

Optimizing radiation therapy dose planning using deep reinforcement learning.

Journal:	<i>IEEE Journal of Biomedical and Health Informatics</i>
Manuscript ID	JBHI-01816-2020
Manuscript Type:	Paper
Date Submitted by the Author:	12-Nov-2020
Complete List of Authors:	Moreau, Grégoire; UCLouvain, ICTeam François-Lavet, Vincent; UCLouvain, ICTeam Desbordes, Paul; Université catholique de Louvain, ICTEAM Macq, Benoit; Université Catholique de Louvain, ICTEAM

SCHOLARONE™
Manuscripts

Optimizing radiation therapy dose planning using deep reinforcement learning

Grégoire Moreau, Vincent François-Lavet, Paul Desbordes, and Benoît Macq, *Fellow, IEEE*

Abstract—Radiation therapy treatments are delivered in fractions to slowly destroy a tumor while avoiding severe side effects in surrounding healthy tissues. Doctors nowadays largely design conventional treatment planning with heuristics. In order to develop a more principled approach to dose planning, this paper investigates how an algorithm can suggest a treatment plan of the doses in order to optimize some preference functions. Given a model of the tumor and surrounding tissues, a treatment is learned through deep reinforcement learning from high-dimensional inputs that show densities of healthy cells and cancer cells on a grid. The developed deep reinforcement learning approaches (based on DQN and DDPG) are able to find treatments that successfully optimize different preference functions.

Index Terms—Reinforcement learning, Automatic treatment planning, Cellular simulation

I. INTRODUCTION

IN external beam radiation therapy (RT), a dose of radiation is sent from outside the body into the tumor in the form of a beam of photons, electrons, or protons. The DNA of irradiated cells will be damaged by the radiation, which will cause some of the cells to die [1]. Because tumor cells have inhibited the DNA repair mechanism in favor of anarchic proliferation, they are more sensitive to radiation than healthy cells. To exploit this, the total radiation dose is delivered in fractions that should slowly destroy the tumor while allowing surrounding healthy tissues to recover from damage over time. The treatment schedule aims to optimize two criteria [2]:

- the Tumor Control Probability (TCP), which is the probability that the tumor is destroyed after the treatment and should be maximized.
- the Normal Tissue Complication Probability (NTCP), which is the probability that tissues and organs close to the tumor will suffer from side effects because of the treatment and should be minimized.

Currently, a dosimetry step is required prior to starting RT. A dose cartography, representing the optimal dose balancing the TCP and NTCP criteria, is defined by the medical physician. Then, a treatment planning system splits this total dose into several fractions (with variable angles, energies, etc.). A conventional treatment is to use daily fractions of equal dosage (e.g., 2 Gy for breast cancers [3]) until the total dose is administered, as this has proved to produce good clinical results. However, recent studies claim that there could be an advantage in varying the fraction size during the treatment [3].

All the authors were with the Institute of Information and Communication Technologies, Electronics and Applied Mathematics, UCLouvain, Louvain-la-Neuve, Belgium.

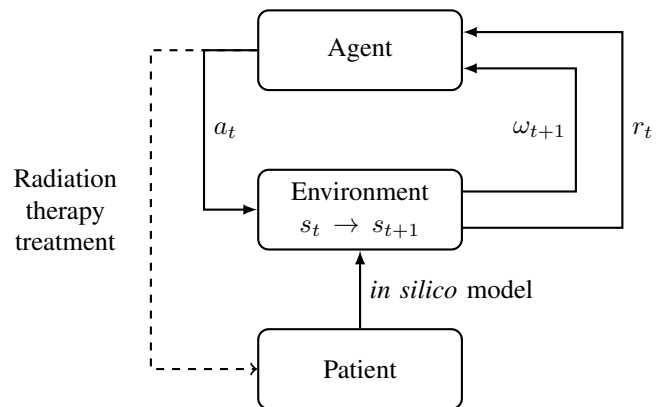


Fig. 1: Agent-environment interaction where the environment is an *in silico* simulation of a tumor within its surrounding tissues. At each time step t , the agent takes an action a_t , which represents the dose and the environment transitions from s_t to s_{t+1} , while providing a reward r_t as well as an observation ω_{t+1} .

Deep Reinforcement Learning (DRL) has been able to solve a large variety of complex sequential decision-making tasks, even from high-dimensional inputs [4]. It has achieved super-human performance in many domains such as in the game of go [5]. In this paper, we apply these techniques in order to automate the tedious task of treatment planning while also possibly obtaining less conventional fractionation schedules that might lead to improvements in current practice.

Ultimately the objective of radiation therapy is to improve the patient's life expectancy. That objective is however complex to assess and scientists and practitioners alike rely on different heuristics when planning treatment. We consider different possible reward functions that take the objectives of (i) maximizing the number of tumor cells killed, (ii) minimizing the number of healthy cells killed, and (iii) minimizing the total dose used in the treatment into account. Different reinforcement learning algorithms are then developed to demonstrate the potential of the approach.

To be optimized, a DRL algorithm needs a large amount of data. For this reason, a model of a tumor growing with surrounding healthy tissues is used as a proxy (see Fig. 1). Given a sufficiently accurate model, the results obtained can immediately be of interest as a more principled way to justify fractionation schedules, which can help understand the current practice. We also believe that this approach might ultimately be used to suggest adaptive treatments for each individual patient who requires radiation therapy. This would require developing

the automatic creation of an *in silico* model from CT scans on which treatment optimization could be performed.

II. RELATED WORKS

A. Mathematical modeling of cells and tumor development

Cell-based computational models can be used to simulate the behaviour of human cells and their possible developments. These *in silico* experiments (in opposition to *in vivo* and *in vitro*) are particularly useful in cancerology as they make it possible to study a tumor's proliferation and reactions to treatments without resorting to expensive and time-consuming clinical experiments.

Metzcar et al. [6] published a review of the different approaches that can be used to model the development of cancer. Those approaches can be split in two groups: lattice-based methods, which place cells on a grid, and off-lattice methods, which allow cells to have free positions. Several articles describing advanced simulations integrating details such as the development of a vascular architecture around the tumor during angiogenesis [7], anaerobic respiration using H+ [8], or interactions between tumor cells and their microenvironments [9] can be found in the literature. Because our approach requires a model efficient enough to provide relatively long simulations with a reasonable computational cost, we focused on the state-of-the-art simpler models.

O'Neil et al. [10] present a simulation of tumoral development in healthy tissue. Healthy cells and cancer cells proliferate semi-independently on a 2D grid containing some glucose. Both healthy and tumor cells go through the cell cycle and can replicate if they have access to a sufficient amount of nutrients. This cycle is split into four phases: eleven hours for Gap 1 (G1), eight hours for synthesis (S), four hours for Gap 2 (G2), and one hour for mitosis (M). Each hour represents a time step for the simulation (discrete time). Healthy cells can also enter quiescence (G0), if the conditions in the environment are not favorable. G0 cells consume less nutrients and do not replicate. Initially, healthy cells are placed randomly on the grid. Step after step, the cancer cells spread through the entire grid to form a tumor.

Radiation therapy is also implemented in the model to enable us to prevent the spread of tumor cells as well simulate side effects on healthy cells. It is represented by applying a survival probability of 37% to each cell along a certain radius of a center of radiation specified by the user. Cells that are not directly killed by the radiation keep count of the doses that they have received and have to repair the radiation damage over time. The radiation dose is then repeated every 24 hours of simulation. To keep the model as simple as possible, the RT is quite naive. Furthermore, the authors did not take dioxygen into account despite its major role in radiation therapy (such as radio-sensitivity of cells).

Later, Jalalimanesh et al. [11] extended this simulation. First, they replaced glucose by dioxygen as the main nutrient on the grid. Second, they used a modified Linear-Quadratic (LQ) model to represent the effects of different intensities of radiation on cells. The LQ model [12] aims to describe the

surviving fraction (SF) of a group of cells subjected to a dose d in grays (Gy):

$$SF(d) = \exp(-\alpha d - \beta d^2) \quad (1)$$

where α and β are parameters corresponding to the radiosensitivity of the irradiated tissue, which varies for each tumor subtype. In Jalalimanesh et al., the model is based on a modified LQ model representing the survival probability S of a cell irradiated with a dose d :

$$S(d) = \exp[\gamma_r(-\alpha \text{OMF } d - \beta(\text{OMF } d)^2)] \quad (2)$$

where γ_r is a factor that reflects the changes in the cell's radiosensitivity in different phases of the cell cycle and OMF depends on the oxygenation of the cell (as hypoxic cells are less sensitive to radiation). The effect of radiation can thus be implemented by going through each cell in the simulation, applying this model, and making it survive or die depending on the predicted survival probability.

While Jalalimanesh et al. [11]'s model is efficient enough to provide enough training samples for a deep learning algorithm, some important effects in the developments of cancers, such as angiogenesis, are omitted.

B. Reinforcement learning in radiation therapy

Reinforcement learning (RL) is a branch of machine learning in which an agent learns to optimize its decision-making from observations of its environment and feedback on the actions that it takes [15].

In RL algorithms, the feedback that the agent receives is a scalar called the reward. The agent aims to maximize the sum of rewards over time. It is important to design how those rewards are computed very carefully, as this defines the objective of the agent. This can be notably hard in cases such as radiation therapy, where multiple objectives need to be taken into account [16].

In Tseng et al. [17], a deep reinforcement learning algorithm, DQN, was used to learn dose fraction adaptation in radiotherapy. Based on an initial dataset of 114 patients who received radiotherapy for lung cancer, a generative adversarial network was used to create a larger amount of training data, while a deep neural network was trained on that data to estimate transition probabilities during the treatment and thus create a radiotherapy artificial environment for the RL agent. Finally a Deep-Q-Network was trained to estimate the value of each possible dose fraction. The reward function used is based on the idea of maximizing the probability of controlling the tumor while minimizing the probability of complications to the surrounding tissues.

In Jalalimanesh [11], the agent uses the tabular Q-learning algorithm where the considered actions are three possible doses of radiation (2.5, 3, and 3.5 Gy). As observations, the number of tumor cells in the simulations are discretized into 200 stages, with a step between two stages corresponding to an increase of 25 cells. As a reward, the agent receives a scalar corresponding to the difference between tumor and healthy cells killed directly by the chosen dose, minus a factor proportional to the dose and the irradiated area.

In this article, we use DRL and study different reward functions that capture different preference functions in the treatments. Deep learning allows the use of high-dimensional features such as images as inputs, instead of making arbitrary decisions to define a state space. For our experiments, we define two different environments: (i) a simple 1D toy environment and (ii) a more realistic 2D simulation. Finally, our source code is publicly available ¹, which we hope will help foster new research in this area.

III. MATERIALS AND METHODS

A. Tumor growth simulation

To simulate tumor development inside healthy tissue and the effect of radiation therapy, a lattice-gas cellular automaton was implemented. This simulation is based on a 2D 50-by-50 grid where each pixel of the grid contains:

- a list of cells,
- an amount of nutrients and
- (potentially) a nutrient source.

The same cell cycle as in [10] is used: eleven hours for G1, eight hours for S, four hours for G2, and one hour for M. Healthy cells also have the ability to enter the G0 phase if the quantity of a nutrient in the pixel is too low ($< q_{g_G0}$ for glucose or $< q_{o_G0}$ for oxygen) or the density of the patch formed by the pixel and its neighbors is greater than 1. The simulation is initialized with one tumor cell in the center of the grid and n_{h_cells} randomly-placed healthy cells.

The cells need two nutrients to develop: glucose and dioxygen, reactants of the aerobic respiration producing ATP that fuels the cells. Initially, a quantity q_{g_ini} glucose and q_{o_ini} dioxygen are placed for each pixel. Nutrients spread uniformly to the neighborhood at a rate of 20% per hour. Each cell's nutrient consumption is determined at the creation for healthy cells and varied every hour for cancer cells. These consumption rates are normally distributed based on averages ($\mu_{g_healthy}$, $\mu_{o_healthy}$, μ_{g_tumor} , μ_{o_tumor}) and standard deviation defined as a third of the average consumption. To generate nutrients during the simulation, $n_{sources}$ sources are randomly placed on the grid at the start of the simulation. Every hour, each of these sources will add q_g mg of glucose and q_o of dioxygen to the pixel on which they are placed. To simulate angiogenesis, the sources follow a random walk with a probability of moving towards the tumor's center proportional to the number of cells forming this tumor.

To represent the effect of radiation on cells during radiation therapy, the modified LQ model presented in Jalalimanesh et al. [11] is used (see Eq. 2). The dose is deposited using an accurate profile corresponding to the convolution product of a window function and a Gaussian. Thus, each cell is irradiated with a dose depending on its distance from the tumor's center of mass, with cells one tumor radius away from this center receiving 95% of the total dose. The tumor's current center of mass and radius are automatically computed after each irradiation.

In order to study first a simple environment for the RL algorithm, we built a toy environment where we collapsed

TABLE I: Growth model parameters.

Parameter	Notation	Value
Glucose quiescence level	q_{g_G0}	1.728 E-7 mg/cell
Oxygen quiescence level	q_{o_G0}	10.37 E-8 ml/cell
Starting healthy cells	n_{h_cells}	1000
Starting nutrient sources	$n_{sources}$	100
Starting glucose	q_{g_ini}	1 E-6 mg
Starting oxygen	q_{o_ini}	1 E-6 ml
Glucose average consumption	$\mu_{g_healthy}$	3.6 E-9 mg/cell/hour
	μ_{g_tumor}	5.4 E-9 mg/cell/hour
Oxygen average consumption	$\mu_{o_healthy}$	2.16 E-9 ml/cell/hour
	μ_{o_tumor}	2.16 E-9 ml/cell/hour
Glucose inputs per step	q_g	1.3 E-6 mg/source/hour
Oxygen inputs per step	q_o	4.86 E-7 ml/source/hour

the different elements of the grid to 0 dimensions. The same amounts of glucose and dioxygen distributed to the full grid initially and every hour in the 2D model are given to one pixel representing the whole environment. Initially, this pixel contains one tumor cell and healthy cells n_{h_cells} . At each step, cells consume the nutrients to spread. During the treatment, the same dose of radiation will be applied to every cell in the simulation.

Table I presents the parameters used in our models.

B. Reinforcement Learning

Reinforcement learning (RL) is a branch of machine learning that aims to solve the problem of sequential decision-making. In this framework, an agent receives observations of the state of the environment in which it evolves and can choose actions to affect this environment. If we can assume that the problem is fully observable, then it can be formalized as a Markov Decision Problem (MDP), which is defined by the following elements:

- States $s_t \in \mathcal{S}$, where the state space \mathcal{S} is,
- Actions $a_t \in \mathcal{A}$, where the action space \mathcal{A} is,
- Rewards $r_t \in [0, R_{max}]$, where R_{max} is set to 1 by scaling the rewards,
- A transition function $T(s, a, s') : \mathcal{S} \times \mathcal{A} \rightarrow P(\mathcal{S})$ that gives the probability of the next state s' given the current state s and action a ,
- A discount factor $\gamma \in [0, 1)$.

At each time step, the agent observes the state of the environment and chooses an action. The effect of this action causes a transition of the environment to a new state and the agent receives a reward in the form of a scalar that indicates the quality of this transition. The agent's goal is to maximize the sum of the rewards that it receives over time.

During its training, the agent will use trial and error to learn policies. The policy that an agent follows maps a state observed by the agent to the action that it should take (or to a distribution of actions).

The value of a state-action pair is the total amount of reward an agent can expect to accumulate over the future steps, starting from that state and taking this action according to the current policy. It can be defined mathematically with:

¹https://github.com/gregoire-moreau/radio_rl

$$q_{\pi}(s, a) \doteq E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \middle| s_t = s, a_t = a \right],$$

where the discount factor γ determines how much the immediate reward should be favored over future rewards.

A common problem encountered during the training of RL agents is the exploration/exploitation dilemma. To maximize the sum of rewards that it receives, the agent should exploit the best solutions that it has already found, but should also keep exploring the state and action spaces in order potentially to discover better strategies. A common strategy to address this problem is the ϵ -greedy algorithm, in which at each time step, the agent has a given probability ϵ of taking a random action instead of the best estimated action.

1) *Q-learning*: Value-based algorithms aim to solve the RL problem by learning an approximation of the optimal action-value function. This function has the form:

$$Q^*(s, a) = \max_{\pi} q_{\pi}(s, a)$$

Expressed in words, $Q^*(s, a)$ is the maximum discounted sum of rewards that could be expected under an optimal policy, assuming that the agent has just observed state s and taken action a . If the agent has access to this function, it can trivially reach optimal decision-making simply by choosing action $a^* = \arg\max_a Q^*(s, a)$ in each state s .

Tabular Q-learning uses a table containing an entry for each possible state-action pair to approximate Q^* . This means that the algorithm requires the state-action space to be discrete and small to produce accurate estimations. At each time step, the agent observes the current state s of the environment, chooses an action a to take depending on its learning strategy, and receives a reward r once the action has been applied. This process repeats at the next time steps when observing the new state s' that was reached. It will then update recursively the entry in the table corresponding to the initial state and the action that was chosen with:

$$Q[s, a] = (1 - \alpha)Q[s, a] + \alpha(r + \gamma \max_{a'} Q[s', a']),$$

where α is the step-size parameter that can be decreased during the training once the estimations get more accurate.

In deep Q-learning, a neural network is used to approximate Q^* . It receives observations of the environment's state as inputs, and estimates the value of each possible action as an output. At each step in the training algorithm, the values given by the neural network $Q(s, a | \theta^Q)$ for state s , action a , and parameterized θ^Q updated towards a target value given by the formula:

$$y_t = r_t + \gamma \max_{a' \in \mathcal{A}} Q(s', a' | \theta^Q)$$

In Deep Q-Networks [13], two additions help with the stability of the training. First, a large buffer of past experiences is kept so that the agent can train on batches sampled randomly from this buffer instead of using only the most recent examples. This trick, called experience replay, ensures that training

samples are more independent from each other. Second, the targets used to estimate the loss are computed using a second "target" network, to increase the stability of these targets. This network is obtained by cloning the main network once every C training steps where C is usually at least of the order of 10^4 .

2) *Deep deterministic policy gradient*: Policy-based algorithms aim instead to estimate the agent's policy directly, and to improve it over time. The Deep Deterministic Policy Gradient (DDPG) algorithm trains two different networks [14] simultaneously. The actor network approximates the function $\mu(s)$, which defines the agent's policy; it receives observations of the environment's state as inputs, and outputs the action that the agent should take. The critic network approximates the function $Q(s, a)$, which is used to evaluate the actions chosen by the agent. It receives both an observation of the environment's state, and the action that the agent chose in this state, and outputs an estimation of the expected return when taking this action in that state, while following $\mu(s)$ afterwards. Compared with DQN, DDPG has the possibility to work with a continuous action space.

The critic is trained to be increasingly accurate in its estimations of action-values. It uses a training target with the form:

$$y_i = r_i + \gamma Q(s', \mu(s' | \theta^{\mu}) | \theta^Q)$$

The actor is trained with a policy gradient. This gradient uses the critic and has the form:

$$\nabla_{\theta^{\mu}} J = \frac{1}{N} \sum_i \nabla_a Q(s, a | \theta^Q) | s = s_i, a = \mu(s_i) \nabla_{\theta^{\mu}} \mu(s | \theta^{\mu}) | s_i$$

The policy will move towards actions that lead to higher values. As the actions are directly predicted by the actor, the action space is continuous and can be multi-dimensional. This algorithm is therefore more flexible than DQN.

Experience replay and target networks are also used to train the critic in DDPG to improve stability.

C. Integration with the simulation

In order to use RL algorithms with the previously described models, the different elements of an MDP have to be defined in terms of these models.

Observations are used by the agent as the current state of the simulation. In the collapsed scalar model, this state is made out of combination of the number of cancer cells and healthy cells after they have been discretized. The DRL algorithms used with the 2D model allow for richer observations. Here, they are 50-by-50 images where each pixel's value is either the number of cancer cells on this pixel on the grid, if there are any, or the number of healthy cells on this pixel. Figure 2 shows such an observation.

The actions chosen by the agent define the RT treatment schedule applied to the simulation. For all algorithms, the agent will choose the dose with which to irradiate the tumor at each time step. This dose will be chosen in the interval [1-5] Gy. This domain is discrete for Q-learning algorithms, with steps of 1 Gy for tabular Q-learning and 0.5 Gy for DQN. It is continuous for DDPG. Furthermore, agents using

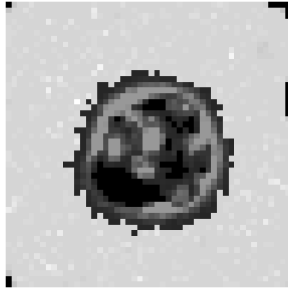


Fig. 2: Observation of the 2D model where light pixels are healthy cells, dark pixels are tumor cells, and black pixels are empty volumes. The intensities are proportional to the density of cells.

the DDPG algorithm will also choose the time between the current fraction and the next fraction, in a domain of [12-72] hours.

After the action chosen by the agent has been applied to the environment, it will receive a reward indicating the quality of this action. The agents have been trained to optimize two different reward functions: "Killed" (K) and "Killed and dose" (KD). On regular (non-terminal) time steps, agents will receive rewards of the form:

- Killed:

$$\frac{\text{cancer cells killed} - x \times \text{healthy cells killed}}{y}$$

- Killed and dose:

$$\frac{\text{cancer cells killed} - x \times \text{healthy cells killed}}{y} - \frac{\text{radiation dose}}{z}$$

In the 2D model, the values of the constants x , y , and z are 5, 100,000, and 200 respectively.

On terminal steps, agents will receive rewards of -1 if the treatment was unsuccessful, or of the form:

- Killed: $0.5 - \text{difference between the initial and final number of healthy cells}/k$
- Killed and dose: $0.5 - \text{difference between the initial and final number of healthy cells}/k - \text{radiation dose}/z$

In the 2D model, the values of the constants k and z are 3,000 and 200 respectively.

A simulated treatment is considered finished in three possible situations. First, the treatment will have succeeded if the tumor has been destroyed and no cancer cells are left in the model. Second, the treatment will have failed if the tumor has completely invaded healthy tissue, and there are fewer than 10 healthy cells in the model. Third, the treatment will have timed out if the two previous conditions are not reached after 1200 hours of treatment.

D. Evaluation of performances

For each combination of environment, algorithm, and reward function, 1000 test treatments are generated. The performances of each treatment are assessed using several indicators:

- TCP
- Total radiation dose (D_{mean})
- Number of fractions (N_{fract})
- Treatment time (T_{mean})

TABLE II: Results of reinforcement learning trained on the collapsed toy environment using the tabular Q-learning algorithm (mean \pm standard error).

Rewards*	TCP (%)	D_{mean} (Gy)	N_{fract}	T_{mean} (h)
Baseline	97	77.9 ± 0.3	39 ± 0.2	934.8 ± 3.8
K	100	66.7 ± 0.3	32.2 ± 0.1	771.9 ± 3.3
KD	100	61.2 ± 0.3	27.7 ± 0.1	664.4 ± 3.6

*K for killed, KD for killed and dose

To put these performance indicators in context, they will be compared with a baseline corresponding to a dose of 2 Gy every 24 hours, as in a conventional RT treatment. The sample mean of each performance indicator is given along with its standard error when relevant. Furthermore, to illustrate the policies found by agents to optimize their reward functions, a graph will give the mean and standard deviation of the radiation dose prescribed at each step of the treatment.

IV. RESULTS

A. Collapsed toy environment

Table II presents the performances of the agents trained on the collapsed toy environment using the tabular Q-learning algorithm.

We observe that for both reward functions (K and KD), a better TCP can be obtained using fewer fractions and lower radiation doses. While the baseline treatment needed 39 fractions and an average dose of 77.9 Gy to reach a TCP of 97%, agents trained to optimize the K (and the KD) reward function need 32.2 (and 27.7) fractions and an average dose of 66.7 Gy (and 61.2 Gy) to reach a TCP of 100% (100%).

Figure 3a shows the shape of the treatments used on the toy environment by the agent that optimized the "Killed" reward function. We can see that three high radiation doses are delivered at the beginning to eliminate a large portion of the cancer cells. Then, lower doses are used to destroy the remaining cancer cells while allowing healthy cells to recover. Figure 3b shows the shape of the treatments used on the toy environment by the agent that optimized the "Killed and dose" reward function. The agent always starts with a dose of 1 Gy. In the few following steps, the observed doses have a very high standard deviation, indicating a high sensitivity to the situation. The agent will then tend to destroy the remaining cancer cells in the model using doses of around 2 Gy.

B. 2D environment

Table III presents the performances of the agents trained on the 2D model using the DQN and DDPG algorithms.

Here again, we observe that for both combinations, a better TCP (100% versus 99%) can be obtained using a lower number of fractions and a lower dose. The shorter treatment is the one obtained using DQN and the "Killed and dose" reward function (8.2 fractions, 32.1 Gy). It uses half the dose sent by the baseline and one-fourth as many fractions are required.

Figure 4d shows the shape of the treatments used on the 2D environment by the agent that optimized the "Killed" and

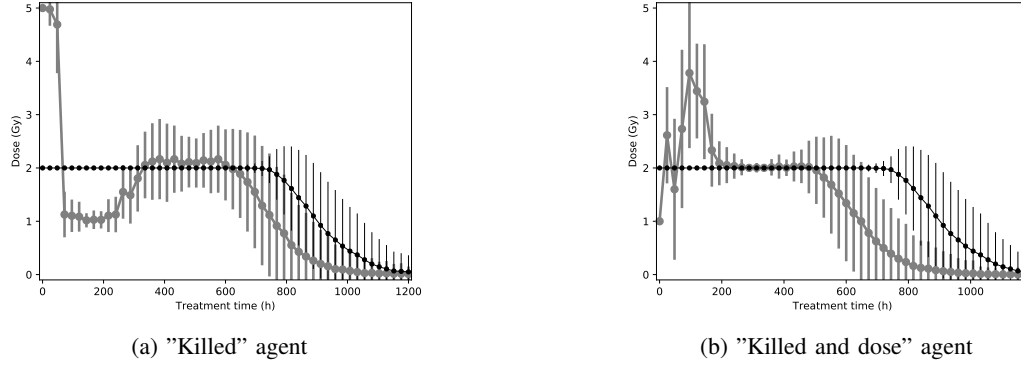


Fig. 3: Shape of treatments found by the agent for the collapsed toy environment (in gray) compared with the baseline treatment (in black). Each point represents a dose delivered to the patients.

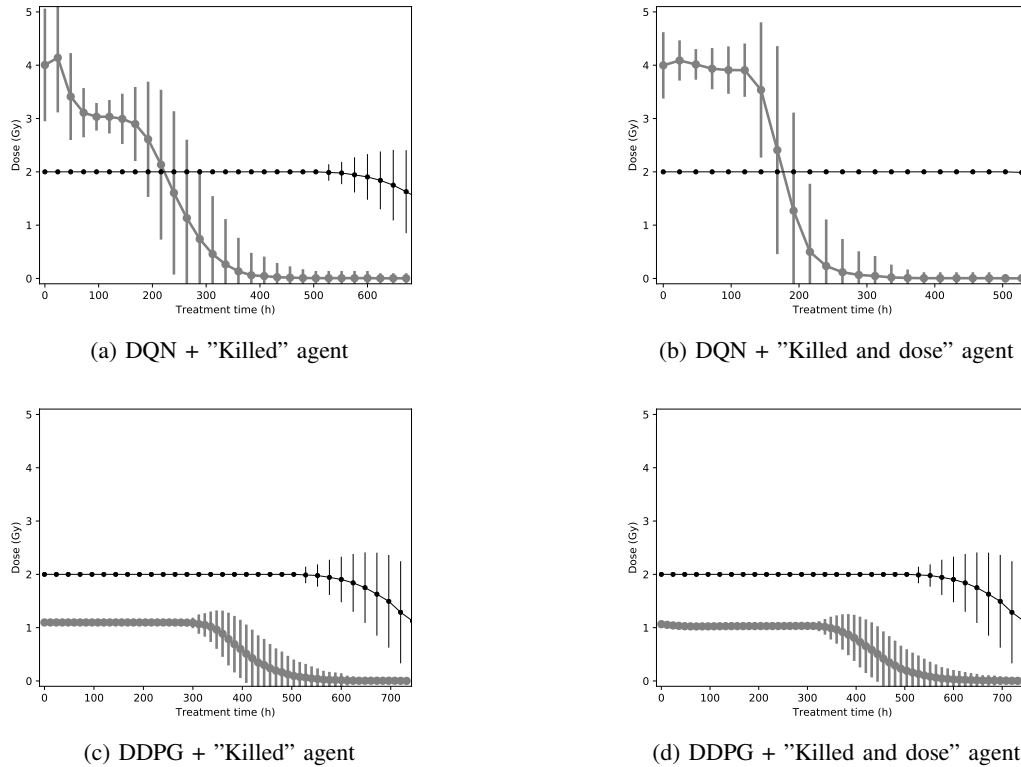


Fig. 4: Shape of treatments found by the agents for the 2D environment (in gray) using the DQN algorithm (top) and the DDPG algorithm (bottom). The baseline is shown (in black) for comparison. Each point represents a dose delivered to the patients.

"Killed and dose" reward functions using the DQN and the DDPG algorithms.

Concerning the DQN algorithm with the "Killed" function (see Fig. 4a), the agent sends a few high doses at the start of the treatment, then destroys the remaining cancer cells with daily doses of around 3 Gy. For the "Killed and dose" function (see Fig. 4b), the agent chooses to send high doses throughout the treatments. Figure 5 shows the effects of such treatments on the simulation.

Figure 5 shows the effects of treatments suggested by the agents controlled by the DQN algorithm on the simulation.

The first few doses of the treatments impact a large share of the tumor and surrounding tissues. After 120 hours, only a few cancer cells remain and healthy tissue is starting to recover from the initial radiation damage. Finally, at the end of the treatments, there are no more cancer cells on the grid while healthy tissue has completely recovered from the treatment.

Concerning the DDPG algorithm with the "Killed" function (see Fig. 4c), the agent chooses to deliver a dose slightly higher than 1 Gy every 12 hours for the entire treatment. For the "Killed and dose" function (see Fig. 4d), like the other DDPG agent, this agent chooses to deliver a dose slightly higher than

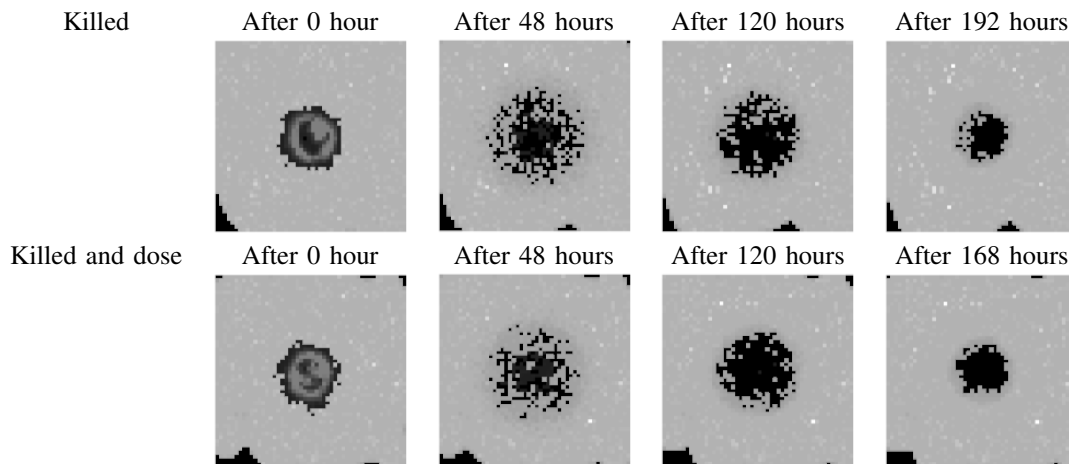


Fig. 5: Effects of the treatments prescribed by the DQN agents on the 2D simulation using the two studied reward functions (K and KD). Light pixels are healthy cells, dark pixels are tumor cells, and black pixels are empty volumes.

TABLE III: Results of reinforcement learning trained on the 2D environment using two different RL algorithms and two different reward functions (mean \pm std).

		TCP (%)	D _{mean} (Gy)	N _{fract}	T _{mean} (h)
Baseline	/	99	65.6 ± .3	32.8 ± .2	787.7 ± 4
DQN	K	100	35.6 ± .2	10.9 ± .1	262.2 ± 1.6
DQN	KD	100	32.1 ± .2	8.2 ± .1	195.8 ± 1.4
DDPG	K	100	38.5 ± .2	35.1 ± .2	420.9 ± 2.1
DDPG	KD	100	39.4 ± .2	38.1 ± .2	457.1 ± 2.2

* K for killed, KD for killed and dose

*K for killed, KD for killed and dose

1 Gy every 12 hours throughout the treatment. Figure 6 shows the effect of such treatments on the simulation.

Figure 6 shows the effects of treatments suggested by agents using the DDPG algorithm on the simulation. The low but frequent doses initially stop the tumor's growth and cause it some damage while mostly sparing healthy tissue. Two-thirds of the way into the treatments, cancer cells have mostly been destroyed and only healthy cells close to the initial borders of the tumor have been impacted by the radiation. At the end of the treatments, healthy tissue has mostly recovered the space invaded by the tumor before these treatments started.

V. DISCUSSION

Our experiments have shown that using reinforcement learning is a promising option for the optimization of treatment planning for external radiation therapy. As suggested by Jalalimanesh et al. [11], a model simulation is a good alternative to model the tumor growth necessary for RL training. Because this kind of model is flexible, we chose to complexify it compared with the previous articles by studying the impact of both glucose and oxygen in the same model and by modeling effects such as angiogenesis. Furthermore, we studied the performances of several RL algorithms (tabular Q-learning, DDPG and DQN). In particular, the different deep RL approaches allowed us to consider different action spaces, namely, (i) a fixed set of actions (DQN) and (ii) a continuous action space (DDPG).

It is necessary to specify that finding the optimal policy is difficult because of the model's complexity, in particular its high-dimensional state space, stochasticity, and nonlinear dynamics. The impact of several random parameters is also quite important (e.g., initial positions of healthy cells, nutrient consumption of cells, positions of daughter cells, probability of linear survival for large numbers but not when only a few cancer cells remain, repair time in case of survival, and sources, positions and displacement). This point explains the significant standard deviation in our performance indicators. But even despite this, our RL algorithms could converge.

Our deep RL agents are able to obtain better treatment schedules than a simple baseline with (i) improved TCP, (ii) shorter treatments, and (iii) lower overall radiations. The improvement is particularly important for the 2D environment. In our simulations, we have observed that starting with relatively high doses and decreasing the radiation doses through the treatment improves the cumulative rewards. We have also observed that hyper-fractionation might lead to better treatments. Further studies with an even more realistic *in silico* model or clinical trials would be of interest to validate these observations.

The two different reward functions produce slightly different results when measured with performance indicators. Indeed, the added dose factor in the "Killed and dose" reward function creates agents that suggest on average shorter treatments with a lower radiation dose and number of fractions than the "Killed" reward function. This effect is observed for the tabular Q-learning and DQN algorithms.

Agents trained with the DDPG algorithm reached the same treatment patterns while optimizing both reward functions. These treatments (sending close to 1 Gy every 12 hours) always use actions very close to the minimum for each dimension of the action space. Different treatments might thus be found if the action space is extended below the current minimum.

In addition, our study has some limitations. First, our models are not realistic enough. For instance, the model overestimates healthy tissue recovery. As shown in Figures 5

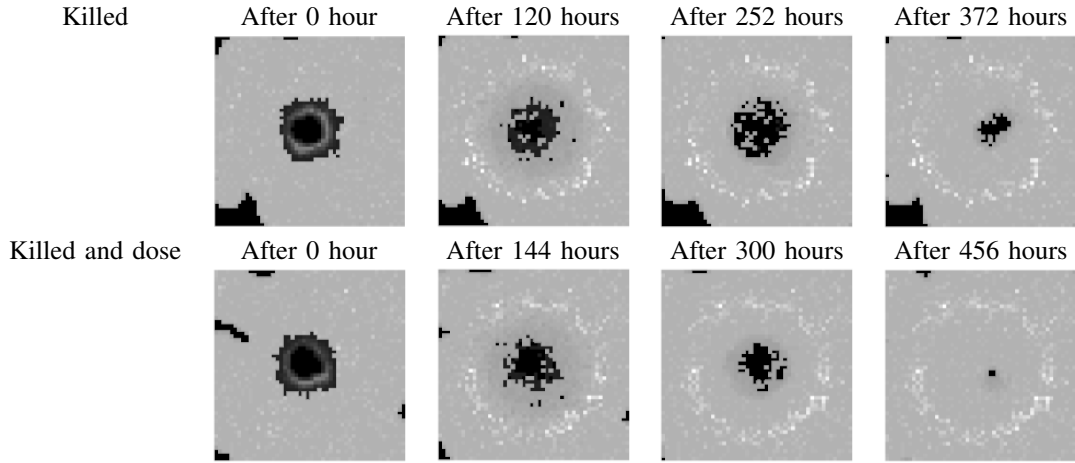


Fig. 6: Effects of the treatments prescribed by the DDPG agents on the 2D simulation using the two studied reward functions (K and KD). Light pixels are healthy cells, dark pixels are tumor cells, and black pixels are empty volumes.

and 6, healthy tissue quickly recovers from the radiation and reclaims space left empty after the destruction of the tumor. The model instantly deletes cells that are considered to have been killed by radiation. We also used an "almost" perfect dose profile. The radiation beam is able to target tumor cells accurately while avoiding healthy cells.

Future work in several directions is needed to bring this research work closer to a useful software for the doctors: (i) The *in silico* model developed for this paper is rather generic and could benefit from an improved cell model that would take the specificities of the different organs and different tumor types into account. (ii) Medical doctors should help develop improved heuristics in terms of the rewards optimized through reinforcement learning, while data from clinical trials should ultimately allow for optimizing directly for the improvement in patient life expectancy. (iii) A transfer learning approach from CT scans to the *in silico* model should be developed in order to have personalized treatment planning in real time.

VI. CONCLUSION

This paper has studied the use of deep reinforcement learning agents on a simulation of curing cancer through radiation therapy. A generic mathematical model of healthy tissues and tumor growth was developed and the dose planning algorithm took damage to both healthy tissues and the tumor into account in order to optimize the fractions. Two different deep RL algorithms were trained, each with two different preference functions. The treatments found by the agents were compared with a baseline corresponding to what is conventionally used in radiation therapy today and the deep RL approaches were able to outperform this baseline in terms of the optimized criteria.

Our approach is aimed at automating parts of treatment planning, which is fully designed by a human expert nowadays. Thanks to radiomics features extracted from CT scans and a biological *in silico* model of the radiation effects, an RL approach could pave the way to adaptive treatments that are able to take the specifics of each patient's case (tumor shrinkage, deformations, etc.) into account.

REFERENCES

- [1] C. Washington, D. Leaver and M. Trad, *Principles and Practice of Radiation Therapy (5th ed.)*, Mosby, 2020.
- [2] M. Baumann and C. Petersen, "TCP and NTCP: a basic introduction", *Rays*, vol. 30, no. 2, pp. 99-104, 2005.
- [3] J. Haviland, J. Owen, J. Dewar, et al., "The UK Standardisation of Breast Radiotherapy (START) trials of radiotherapy hypofractionation for treatment of early breast cancer: 10-year follow-up results of two randomised controlled trials", *Lancet Oncol.*, vol. 14, pp. 1086-94, 2013.
- [4] V. Francois-Lavet, P. Henderson, R. Islam, et al., "Introduction to Deep Reinforcement Learning", *Found. Trends Mach. Learn.*, vol. 11, no. 3-4, pp. 219-354, 2018.
- [5] D. Silver, A. Huang, C. Maddison, et al., "Mastering the game of Go with deep neural networks and tree search", *Nature*, no. 529, pp. 484-489, 2016.
- [6] J. Metzcar, Y. Wang, R. Heiland, et al., "A Review of Cell-Based Computational Modeling in Cancer Biology", *JCO Clin. Cancer Inform.*, no. 3, pp. 1-13, 2019.
- [7] S. McDougall, A. Anderson and M. Chaplain, "Mathematical modelling of dynamic adaptive tumor-induced angiogenesis: Clinical implications and therapeutic targeting strategies", *J. Theor. Biol.*, no. 241, vol. 3, pp. 564-589, 2006.
- [8] K. Smallbone, R. Gatenby, R. Gillies, et al., "Metabolic changes during carcinogenesis: Potential impact on invasiveness", *J. Theor. Biol.*, no. 244, vol. 4, pp. 703-713, 2007.
- [9] A. Anderson, K. Rejniak, P. Gerlee, et al., "Microenvironment driven invasion: a multiscale multimodel investigation", *J. Math. Biol.*, vol. 58, no. 4, pp. 579, 2008.
- [10] N. O'Neil, *An Agent Based Model of Tumor Growth and Response to Radiotherapy*, Virginia Commonwealth University, Richmond, Virginia, 2012.
- [11] A. Jalalimanesh, H. Shahabi Haghighi, A. Ahmadi, et al., "Simulation-based optimization of radiotherapy: Agent-based modeling and reinforcement learning", *Math. Comput. Simulat.*, vol. 133, pp. 235-248, 2017.
- [12] C. van Leeuwen, A. Oei, J. Crezee, et al., "The alfa and beta of tumors: A review of parameters of the linear-quadratic model, derived from clinical radiotherapy studies", *Radiat. Oncol.*, vol. 13, num. 1, pp. 96, 2018.
- [13] V. Mnih, K. Kavukcuoglu, D. Silver, et al., "Human-level control through deep reinforcement learning", *Nature*, vol. 518, no. 7540, pp. 529-533, 2015.
- [14] T. P. Lillicrap, J. J. Hunt, A. Pritzel, et al., "Continuous control with deep reinforcement learning", *arXiv. 1509.02971*, 2015.
- [15] R. Sutton, Richard S and A.G. Barto, *Introduction to reinforcement learning*, MIT press Cambridge, 135, 1998.
- [16] C. Yu, J. Liu and S. Nemati, "Reinforcement Learning in Healthcare: A Survey", *arXiv. 1908.08796*, 2020.
- [17] HH. Tseng, Y. Luo, S. Cui, et al., "Deep reinforcement learning for automated radiation adaptation in lung cancer", *Med. Phys.*, vol. 44, no. 12, pp. 6690-6705, 2017.