implementation (figure 4a), (2) L-BFGS is equivalent to simulated annealing, and is not trapped in local minima (figure 2), (3) gL-BFGS is able to compute 10 000 plans within 9 s (optimization time in figure 4a), and (4) the multi-GPU approach is not necessary, considering the fact that the mean GPU memory usage to generate 10 000 plans was 2.6 GB out of 12 GB (figure 4b).

A new patient-specific approach called GPU-based MCO (gMCO) was proposed as an upgrade of our prior studies (Cui *et al* 2018a, Cui *et al* 2018b). gMCO can increase the RTOG acceptance rate from 97.5% (Cui *et al* 2018b) to 99.8%, and can decrease the planning time from 1 h (300 plans) (Cui *et al* 2018a), to 41 s (14 relevant plans) (Cui *et al* 2018b), to 9.4 s (1000 plans). Compared with the IPSA physician approved plans, gMCO can increase the RTOG+ acceptance rate by 27.3%, eliminating around 10 manual tweaking needed to achieve the observed clinical level based on the results presented in figure 5. For example, a RTOG invalid plan (urethra D_{10} above 118%) can be escalated to a RTOG+ valid plan by using gMCO. This has been made possible by relaxing the the bladder V_{75} dose (still below 1 cc), while still meeting all requirements for target, urethra and rectum dose parameters as shown in figure 7a. Such information can also be seen from the isodose curves in figure 7b-c. Note that in this study, the trade-off involved in the automatic selection scenario is based on selecting the highest target V_{100} while satisfying all the other RTOG criteria (figure 6). However, a high quality gMCO plan pool is available for the user to pick a plan that best suits the patient-specific conditions.

KBP and MCO are widely used patient-specific inverse planning algorithms. In KBP, clinical plans were used to extract the regression models based on geometric features. However, clinical plans are user-dependent (Das *et al* 2008), and may be inconsistent between centers (Chung *et al* 2008). On the other hand, gMCO is independent of these issues. In MCO, even though interpolations between calculated plans were usually used to achieve a high planning efficiency, ultra-fast planning remains a challenge since no parallelization scheme was implemented. In this study, it only takes 9.4 s to generate a high quality plan pool with gMCO. However, it is admitted that these comparisons were made by ignoring that the dwell times optimization in HDR brachytherapy is a relatively small scale problem compared to the fluence map optimization in EBRT.

Note that for the objective function considered in this work, the solution space is convex and it would be easy to dismiss SA in favor of the more computationally efficient gL-BFGS algorithm. While this objective function is popular in the field, other types of objective function might have more complex solution spaces. Therefore, having a

robust, albeit slower, MCO algorithm based on SA remains an essential tool.

We anticipate that the approach proposed in this study will be implemented in clinical systems as an adjunct tool. In future work, the application of gL-BFGS as well as gMCO to other HDR brachytherapy sites will be investigated.

5 Conclusion

Two GPU-based optimization engines were designed to calculate multiple plans in parallel. With the preferred engine, an ultra-fast patient-specific planning tool that is able to generate a high quality plan without any user interventions was proposed. After a validation over a large-scale patient cohort, both plan quality and planning efficiency can be significantly improved compared with the traditional planning in clinic.

Acknowledgement

This work was supported in part by the National Sciences and Engineering Research Council of Canada (NSERC) via the NSERC-Elektta Industrial Research Chair Grant (#484144-15), via the NSERC Discovery Grants (#355493 and #435510), and via the CREATE Medical Physics Research Training Network Grant (#432290). The authors acknowledge a scholarship from the Chinese Scholarship Council and partial support by the Canada Foundation for Innovation (#CFI30889). The authors acknowledge the supports from their colleagues, especially Louis Archambault, Paul Edimo, Andrea Frezza, Charles Joachim-Paquet, Frédéric Lacroix, Marie-Claude Lavallée, Ghyslain Leclerc, Loïc Paradis-Laperrière, Éric Poulin, and Nicolas Varfalvy.

References

- Chung H T, Lee B, Park E, Lu J J and Xia P 2008 Can All Centers Plan Intensity-Modulated Radiotherapy (IMRT) Effectively? An External Audit of Dosimetric Comparisons Between Three-Dimensional Conformal Radiotherapy and IMRT for Adjuvant Chemoradiation for Gastric Cancer *International Journal of Radiation Oncology*Biophysics* **71**(4), 1167 – 1174.
- Citrin D E 2017 Recent Developments in Radiotherapy *New England Journal of Medicine* **377**(11), 1065–1075.
- Craft D L, Halabi T F, Shih H A and Bortfeld T R 2006 Approximating convex Pareto surfaces in multiobjective radiotherapy planning *Medical Physics* **33**(9), 3399–3407.
- Cui S, Desprs P and Beaulieu L 2018a A multi-criteria optimization approach for HDR prostate brachytherapy: I. Pareto surface approximation *Physics in Medicine and Biology* **63**(20), 205004.
- Cui S, Desprs P and Beaulieu L 2018b A multi-criteria optimization approach for HDR prostate brachytherapy: II. Benchmark against clinical plans *Physics in Medicine and Biology* **63**(20), 205005.
- Das I J, Cheng C W, Chopra K L, Mitra R K, Srivastava S P and Glatstein E 2008 Intensity-Modulated Radiation Therapy Dose Prescription, Recording, and Delivery: Patterns of Variability Among Institutions and Treatment Planning Systems *JNCI: Journal of the National Cancer Institute* **100**(5), 300–307.

- Delaney G, Jacob S, Featherstone C and Barton M 2005 The role of radiotherapy in cancer treatment *Cancer* **104**(6), 1129–1137.
- Després P and Jia X 2017 A review of GPU-based medical image reconstruction *Physica Medica* **42**, 76 – 92.
- DeVita V T, Lawrence T S and Rosenberg S A 2015 *DeVita, Hellman, and Rosenberg extquotetelefts cancer: Principles and practice of oncology: Tenth edition* Wolters Kluwer Health Adis (ESP). ISBN 9781451192940.
- Edimo P, Kroshko A, Beaulieu L and Archambault L 2019 A stochastic frontier analysis for enhanced treatment quality of high-dose-rate brachytherapy plans *Physics in Medicine and Biology* **64**(6), 065012.
- Hazell I, Bzdusek K, Kumar P, Hansen C R, Bertelsen A, Eriksen J G, Johansen J and Brink C 2016 Automatic planning of head and neck treatment plans *Journal of Applied Clinical Medical Physics* **17**(1), 272–282.
- Jia X, Ziegenhein P and Jiang S B 2014 GPU-based high-performance computing for radiation therapy *Physics in medicine and biology* **59**(4), R151.
- Karabis A, Giannouli S and Baltas D 2005 HIPO: A hybrid inverse treatment planning optimization algorithm in HDR brachytherapy *Radiotherapy and Oncology* **76**, S29.
- Lahanas M, Baltas D and Giannouli S 2003 Global convergence analysis of fast multiobjective gradient-based dose optimization algorithms for high-dose-rate brachytherapy *Physics in medicine and biology* **48**(5), 599.
- Lahanas M, Schreiber E and Baltas D 2003 Multiobjective inverse planning for intensity modulated radiotherapy with constraint-free gradient-based optimization algorithms *Physics in Medicine and Biology* **48**(17), 2843.
- Lessard E and Pouliot J 2001 Inverse planning anatomy-based dose optimization for HDR-brachytherapy of the prostate using fast simulated annealing algorithm and dedicated objective function *Medical Physics* **28**(5), 773–779.
- Liu D C and Nocedal J 1989 On the limited memory BFGS method for large scale optimization *Mathematical Programming* **45**(1), 503–528.
- Men C, Gu X, Choi D, Majumdar A, Zheng Z, Mueller K and Jiang S B 2009 GPU-based ultrafast IMRT plan optimization *Physics in Medicine and Biology* **54**(21), 6565.
- Milickovic N, Lahanas M, Papagiannopoulou M, Zamboglou N and Baltas D 2002 Multiobjective anatomy-based dose optimization for HDR-brachytherapy with constraint free deterministic algorithms *Physics in Medicine and Biology* **47**(13), 2263.
- Moore K, Brame R S, Low D A and Mutic S 2011 Experience-Based Quality Control of Clinical Intensity-Modulated Radiotherapy Planning *International Journal of Radiation Oncology*Biography*Physics* **81**(2), 545 – 551.
- Nelms B E, Robinson G, Markham J, Velasco K, Boyd S, Narayan S, Wheeler J and Sobczak M L 2012 Variation in external beam treatment plan quality: An inter-institutional study of planners and planning systems *Practical Radiation Oncology* **2**(4), 296 – 305.
- Petit S F, Wu B, Kazhdan M, Dekker A, Simari P, Kumar R, Taylor R, Herman J M and McNutt T 2012 Increased organ sparing using shape-based treatment plan optimization for intensity modulated radiation therapy of pancreatic adenocarcinoma *Radiotherapy and Oncology* **102**(1), 38 – 44.
- Prax G and Xing L 2011 GPU computing in medical physics: A review *Medical Physics* **38**(5), 2685–2697.
- Radiation Therapy Oncology Group 2016 RTOG 0924 *Androgen deprivation therapy and high dose radiotherapy with or without whole pelvic radiotherapy in unfavourable intermediate or favourable high risk prostate cancer: A Phase III randomised trial* .
- Shen C, Gonzalez Y, Klages P, Qin N, Jung H, Chen L, Nguyen D, Jiang S B and Jia X 2018 Intelligent Inverse Treatment Planning via Deep Reinforcement Learning, a Proof-of-Principle Study in High Dose-rate Brachytherapy for Cervical Cancer *arXiv e-prints* p. arXiv:1811.10102.
- Teichert K, Sss P, Serna J I, Monz M, Kfer K H and Thieke C 2011 Comparative analysis of Pareto

- surfaces in multi-criteria IMRT planning *Physics in Medicine and Biology* **56**(12), 3669.
- van der Meer M, Alderliesten T, Pieters B, Bel A, Niatsetski Y and Bosman P 2018 in 'GECCO 2018 - Proceedings of the 2018 Genetic and Evolutionary Computation Conference' pp. 1387–1394.
- Wetzl J and Taubmann O 2013 'CudaLBFGS' <https://github.com/jwetzl/CudaLBFGS/>. Accessed: 2019-3-4.
- Wetzl J, Taubmann O, Haase S, Köhler T, Kraus M and Hornegger J 2013 in H.-P Meinzer, T. M Deserno, H Handels and T Tolxdorff, eds, 'Bildverarbeitung für die Medizin 2013' Springer Berlin Heidelberg Berlin, Heidelberg pp. 21–26.
- Wu B, Ricchetti F, Sanguineti G, Kazhdan M, Simari P, Chuang M, Taylor R, Jacques R and McNutt T 2009 Patient geometry-driven information retrieval for IMRT treatment plan quality control *Medical Physics* **36**(12), 5497–5505.
- Wu B, Ricchetti F, Sanguineti G, Kazhdan M, Simari P, Jacques R, Taylor R and McNutt T 2011 Data-Driven Approach to Generating Achievable DoseVolume Histogram Objectives in Intensity-Modulated Radiotherapy Planning *International Journal of Radiation Oncology*Biography*Physics* **79**(4), 1241 – 1247.
- Yamada Y, Rogers L, Demanes D J, Morton G, Prestidge B R, Pouliot J, Cohen G N, Zaider M, Ghilezan M and Hsu I C 2012 American Brachytherapy Society consensus guidelines for high-dose-rate prostate brachytherapy *Brachytherapy* **11**(1), 20 – 32.
- Zhou Y, Klages P, Tan J, Chi Y, Stojadinovic S, Yang M, Hrycushko B, Medin P, Pompos A, Jiang S, Albuquerque K and Jia X 2017 Automated high-dose rate brachytherapy treatment planning for a single-channel vaginal cylinder applicator *Physics in Medicine and Biology* **62**(11), 4361.