



**Citation:** Wang M, Scott JG, Vladimirsky A (2024) Threshold-awareness in adaptive cancer therapy. PLoS Comput Biol 20(6): e1012165. https://doi.org/10.1371/journal.pcbi.1012165

**Editor:** Ville Mustonen, University of Helsinki, FINLAND

Received: April 28, 2023
Accepted: May 9, 2024
Published: June 14, 2024

Copyright: © 2024 Wang et al. This is an open access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data Availability Statement:** The full source code for all computational experiments can be found in the following public repository: https://github.com/eikonal-equation/Stochastic-Cancer.

**Funding:** MW and AV are grateful to the National Science Foundation, Division of Mathematical Sciences (1645643, 1738010, and 2111522) for supporting this project. JGS would like to thank the National Institutes of Health, Division of Cancer Epidemiology and Genetics, National Cancer Institute (5R37CA244613, 5U54CA274513, U01CA280829) and the American Cancer Society (132691-RSG-20-096-01) for their generous

RESEARCH ARTICLE

# Threshold-awareness in adaptive cancer therapy

MingYi Wang<sup>1</sup>, Jacob G. Scott<sup>2</sup>, Alexander Vladimirsky<sub>0</sub><sup>3</sup>\*

- 1 Center for Applied Mathematics, Cornell University, Ithaca, New York, United States of America, 2 Department of Translational Hematology and Oncology Research, Cleveland Clinic, Cleveland, Ohio, United States of America, 3 Department of Mathematics and Center for Applied Mathematics, Cornell University, Ithaca, New York, United States of America
- \* vladimirsky@cornell.edu

# **Abstract**

Although adaptive cancer therapy shows promise in integrating evolutionary dynamics into treatment scheduling, the stochastic nature of cancer evolution has seldom been taken into account. Various sources of random perturbations can impact the evolution of heterogeneous tumors, making performance metrics of any treatment policy random as well. In this paper, we propose an efficient method for selecting optimal adaptive treatment policies under randomly evolving tumor dynamics. The goal is to improve the cumulative "cost" of treatment, a combination of the total amount of drugs used and the total treatment time. As this cost also becomes random in any stochastic setting, we maximize the probability of reaching the treatment goals (tumor stabilization or eradication) without exceeding a prespecified cost threshold (or a "budget"). We use a novel Stochastic Optimal Control formulation and Dynamic Programming to find such "threshold-aware" optimal treatment policies. Our approach enables an efficient algorithm to compute these policies for a range of threshold values simultaneously. Compared to treatment plans shown to be optimal in a deterministic setting, the new "threshold-aware" policies significantly improve the chances of the therapy succeeding under the budget, which is correlated with a lower general drug usage. We illustrate this method using two specific examples, but our approach is far more general and provides a new tool for optimizing adaptive therapies based on a broad range of stochastic cancer models.

# Author summary

Tumor heterogeneities provide an opportunity to improve therapies by leveraging complex (often competitive) interactions of different types of cancer cells. These interactions are usually stochastic due to both individual cell differences and random events affecting the patient as a whole. The new generation of cancer models strive to account for this inherent stochasticity, and *adaptive* treatment plans need to reflect it as well. In optimizing such treatment, the most common approach is to maximize the probability of eventually stabilizing or eradicating the tumor. In this paper, we consider a more nuanced version of success, maximizing the probability of reaching these therapy goals before the

support. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

**Competing interests:** The authors have declared that no competing interests exist.

cumulative burden from the disease and treatment exceed a chosen threshold. Importantly, our method allows computing such optimal treatment plans efficiently and for a range of thresholds at once. If used on a high-fidelity personalized model, our general approach could potentially be used by clinicians to choose the most suitable threshold after a detailed discussion of a specific patient's goals (e.g., to include the trade-offs between toxicity and quality of life).

### 1 Introduction

Optimizing the schedule and composition of drug therapies for cancer patients is an important and active research area, with mathematical tools often employed to improve the outcomes and reduce the negative side effects. Tumor heterogeneity is increasingly viewed as a key aspect that can be leveraged to improve therapies through the use of optimal control theory [1]. Most researchers using this perspective focus on deterministic models of tumor evolution, with typical optimization objectives of maximizing the survival time [2], minimizing the tumor size [3], or minimizing the time until the tumor size is stabilized [4]. In models that address the stochasticity in tumor evolution, a typical optimization goal is to find treatment policies that maximize the likelihood of patient's eventual cure (e.g., [5–7]) or minimize the likelihood of negative events (e.g., metastasis) after specified time [8]. However, this ignores the need for more nuanced treatment policies that maximize the likelihood of different levels of success—e.g., the probability of reaching remission or tumor stabilization without exceeding the specified amount of drugs and/or the specified treatment duration. The primary goal of this paper is to introduce a rigorous and computationally efficient approach that addresses such challenging objectives in stochastic cancer models.

The advent of personalized medicine in cancer has changed the way we think about therapy for patients whose tumors have actionable mutations. This has been a game changer for some patients, drastically increasing life spans, reducing toxicity and improving quality of life. Frustratingly, however, this population of patients is still small; it was estimated in 2020 that only ≈5% of patients benefit from these targeted therapies [9]. Further, despite the many advantages of personalized therapies, they rarely, if ever, lead to a complete cure since tumors develop resistance through the process of Darwinian evolution [10]. In response to this realization, a new approach called "evolutionary therapy" seeks to use the evolutionary dynamics of diseases to alter therapeutic schedules and drug choices. Through a combination of mathematical and experimental modeling, investigators have worked to understand a range of theoretical questions of practical importance. E.g., how does the emergence of resistance to one drug affect the sensitivity to another? Do heterogeneous (phenotypically or genotypically mixed) populations within tumors respond to drugs differently depending on their current state? The insights gained in these investigations have already led to progress in rational drug ordering/ cycling for bacterial infections [11–14] as well as for a number of cancers [15, 16]. In the study of therapeutic scheduling, adaptive therapy, which uses mathematical tools from Evolutionary Game Theory (EGT), has shown promise not only in theory [17], but also in a phase 2 trial for men with metastatic prostate cancer [18]. Experimentally, there have been confirmations of EGT principles in vivo [19] as well as more quantitatively focused assay development in vitro [20], and observations of game interactions using these methods [21]. There are also many other models capturing the competition within heterogeneous tumors without using gametheoretic derivations; e.g., [3, 22-24]. The majority of theoretical work in this space has focused on optimization of different drug regimens for deterministic models of cancer

evolution [25–29]. In contrast, our goal here is to provide efficient computational tools for nuanced therapeutic scheduling in cancer models that directly account for stochastic perturbations.

Cancers (and other populations of living things) are comprised of individual cells (or organisms) with their own behaviours and evolutionary histories. Stochastic phenomena are ubiquitous in their interactions and life histories. These include individual genetic differences, fate transitions [30], varying reactions to drugs [31], differences in signalling, and small-scale variations in the tumor environment. Many of these are instances of demographic stochasticity [32], which often can be "averaged-out" when dealing with a sufficiently large population. Indeed, this notion is crucial for any description of tumor heterogeneity through splitting the cells into sub-populations. Such splitting is natural if the mutation-selection balance is tuned so that only closely related genotypes, encoding the same phenotype, will stably exist. These groups are also referred to as quasispecies [33, 34] and exist as distributions around a central genotype, with all cells in the group behaving in a similar manner despite random birth/death events [32, 35] and small within-the-group genetic heterogeneities [36]. In contrast, our focus here is on environmental stochasticity, which cannot be ignored even in large populations since it describes random events that simultaneously affect the entire groups. Such perturbations are typically external [32, 35]; e.g., for cancer they might result from therapy-unrelated drugs or from frequent small changes in the host's physiology. Of course, any such event will also cause varying responses of individual cells within each subpopulation; so, our use of the term "environmental stochasticity" should be interpreted as direct modeling of subpopulation-averaged responses to such system-wide perturbations.

Modeling such perturbations in continuous-time usually results in Stochastic Differential Equations (SDEs) [37, 38], whose behavior can be optimized using Stochastic Optimal Control Theory [39]. The latter provides a mathematical framework for handling sequential decision making (e.g., how much drug to administer at each point in time) under random perturbations (e.g., stochastic changes in respective fitness of competing subpopulations of cancer cells). Any fixed treatment strategy will result in a random tumor-evolutionary trajectory and a random cumulative "cost" (e.g., cumulative amount of drugs used, or time to recovery, or a combination of these two metrics). The key idea of *Dynamic Programming* (DP) is to pose equations for the cumulative cost of the optimal strategy and to recover that strategy in feedback form: i.e., decisions about the dose and duration of therapy are frequently re-evaluated based on the current state of the tumor instead of selecting a fixed time-dependent treatment schedule in advance. This idea is applicable across a wide range of cancer models and therapy types, including those intended to stabilize the tumor and those aiming to eradicate it. We follow this approach here, but with an important caveat: instead of selecting an on-average optimal strategy (e.g., the one which minimizes the expected cost of treatment) as would be usual in stochastic DP, we select a strategy maximizing the probability of some desirable outcome (e.g., reaching the goals of the therapy without exceeding a specific cost threshold). The resulting riskaware (or, more precisely, "threshold-aware") policies are designed to be adaptive, adjusting the treatment plan along the way based on the responsiveness of tumor to drugs already used (and the cost already incurred) so far. In contrast to standard methods of constrained stochastic optimal control, our approach makes it easy to compute such threshold-aware policies for a range of thresholds simultaneously.

As is often the case, there remains a significant gap between simplified mathematical models and clinical applications. Much work remains in refining and calibrating EGT models, and also in measuring different aspects of biological stochasticity. But once high-fidelity personalized models become available, our general approach could potentially be used by clinicians to

choose the most suitable threshold after a detailed discussion of a specific patient's goals (to include the trade-offs between toxicity and quality of life, for example).

#### 2 Methods and models

To emphasize the broad applicability of our "risk-aware" adaptive therapy optimization approach, we first describe it for a fairly generic cancer model. Two specific examples are then studied in detail in §2.2 and §2.3.

# 2.1 Traditional and risk-aware control in drug therapy optimization

We note that most of the literature on dynamic programming in cancer models starts with positing a specific known/fixed treatment horizon T, with the success or failure of therapy assessed after that time (or earlier, in case of the modeled patient's death). This makes it easier to use the standard equations and algorithms of "finite-horizon" optimal control theory. But such a pre-determined T is not well-aligned with the notion of adaptive therapies. Instead, we adopt the *indefinite-horizon* framework, in which the process terminates as soon as the tumor's state satisfies some predefined conditions, with the terminal time T thus dependent on the chosen treatment policy. We use this framework in all of the control approaches described below, even though many of them have direct finite-horizon analogs as well.

We begin by describing several "traditional" optimal control formulations, followed by the threshold-aware version, which addresses some of their shortcomings in cancer applications. Starting with the deterministic setting summarized in Box 1, we use  $x(t) \in \mathbb{R}^n$  to encode the time-dependent state of a tumor (e.g., this could be the size or the relative abundance of n different sub-types of cancer cells). Tumor dynamics are modeled by an ODE system  $\dot{x} = f(x, d)$ , where the rate function f takes as inputs both the current state x(t) and the current control, the "therapy intensity" d(t). In models with a single drug, this is just a scalar  $d(t) \in [0, 1]$  $d_{\text{max}}$ ] indicating the current rate of that drug's delivery, where  $d_{\text{max}}$  encodes the Maximum Tolerated Dose (MTD), which can in principle be patient-specific. But the same framework can also be used for multiple drugs, with a separate upper bound specified for each element of d(t). Given an initial tumor configuration  $x(0) = \xi$ , a successful therapy aims to drive the tumor state to a set  $\Delta_{\text{succ}}$  while ensuring that the set  $\Delta_{\text{fail}}$  is avoided. E.g., in eradication therapy models,  $\Delta_{\text{succ}}$  might correspond to tumors below the detection level, while  $\Delta_{\text{fail}}$  might specify a much larger size that effectively kills a patient; see §2.3. On the other hand, for models that only track the relative abundance of cancer subpopulations,  $\Delta_{\text{succ}}$  might be defined in terms of the desired low abundance of specific subpopulations affected by d(t), with the idea that the tumor size stabilizes or an entirely different therapy strategy is adopted after x(t) enters  $\Delta_{\text{succ}}$ ; see §2.2.

If the therapy manages to reach  $\Delta_{\text{succ}}$  while avoiding  $\Delta_{\text{fail}}$ , its overall "cost"  $\mathcal{J}$  is assessed by integrating some running cost K = K(x, d) along the "trajectory" from  $\boldsymbol{\xi}$  to  $\Delta = \Delta_{\text{succ}} \cup \Delta_{\text{fail}}$  and adding the "terminal cost" g depending on its final state. E.g., g might be defined as  $+\infty$  on  $\Delta_{\text{fail}}$  to make such outcomes unacceptable. The running cost K depends on the current state of the tumor and the current drug usage levels and can be used to model the direct impact on the patient of the tumor size and composition as well as the side effects of the therapy. The key idea of dynamic programming [40] is to define a *value function*  $u(\boldsymbol{\xi})$  encoding the minimal overall cost for each specific initial tumor state and to show that this u must satisfy a stationary Hamilton-Jacobi-Bellman equation (2.4). Once that partial differential equation (PDE) is solved numerically, the globally optimal rate of treatment can be obtained in *feedback form* for all cancer states (i.e.,  $d = d_{\star}(\boldsymbol{\xi})$ ), which makes it suitable for the adaptive therapy framework. Throughout this paper, we will use  $d_{\star}(\boldsymbol{\xi})$  to denote an optimal feedback policy for the

# Box 1: Problem setup of a typical deterministic optimal cancercontrol problem

Evolutionary dynamics with control on therapy intensity:

$$\begin{cases} \dot{\mathbf{x}} = f(\mathbf{x}, \mathbf{d}), \\ \mathbf{x}(0) = \xi. \end{cases}$$
 (2.1)

 $\textbf{Process } \textit{terminates} \text{ as soon as either } \begin{cases} \textit{\textbf{x}}(t) \in \Delta_{\text{succ}}, & \quad \text{if therapy succeeds;} \\ \textit{\textbf{x}}(t) \in \Delta_{\text{fail}}, & \quad \text{if therapy fails.} \end{cases}$ 

#### **Definitions and Parameters:**

- $x \in \mathbb{R}^n$ , *n*-dimensional cancer state;
- $d : \mathbb{R}_{\perp} \to \mathcal{D}$  ( $\mathcal{D}$  compact), time-dependent intensity of the therapy (control);
- $\Delta_{\text{succ}} \subset \mathbb{R}^n$ , success region;
- $\Delta_{\text{fail}} \subset \mathbb{R}^n$ , failure region;
- $\Delta = \Delta_{\text{succ}} \cup \Delta_{\text{fail}}$ , terminal set.

#### **Total treatment time:**

$$T(\boldsymbol{\xi}, \boldsymbol{d}(\cdot)) = \inf \left\{ t \in \mathbb{R}_+ \mid \boldsymbol{x}(t) \in \Delta, \ \boldsymbol{x}(0) = \boldsymbol{\xi} \right\}. \tag{2.2}$$

Treatment cost function:

$$\mathcal{J}(\boldsymbol{\xi}, \boldsymbol{d}(\cdot)) = \int_{0}^{T} K(\boldsymbol{x}(\tau), \boldsymbol{d}(\tau)) d\tau + g\left(\boldsymbol{x}(T)\right), \tag{2.3}$$

where  $T := T(\xi, d(\cdot))$  is the terminal time, K(x, d) is the running cost, and the terminal cost is

$$g(\mathbf{x}) = \left\{ egin{aligned} +\infty, & ext{if } \mathbf{x}(T) \in \Delta_{ ext{fail}}, \ 0, & ext{if } \mathbf{x}(T) \in \Delta_{ ext{succ}}. \end{aligned} 
ight.$$

Value function:

$$u(\xi) = \inf_{d(\cdot)} \mathcal{J}(\xi, d(\cdot))$$

is found by numerically solving a first-order HJB PDE

$$\min_{\boldsymbol{d} \in \mathcal{D}} \left\{ K(\boldsymbol{\xi}, \boldsymbol{d}) + \nabla u(\boldsymbol{\xi}) \cdot \boldsymbol{f}(\boldsymbol{\xi}, \boldsymbol{d}) \right\} = 0, \tag{2.4}$$

with the boundary condition u = g on  $\Delta$ . See S1 Text C.1 for the derivation and D.2 for the numerics.

Box 2: Problem setup of the stochastic optimal control problem Stochastic evolution dynamics with control on therapy intensity (a *drift-diffusion* process):

$$\begin{cases} dX = a(X, d) dt + \Sigma(X, d) dW, \\ X(0) = \xi. \end{cases}$$
 (2.5)

 $\textbf{Process } \textit{terminates} \text{ as soon as either } \begin{cases} \textit{X}(t) \in \Delta_{\text{succ}}, & \text{ if therapy succeeds;} \\ \textit{X}(t) \in \Delta_{\text{fail}}, & \text{ if therapy fails.} \end{cases}$ 

#### **Definitions and Parameters:**

- $X \in \mathbb{R}^n$ , n-dimensional cancer state;
- W, standard m-dimensional Brownian motion;
- $a(X, d) \in \mathbb{R}^n$ , the drift function;
- $\Sigma(X, d) \in \mathbb{R}^{n \times m}$ , the diffusion function.

**Note:** Definitions of the *total treatment time*  $T := T(\xi, d(\cdot))$  and the *treatment cost function*  $\mathcal{J}(\xi, d(\cdot))$  stay the same as in Box 1. But they are now *random variables* as we will replace x(t) by X(t).

deterministic version of each control problem. If K is chosen so that the overall cost of a successful therapy  $\mathcal J$  reflects a weighted sum of the total therapy duration and the cumulative use of each drug, the weights can be adjusted to reflect the relative importance of these optimization criteria. In this case, if f is also a linear function of d, it is easy to show that the optimal treatment policy  $d_{\star}(\xi)$  is generally bang-bang; i.e., for each drug, it prescribes either no usage or the maximal (MTD) usage in every cancer state  $\xi$ .

In a generic continuous-time stochastic cancer model (summarized in Box 2), the tumor state X(t) becomes a random variable, with the dynamics specified by a Stochastic Differential Equation (SDE) (2.5), which replaces the deterministic Ordinary Differential Equation (ODE) (2.1). The definitions of the *total treatment time*  $T(\xi, d(\cdot))$  and the *overall treatment cost*  $\mathcal{J}(\xi, d(\cdot))$  remain the same, but both of them become random variables. The standard (*risk-neutral*) approach of stochastic optimal control is to find a feedback-form treatment policy that minimizes the expected treatment cost  $\mathbb{E}[\mathcal{J}]$ . As explained in Box 3, the resulting value function satisfies another stationary HJB PDE (2.7). The choice of suitable boundary conditions is more subtle here: setting  $g = +\infty$  on  $\Delta_{\text{fail}}$  is no longer an option since the probability of entering  $\Delta_{\text{fail}}$  before  $\Delta_{\text{succ}}$  is usually positive under every treatment policy, which would result in  $\mathcal{J} = +\infty$  for every occasional failure and the overall  $\mathbb{E}[\mathcal{J}] = +\infty$ . This makes it necessary to either choose a specific finite "cost" of failure (which can be problematic both for practical and ethical reasons) or switch to an entirely different optimization objective. For example, one can try to simply maximize the probability of fulfilling the therapy goals (i.e., eventually reaching  $\Delta_{\text{succ}}$  while avoiding  $\Delta_{\text{fail}}$ ) by solving the Eq (2.9). But this latter

# Box 3: Standard stochastic dynamic programming approaches A risk-neutral (expectation-minimizing) approach [41]:

Value function:

$$w(\xi) = \inf_{\boldsymbol{d}(\cdot)} \mathbb{E}[\mathcal{J}(\xi, \boldsymbol{d}(\cdot))]$$
 (2.6)

can be found by solving a second-order Hamilton-Jacobi-Bellman (HJB) equation:

$$\min_{\boldsymbol{d}\in\mathcal{D}}\left\{K(\boldsymbol{\xi},\boldsymbol{d}) + \nabla w(\boldsymbol{\xi}) \cdot \boldsymbol{a}(\boldsymbol{\xi},\boldsymbol{d}) + \frac{1}{2}\sum_{i,j=1}^{n} \frac{\partial^{2}}{\partial \xi_{i} \partial \xi_{j}} w(\boldsymbol{\xi}) \boldsymbol{B}(\boldsymbol{\xi},\boldsymbol{d})_{i,j}\right\} = 0, \tag{2.7}$$

where  $\mathbf{B} = \Sigma \Sigma^{\top}$  and w = g on  $\Delta = \Delta_{\text{succ}} \cup \Delta_{\text{fail}}$ .

**Note:** If one uses  $g(X(T)) = +\infty$  when therapy fails (i.e., when  $X(T) \in \Delta_{\text{fail}}$ ), the diffusion will generally result in  $w = +\infty$  for most if not all initial tumor configurations outside of  $\Delta_{\text{succ}}$ .

An alternative is to maximize the probability of eventual goal attainment:

Value function:

$$\tilde{w}(\xi) = \sup_{d(\cdot)} \mathbb{P}(X(T) \in \Delta_{\text{succ}})$$
 (2.8)

can be found by solving a second-order Hamilton-Jacobi-Bellman (HJB) equation:

$$\max_{\boldsymbol{d}\in\mathcal{D}}\left\{\nabla \tilde{w}(\boldsymbol{\xi})\cdot\boldsymbol{a}(\boldsymbol{\xi},\boldsymbol{d}) + \frac{1}{2}\sum_{i,j=1}^{n}\frac{\partial^{2}}{\partial\xi_{i}\partial\xi_{j}}\tilde{w}(\boldsymbol{\xi})\boldsymbol{B}(\boldsymbol{\xi},\boldsymbol{d})_{i,j}\right\} = 0,$$
(2.9)

with the boundary condition  $\tilde{w}=1$  on  $\Delta_{\rm succ}$  and  $\tilde{w}=0$  on  $\Delta_{\rm fail}$ .

formulation ignores many important practical considerations: e.g., it can easily result in an unreasonably long treatment time or in significant side effects from a prolonged MTD-level drug administration.

In contrast, the approach we are pursuing here allows for a more nuanced definition of success (e.g., taking into account the total drug usage, the treatment duration, and/or the cumulative burden from the tumor). Choosing a running cost K to reflect the above factors, we define

the overall cost as 
$$\mathcal{J}(\boldsymbol{\xi}, \boldsymbol{d}(\cdot)) = \int_{0}^{T} K(\boldsymbol{X}(\tau), \boldsymbol{d}(\tau)) d\tau + g(\boldsymbol{X}(T))$$
, which might be infinite if  $\boldsymbol{X}$ 

 $(T) \in \Delta_{\mathrm{fail}}$ . We then maximize the probability of reaching the policy goals, but constraining the overall cost by some pre-specified threshold  $\bar{s}$ . I.e., we need to find an adaptive therapy that maximizes  $\mathbb{P}(\mathcal{J} \leq \bar{s})$ . Our goal is to compute such *threshold-aware* policies efficiently for all starting tumor configurations  $\xi$  and a broad range of threshold levels simultaneously. It is easy to see that here good treatment policies will have to also take into account the cost accumulated so far. This makes it natural to treat our chosen threshold  $\bar{s}$  as an *initial cost budget*, tracking the remaining budget s(t) by solving Eq. (2.13) in Box 4. The value function can be found

# Box 4: Our threshold-aware approach

Value function:

$$\nu(\xi, \bar{s}) = \sup_{d(\cdot)} \mathbb{P}\left(\mathcal{J}\left(\xi, d(\cdot)\right) \le \bar{s}\right)$$
 (2.10)

can be found by solving a different second-order HJB equation:

$$\max_{\boldsymbol{d} \in \mathcal{D}} \left\{ -\frac{\partial}{\partial s} \nu(\boldsymbol{\xi}, s) K(\boldsymbol{\xi}, \boldsymbol{d}) + \nabla \nu(\boldsymbol{\xi}, s) \cdot \boldsymbol{a}(\boldsymbol{\xi}, \boldsymbol{d}) + \frac{1}{2} \sum_{i,j=1}^{n} \frac{\partial^{2}}{\partial \xi_{i} \partial \xi_{j}} \nu(\boldsymbol{\xi}, s) \boldsymbol{B}(\boldsymbol{\xi}, \boldsymbol{d})_{i,j} \right\} = 0, (2.11)$$

where  $s \in [0, \bar{S}]$ . See the detailed derivation in S1 Text §C.2.

The boundary conditions of HJB equation:

$$\begin{cases} \nu(\xi, s) = 1, & \text{if } \xi \in \Delta_{\text{succ}}, s \in [0, \bar{\mathcal{S}}]; \\ \nu(\xi, s) = 0, & \text{if } \xi \in \Delta_{\text{fail}}, s \in [0, \bar{\mathcal{S}}]; \\ \nu(\xi, 0) = 0, & \text{if } \xi \notin \Delta_{\text{succ}}. \end{cases}$$

$$(2.12)$$

The (random) ODE describing the reduction of budget:

$$\dot{s} = -K(\mathbf{X}(t), \mathbf{d}(t)), \qquad s(0) = \bar{s} \in (0, \bar{\mathcal{S}}]. \tag{2.13}$$

by solving the parabolic PDE (2.11) numerically, and the optimal feedback policy  $d_*^{\bar{s}}(\xi, s)$  is recovered in the process for all  $\bar{s} \in (0, \bar{S}]$ .

We note that, in classical stochastic optimal control, parabolic HJB equations are usually encountered when dealing with finite horizon problems, where the terminal time T is specified in advance. See [8, 42, 43] for typical examples in cancer-related literature. In contrast, the parabolicity in PDE (2.11) arises because of the monotone decrease in the remaining budget s (t).

The details of our numerical method based on Box 4 are included in S1 Text \$D.1. In the interest of computational reproducibility, we provide the source code for approximating value functions and computing threshold-aware policies for all the examples from \$3 at <a href="https://github.com/eikonal-equation/Stochastic-Cancer">https://github.com/eikonal-equation/Stochastic-Cancer</a>.

While this threshold-aware framework has important advantages illustrated below, it also brings to the forefront several subtleties avoided in the more traditional stochastic optimal control approaches. First, an adaptive treatment policy optimal for one specific threshold is usually not optimal for another. (The starting budget in (2.13) is important for deciding when to administer drugs.) This would make it necessary for a practitioner to have a detailed discussion with their patient to choose a suitable threshold value before the treatment is started. Second, stochastic perturbations make the outcome random, and the budget might run out under any treatment policy; i.e., we might see  $s(t_{\sharp})=0$  at some random time  $t_{\sharp}$ . But this scenario is only a failure in the sense that the overall cost  ${\mathcal J}$  will now definitely exceed the threshold value

 $\bar{s}$ . If the patient is still alive  $(X(\tau) \notin \Delta_{\text{fail}}, \ \forall \tau \in [0, t_{\sharp}])$  and interested in continuing treatment, one has to make a decision on the new strategy. This can be done either by posing a new threshold for future treatment costs or by switching to an entirely different policy—e.g., either by employing some traditional stochastic optimal control approach (based on Eqs (2.7) or (2.9)) or by using a deterministic-optimal policy based on Eq (2.4). The latter version is used in all stochastic simulations in the following sections.

Informally, threshold-aware policies reflect a tension between two objectives which are often (but not always) in conflict. Maximizing the probability of treatment attaining its primary goals (e.g., tumor stabilization or eradication) is balanced against reducing the cost (a combination of tumor and treatment burdens) suffered along the way. The former is optimized but only over the scenarios where the latter stays below the prescribed threshold. We close this subsection by highlighting connections of our approach to general multiobjective optimal control and optimal control with integral constraints.

In deterministic optimal control theory, the idea of treating some version of cumulative cost as an additional state variable is well-known. But the resulting ODE systems are typically treated within the framework of *Pontryagin's Maximum Principle* (PMP) [44], which has an important advantage (its suitability for high-dimensional problems) but also a number of serious drawbacks: the fact that policies are not recovered in feedback form, the fact that these policies are generally not guaranteed to be *globally* optimal, occasional difficulties in ensuring the convergence of numerical methods needed to find such policies, and challenges in handling non-trivial state constraints. In cancer literature, this PMP-based approach has been used to impose "isoperimetric constraints" on the amount of administered chemotherapy [45] or immunotherapy [46]. In addition to the issues listed above, we note that the suitability of equality (isoperimetric) constraints is not obvious in many cancer applications. Indeed, the fact that a less aggressive treatment may in some cases improve the outcomes is one of the main reasons for the interest in adaptive therapies. Thus, insisting that all available drugs must be used is hard to justify, and inequality constraints (e.g., imposing an upper bound on the cumulative drug use) seem much more reasonable.

The first dynamic programming (HJB-based) formulation for handling such constraints in general deterministic control problems was developed in [47]. It circumvents all these PMP-associated difficulties with an added benefit of finding globally optimal policies for a range of inequality constraint levels simultaneously. The threshold-aware method presented here extends many of the same ideas to a stochastic setting.

#### 2.2 Example 1: An EGT-based competition model

To develop our first example, we adopt the base model of cancer evolution proposed by Kaznatcheev et al. in [48], which describes a competition of 3 types of cancer cells. Glycolytic cells (GLY) are anaerobic and produce lactic acid, which damages the surrounding non-cancerous tissue. The other two types are aerobic and benefit from better vasculature, development of which is promoted by production of the VEGF signaling protein. Thus, the VEGF (over)-producing cells (VOP) devote some of their resources to vasculature development, while the remaining aerobic cells are essentially free-riders or *defectors* (DEF) in game-theoretic terminology. If  $(z_G(t), z_D(t), z_V(t))$  encode the time-dependent subpopulation sizes of these three cancer types, their dynamics are given by  $\dot{z}_i = \psi_i z_i$ , where  $i \in \{G, D, V\}$  and  $(\psi_G, \psi_D, \psi_V)$  are the respective type fitnesses. The actual expressions for these  $\psi_i$  are derived from the inter-population competition in the usual EGT framework; see S1 Text §A.1. This competition of cells in the tumor is modeled as a "public goods" / "club goods" game: VEGF is a "club good" since it benefits only VOP and DEF cells, while the acid generated by GLY is a "public good" since

the damage to healthy tissue is assumed to benefit all cancer cells. The base model in [48] assumes that each cell interacts with n others nearby. How much it benefits from these interactions depends on its own type and the proportions of different cell types among those nearby cells. Assuming that all participants are drawn uniformly at random from a large well-mixed population, one can derive all fitnesses  $\psi_i$  as expected payoffs in this game of (n+1) players. Those expected payoffs will naturally depend on the current subpopulation fractions (or relative abundances)  $x_G = \frac{z_G}{z_G + z_D + z_V}$ ,  $x_D = \frac{z_D}{z_G + z_D + z_V}$ , and  $x_V = \frac{z_V}{z_G + z_D + z_V}$ . A Replicator Ordinary Differential Equation (ODE) [49, 50] is a standard EGT model for predicting the changes in these subpopulation-fractions as a function of time.

In both the original deterministic case and its stochastic extension, it is easier to view the

replicator equation as a 2-dimensional system (e.g., by noting that  $x_D = 1 - x_G - x_V$ ). Following [48], we use a slightly different reduction, rewriting everything in terms of the proportion of glycolytic cells in the tumor  $p(t) = x_G(t)$  and the proportion of VOP among aerobic cells  $q(t) = \frac{x_{\rm V}(t)}{x_{\rm V}(t) + x_{\rm D}(t)}$ . A drug therapy (in this example, affecting the fitness of GLY cells only) is similarly easy to encode by modifying the Replicator ODE; see Eq (2.14) in Box 5 and the Supplementary Materials in [48] for the derivation. The goal of the drug therapy here is to drive the GLY fraction  $p(t) = x_G(t)$  down below a specified "stabilization barrier"  $\gamma_t$ . (In [48], this goal is justified by noting that, with GLY gone, the DEF cells will then quickly overcome VOP, leading to "an aerobic tumor with no—or significantly diminished—ability to recruit blood vessels," which stabilizes (or at least significantly slows down the growth of) the tumor.) For a range of parameter values, this model yields periodic behavior of cancer subpopulations: without drugs,  $x_G(t)$ ,  $x_D(t)$ , and  $x_V(t)$  alternate in being dominant in the tumor, with the amplitude of oscillations determined by the initial conditions [48]. This highlights the importance of proper timing in therapies: starting from the same initial tumor composition  $(q_0, p_0)$ , the same MTD therapy of a fixed duration could lead to either a stabilization (p(t) falling below  $\gamma_r$ ) or a death (p(t) rising above the specified "failure barrier"  $\gamma_f$ ) depending on how long we wait until this therapy starts; see Fig 2 in Kaznatcheev et al. [48].

This strongly suggests the advantage of *adaptive therapies*, which prescribe the amount of drugs based on continuous or occasional monitoring of (q(t), p(t)) or some proxy (non-invasively measured) variables. A natural question is how to optimize such policies to reduce the total amount of drugs used and the total duration of treatment until  $p(t) < \gamma_r$ . Gluzman et al. have addressed this in [29] using the framework of deterministic optimal control theory [39]. A time-dependent intensity of the therapy d(t) (ranging from 0 to the MTD level  $d_{max}$ ) was

chosen to minimize the overall cost of treatment  $\mathcal{J}(q_0, p_0, d(\cdot)) = \int\limits_0^T d(t) dt + \delta T + \delta T$ 

g(q(T),p(T)), where T is the time till stabilization (or failure, if  $(q(T),p(T))\in\Delta_{\mathrm{fail}}$  and  $g=+\infty$ ) while the value of  $\delta>0$  reflects the relative importance of two optimization goals (total drugs vs total time). In the framework of deterministic dynamic programming [40] summarized in Box 1, this corresponds to minimizing the integral of the running cost  $K=d(t)+\delta$ . In [29], the deterministic-optimal policy is obtained in *feedback form* (i.e.,  $d=d_{\star}(q,p)$ ) by numerically solving the Hamilton-Jacobi-Bellman (HJB) PDE (2.4). As explained in §2.1, this policy is bang-bang. Fig 1a summarizes it (showing in yellow the MTD region where  $d_{\star}(q,p)=d_{\max}$ ) and illustrates the corresponding trajectory for one specific initial  $(q_0,p_0)$ .

A natural way to introduce stochastic perturbations into this base model is to assume that the rates of subpopulation growth/decay are actually random and normally distributed at any instant, with the fitness functions ( $\psi_G$ ,  $\psi_D$ ,  $\psi_V$ ) encoding the expected values of those rates and

Box 5: Example 1 (an EGT-based competition model adopted from [29, 48]) The deterministic base model (components for the approach in Box 1):

$$\mathbf{x} = \begin{bmatrix} q \\ p \end{bmatrix} := \begin{bmatrix} \frac{x_{\text{V}}}{x_{\text{V}} + x_{\text{D}}} \end{bmatrix} \text{ and } \dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}, d) = \begin{bmatrix} q(t) \left( 1 - q(t) \right) \left( \frac{b_{\text{V}}}{n+1} \sum_{k=0}^{n} (p(t))^{k} - c \right) \\ p(t) \left( 1 - p(t) \right) \left( \frac{b_{\text{d}}}{n+1} - (b_{\text{V}} - c)q(t) - d(t) \right) \end{bmatrix} (2.14)$$

The above reflects the formulas for subpopulation fitnesses ( $\psi_G$ ,  $\psi_D$ ,  $\psi_V$ ); see details in S1 Text §A.1.

The stochastic model (components for the approaches in Boxes 2-4):

$$\begin{split} \boldsymbol{X} &= \begin{bmatrix} Q \\ P \end{bmatrix} \coloneqq \begin{bmatrix} \frac{X_{\text{V}}}{X_{\text{V}} + X_{\text{D}}} \end{bmatrix}, \quad \boldsymbol{W} = \begin{bmatrix} W_{\text{G}} \\ W_{\text{D}} \\ W_{\text{V}} \end{bmatrix}, \\ \boldsymbol{a}(\boldsymbol{X}, d) &= \begin{bmatrix} Q(1 - Q) \left\{ \left( \frac{b_{\text{v}}}{n+1} \left[ \sum_{k=0}^{n} P^{k} \right] - c \right) + \left[ (1 - Q)\sigma_{2}^{2} - Q\sigma_{3}^{2} \right] \right\} \\ P(1 - P) \left\{ \left( \frac{b_{a}}{n+1} - (b_{v} - c)Q - d \right) - \left[ \sigma_{1}^{2}P - \sigma_{2}^{2}(1 - P)(1 - Q)^{2} - \sigma_{3}^{2}(1 - P)Q^{2} \right] \right\} \end{bmatrix}, \end{split}$$

$$\boldsymbol{\Sigma}(\boldsymbol{X}, d) &= \begin{bmatrix} 0 & -\sigma_{\text{D}}Q(1 - Q) & \sigma_{\text{V}}Q(1 - Q) \\ \sigma_{\text{G}}P(1 - P) & \sigma_{\text{D}}P(1 - P)(1 - Q) & \sigma_{\text{V}}P(1 - P)Q \end{bmatrix}. \end{split}$$

#### **Definitions and Parameters:**

- $d: \mathbb{R}_+ \to [0, d_{\text{max}}]$ , time-dependent intensity of GLY-targeting therapy;
- $\Delta_{\text{succ}} = \{(q, p) \in [0, 1]^2 \mid p < \gamma_r\}$ , success region where  $\gamma_r$  is the *stabilization barrier*;
- $\Delta_{\text{fail}} = \{(q, p) \in [0, 1]^2 \mid p > \gamma_f\}$ , failure region where  $\gamma_f$  is the *failure barrier*;
- $K(X, d) = d + \delta$ , running cost function where  $\delta$  is the treatment time penalty;
- $g = \begin{cases} +\infty, & \text{if } \mathbf{x}(T) \in \Delta_{\text{fail}}, \\ 0, & \text{if } \mathbf{x}(T) \in \Delta_{\text{succ}}, \end{cases}$  terminal cost;
- $W = (W_G, W_D, W_V)$ , standard 3D Brownian motion for (GLY, DEF, VOP) cells;
- $(\sigma_G, \sigma_D, \sigma_V)$ , volatilities for (GLY, DEF, VOP) cells;
- $b_a$ , the benefit per unit of acidification;
- $b_{\nu}$ , the benefit from the oxygen per unit of vascularization;
- *c*, the cost of production of VEGF;
- (n + 1), the number of cells in the interaction group.

Conditions for the heterogeneous regime (coexistence of all cell types):

$$\frac{b_a}{n+1} < b_v - c < cn. {(2.16)}$$

The optimal threshold-aware policy in feedback form:

$$d_{*}(q, p, s) = \begin{cases} d_{\text{max}}, & \text{if } \left(\frac{\partial \nu}{\partial p} p(1 - p) + \frac{\partial \nu}{\partial s}\right) < 0, \\ 0, & \text{otherwise.} \end{cases}$$
 (2.17)

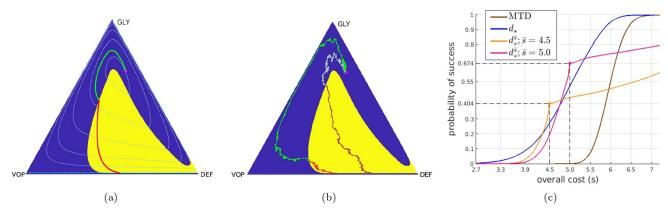


Fig 1. Deterministic-optimal policy in the EGT-model. The (GLY-VOP-DEF) triangle represents all possible relative abundances of respective subpopulations. Since the optimal policy is bang-bang, we show it by using the yellow background where drugs should be used at the MTD rate and the blue background where no drugs should be used at all. Starting from an initial state ( $q_0$ ,  $p_0$ ) = (0.26, 0.665) (magenta dot), the subfigures show (a) the optimal trajectory found from the truly deterministically driven system (2.14) with cost 5.13; (b) two representative sample paths generated under the deterministic-optimal policy but subject to stochastic fitness perturbations (the brighter one incurs a total cost of 3.33, whereas the duller-colored path incurs a much higher 6.23); (c) CDFs of the cumulative cost  $\mathcal{J}$  approximated using  $10^5$  random simulations. In both (a) & (b), the green parts of trajectories correspond to not prescribing drugs and the red parts of trajectories correspond to prescribing drugs and the red parts of trajectories correspond to prescribing drugs and the red parts of trajectories correspond to prescribing drugs and the red parts of trajectories correspond to not prescribing drugs and the red parts of trajectories correspond to prescribing drugs and the magenta of the value function in the deterministic case are shown in light blue. In (c), the blue curve is the CDF generated with the deterministic-optimal policy  $d_{\star}$ . Its observed median and mean conditioning on success are 4.95 and 4.91 respectively. The brown curve is the CDF generated with the MTD-based therapy, which in this example also maximizes the chances of "budget-unconstrained" tumor stabilization. Its observed median and mean conditioning on success are 5.95 and 5.96 respectively. Orange and pink curves show the CDFs for two different threshold-aware policies (with  $\bar{s} = 4.5$  and  $\bar{s} = 5$  respectively). The large dot on each of them represents the maximized probability of not exceeding the corresponding threshold. The t

the scale of random perturbations specified by  $(\sigma_G, \sigma_D, \sigma_V)$ . This approach, originating from Fudenberg and Harris paper [51], is suitable for modeling heterogeneous tumors, in which subpopulations not only interact [52] but can also vary in their growth rates over time [53]. Adopting the usual probabilistic notation of using capital letters for random variables, we can again start with the subpopulation sizes  $(Z_G, Z_D, Z_V)$  evolving based on the *Stochastic Differential Equations* (SDEs)  $dZ_i = (\psi_i dt + \sigma_i dW_i)Z_i$ , where  $i \in \{G, D, V\}$  and each  $W_i$  is a standard one-dimensional Brownian motion, modeling independent perturbations to the fitness of the respective subpopulation. This can be used to derive the SDEs for the corresponding fractions  $(X_G, X_D, X_V)$  We note that similar Stochastic Replicator Equations arise naturally in ecology, where they have been studied in depth to address a possible coexistence of species in randomly perturbed environments [54, 55].

The summary of Replicator SDEs for the reduced (Q, P) coordinates is provided in Box 5; the derivation can be found in S1 Text §B.2. The terminal set  $\Delta$  is still the same: the process terminates as soon as P(t) crosses a stabilization barrier (GLY's are low, leaving mostly aerobic cells in the tumor) or the failure barrier (GLY's are high, the patient dies). But the terminal time T and the incurred cumulative cost  $\mathcal J$  will also be random even if we fix the initial tumor configuration  $(q_0, p_0)$  and choose a specific treatment policy  $d(\cdot)$ . Fig 1b shows one example of using the deterministic-optimal policy  $d(t) = d_{\star} (Q(t), P(t))$  in this stochastic setting. Gathering statistics from many random simulations that start from the same  $(q_0, p_0)$ , we can approximate the *Cumulative Distribution Function* (CDF), measuring the probability of keeping  $\mathcal J$  below any given threshold s if the deterministic-optimal policy is employed:

$$F_{d_{\bullet}}(s) = \mathbb{P}(\mathcal{J} \leq s),$$

whose graph is shown in blue in Fig 1c. If one instead opts to solve the PDE (2.9) to maximize

the probability of reaching  $\Delta_{\text{succ}}$  while avoiding  $\Delta_{\text{fail}}$ , this yields a simple MTD-policy  $d = d_{\text{max}}$  whose CDF (shown in brown in Fig 1c) is strictly worse than that of  $d_{\star}$ . This is not surprising since the more selective  $d_{\star}$  is quite safe for this particular  $(q_0, p_0)$ , with  $\Delta_{\text{fail}}$  avoided in all of our  $10^5$  simulations. However, its resulting "cost" can be still high in many scenarios. E.g., in 47.4% of the  $d_{\star}$ -based simulations,  $\mathcal{J}$  exceeded 5; in 72.6% of all cases it exceeded 4.5.

This motivates our optimization approach: deriving a *threshold-aware optimal policy*  $d_*^{\bar{s}}$  to maximize the probability of stabilization without exceeding a specific cost threshold  $\bar{s}$ . As explained in §2.1 and summarized in Box 4, this is accomplished for a range of threshold values and all initial cancer configurations simultaneously. Fig 1c already shows that such policies can provide significant threshold-specific advantages over the deterministic-optimal therapy. Additional simulation results and the actual policies are illustrated in §3.1.

# 2.3 Example 2: A Sensitive-Resistant competition model

We also illustrate our approach by extending a model proposed by Carrère [3], which focuses on the actual size of lung cancer cell populations studied *in vitro*. They consider a heterogeneous tumor that consists of two types of lung cancer cells: the sensitive (*S*) "A549" (sensitive to the drug "Epothilene") and the resistant (*R*) "A549 Epo40". This was based on the data from a series of experiments conducted by Manon Carrè at the Center for Research in Oncobiology and Oncopharmacology, Aix-Marseille Université. Mutation events were neglected due to their rarity at the considered dosages of Epothilene and due to relatively short treatment durations. The competition model presented below was derived based on phenotypical observations, with fluorescent marking used to trace and differentiate the cells.

Considered separately, both of these types obey a logistic growth model with respective intrinsic growth rates  $g_S$  and  $g_R$ . The carrying capacity of the Petri dish (C) is assumed to be shared, with the resistant cells assumed to be m times bigger than the sensitive; so, the fraction of space used at the time t is  $\frac{z_S(t)+mz_R(t)}{C}$ . When cultivated together, it was observed that the sensitive cells quickly outgrow the resistant ones despite the fact that their intrinsic growth rates are similar [3]. To model this competitive advantage, they have used an additional competition term  $-\beta z_S z_R$  to describe the rate of change of  $z_R(t)$  with the coefficient  $\beta$  calibrated based on experimental data. It was further assumed that R cells are completely resistant to a specific drug, which reduces the population of S cells at the rate of  $\alpha z_S(t)d(t)$  with d(t) reflecting the current rate of drug delivery and the constant coefficient  $\alpha$  reflecting that drug's effectiveness. With a normalization  $z_S(t) \to z_S(t)/C$ ,  $z_R(t) \to z_R(t)/C$ , the resulting dynamics are summarized by

$$\dot{z}_{S}(t) = g_{S}(1 - z_{S}(t) - mz_{R}(t))z_{S}(t) - \alpha z_{S}(t)d(t), 
\dot{z}_{R}(t) = g_{R}(1 - z_{S}(t) - mz_{R}(t))z_{R}(t) - \beta Cz_{S}(t)z_{R}(t).$$
(2.18)

In both the original deterministic case and its stochastic extension, it is more convenient to restate the dynamics in terms of the *effective tumor size*  $p(t) = z_S(t) + mz_R(t)$  and the fraction of effective tumor size comprised of the sensitive cells  $q(t) = z_S(t)/p(t)$ . Note that the proportion of sensitive cells in the tumor, by number, is  $\frac{mq(t)}{1+(m-1)q(t)}$  instead of just q(t) due to the size ratio m between S and R cells. This change of coordinates yields an ODE model (2.19) summarized in Box 6; see S1 Text §A.2 for the derivation. In this case, the goal of our adaptive therapy is eradication: i.e., driving the total tumor size p(t) below some remission barrier  $\gamma_r$  (e.g., a physical detection level) while ensuring that throughout the treatment this p(t) stays below a significantly higher failure barrier  $\gamma_f < 1$ .

Fig 2a illustrates the natural dynamics of this model with no drug use. In this case, the competitive pressure reduces the population R, which at first decreases the tumor size for many initial conditions. But a rapid growth in S eventually increases the overall tumor, leading to an inevitable failure  $(p(t) > \gamma_f)$ . The deterministic-optimal drug therapy is again sought to minimize a weighted sum of total drugs used and the time of treatment (with the running cost  $K = d(t) + \delta$ ) until the eradication. It is obtained in feedback form  $d = d_{\star}(q, p)$  after solving the PDE (2.4). Fig 2b shows that, for smaller tumor sizes, this  $d_{\star}$  prescribes MTD-level treatment only after this initial tumor reduction is over, once S gets rid of most R cells which are not sensitive to Epothilene. However, for larger initial p, this deterministic-optimal policy starts using the drugs much earlier, planning to keep S cells in check as soon as they are numerous enough to control R.

Stochastic perturbations can be similarly introduced here by assuming that the intrinsic growth rates are actually random and normally distributed at any instant. (This approach was also used in modeling persistence strategies among bacteria in [42].) In Example 1, we assumed that the fitness function of each subpopulation was affected by its own random perturbations. The Brownian motion in Box 5 was three-dimensional, corresponding to subpopulations impacted by three separate and uncorrelated aspects of the fluctuating environment. Depending on the nature of perturbations, a similar assumption might be reasonable in the current example as well. But this is not a necessary feature for the threshold-aware optimization approach to be applicable. To demonstrate this, we will instead assume here that the same aspect of fluctuating environment impacts both subpopulations, and thus a single (1D) Brownian motion perturbs the intrinsic growth rates of both S and R. We will use  $(g_S, g_R)$  to represent their expected growth rates and  $(\sigma_S, \sigma_R)$  to denote their respective volatilities. This yields SDEs for the stochastic evolution of (Q, P), which are derived in S1 Text  $\S B.3$  and summarized in Box 6. As shown in Fig 2c, if the deterministic-optimal policy  $d_{\star}$  is used in this stochastic setting, the initiation time of the MTD-based therapy (and the resulting overall cost  $\mathcal{J}$ ) can vary significantly. This motivates us again to use the threshold-aware approach based on the PDE (2.11), with the policies illustrated and advantages quantified in §3.2.

#### 3 Results

# 3.1 Policies, trajectories, and CDFs for the EGT-based model

We explore the structure and performance of threshold-aware policies computed for the system described in §2.2. The parameter values  $d_{\rm max}=3$ ,  $b_a=2.5$ ,  $b_v=2$ , c=1, n=4 are the same ones provided in Kaznatcheev et al. [48] and Gluzman et al. [29]. However, we use  $\gamma_{\rm r}=1-\gamma_{\rm f}=10^{-2}$  and  $\delta=0.05$  as opposed to  $\gamma_{\rm r}=1-\gamma_{\rm f}=10^{-1.5}$  and  $\delta=0.01$  in [29]. Additionally, we consider small uniform constant volatilities  $\sigma_{\rm G}=\sigma_{\rm D}=\sigma_{\rm V}=0.15$ , characterizing the scale of random perturbations in fitness function for all 3 cancer subpopulations. The details of our Monte-Carlo simulations used to build all CDFs can be found in \$1 Text \$D.3. Additional examples, including those with higher volatilities, in which the threshold-performance advantages are even more significant, can be found in \$1 Text \$E.

In Fig 3, we present some representative *s*-slices of threshold-aware optimal policies and their corresponding optimal probability of success for respective threshold values. Since these policies are also bang-bang, the drugs-on region (at the MTD level) is shown in yellow and the drugs-off region is shown in blue in all of our figures, following the convention from [29]. We observe that this drugs-on region is strongly *s*-dependent and completely different from the one in the deterministic-optimal case shown in Fig 1a. Since the cancer evolution considered here has stochastic dynamics given in (2.5) and (2.15), different realizations of random

Box 6: Example 2 (a Sensitive-Resistant competition model adopted from [3]) The deterministic base model (components for the approach in Box 1):

$$\mathbf{x} = \begin{bmatrix} q \\ p \end{bmatrix} \coloneqq \begin{bmatrix} \frac{z_{S}}{z_{S} + mz_{R}} \\ z_{S} + mz_{R} \end{bmatrix}$$

and

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}, d) = \begin{bmatrix} (1-p)q(1-q)(g_{S} - g_{R}) + \beta Cp^{2}q^{2}(1-q) - \alpha q(1-q)d(t) \\ p(1-p)[g_{S}q + g_{R}(1-q)] - \beta Cp^{2}q(1-q) - \alpha qpd(t) \end{bmatrix}$$
(2.19)

The stochastic model (components for the approaches in Boxes 2-4):

$$\begin{split} \boldsymbol{X} &= \begin{bmatrix} Q \\ P \end{bmatrix} \coloneqq \begin{bmatrix} \frac{Z_{\text{S}}}{Z_{\text{S}} + mZ_{\text{R}}} \\ Z_{\text{S}} + mZ_{\text{R}} \end{bmatrix}, \quad \boldsymbol{W} = \begin{bmatrix} B_{\text{t}} \end{bmatrix}, \\ \boldsymbol{a}(\boldsymbol{X}, d) &= \begin{bmatrix} Q(1-Q) \Big\{ (1-P)(g_{\text{S}} - g_{\text{R}}) - \alpha d + \beta CQP + (1-P)^2 [\sigma_{\text{R}}^2 (1-Q) - \sigma_{\text{S}}^2 Q + \sigma_{\text{S}} \sigma_{\text{R}}] \Big\} \\ P(1-P)(g_{\text{S}}Q + g_{\text{R}}(1-Q)) - \alpha QPd - \beta CP^2 Q(1-Q) \\ \boldsymbol{\Sigma}(\boldsymbol{X}, d) &= \begin{bmatrix} (1-P)Q(1-Q)(\sigma_{\text{S}} - \sigma_{\text{R}}) \\ P(1-P)[\sigma_{\text{S}}Q + \sigma_{\text{R}}(1-Q)] \end{bmatrix}. \end{split}$$

#### **Definitions and Parameters:**

- $d: \mathbb{R}_+ \to [0, d_{\text{max}}]$ , time-dependent intensity of S-targeting therapy;
- $\Delta_{\text{succ}} = \{(q, p) \in [0, 1]^2 \mid p < \gamma_r\}$ , success region where  $\gamma_r$  is the *remission barrier*;
- $\Delta_{\text{fail}} = \{(q, p) \in [0, 1]^2 \mid p > \gamma_f\}$ , failure region where  $\gamma_f$  is the *failure barrier*;
- $K(X, d) = d + \delta$ , running cost function where  $\delta$  is the treatment time penalty;

• 
$$g = \begin{cases} +\infty, & \text{if } \mathbf{x}(T) \in \Delta_{\text{fail}}, \\ 0, & \text{if } \mathbf{x}(T) \in \Delta_{\text{succ}}, \end{cases}$$
 terminal cost;

- $(g_S, g_R)$ , growth rate for the sensitive and resistant cells, respectively;
- *B<sub>t</sub>*, standard 1D Brownian motion;
- $(\sigma_S, \sigma_R)$ , volatilities for the sensitive and resistant cells, respectively;
- *m*, size ratio between *S* and *R* cells;
- *C*, Petri dish carrying capacity;
- $\alpha$ , drug efficiency;
- $\beta$ , action of sensitive on resistant.

Parameter values are specified in S1 Text §E.2.

The optimal threshold-aware policy in feedback form:

$$d_{*}(q, p, s) = \begin{cases} d_{\text{max}}, & \text{if } \left(\frac{\partial v}{\partial q} \alpha q (1 - q) + \frac{\partial v}{\partial p} \alpha q p + \frac{\partial v}{\partial s}\right) < 0, \\ 0, & \text{otherwise.} \end{cases}$$
(2.21)

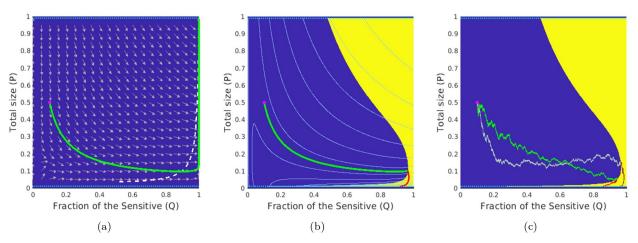


Fig 2. Deterministic-optimal policy in the Sensitive-Resistant model. Starting from an initial state ( $q_0$ ,  $p_0$ ) = (0.1, 0.5) (magenta dot), the subfigures show (a) the deterministic trajectory without therapy that ends in the  $\Delta_{\text{fail}}$ ; (b) the optimal trajectory found from the deterministically driven system (2.1) and (2.19) with cost 49.30; (c) two representative sample paths generated under the deterministic-optimal policy but subject to stochastic perturbations in ( $g_5$ ,  $g_R$ ) (the brighter one incurs a total cost of 49.43, versus a much higher 70.45 for the duller-colored path); In (a), the white dashed-line is part of the nullcline where  $\dot{p}=0$ ; In both (b)&(c), the green parts of trajectories correspond to not prescribing drugs and the red parts of trajectories correspond to prescribing drugs at the MTD rate. The level sets of the value function u in the deterministic case are shown in light blue.

perturbations will result in entirely different sample paths even if the starting configuration and the feedback policy remain the same. Three such representative sample paths are shown in Fig 4, starting from the same initial tumor configuration  $(q_0, p_0) = (0.26, 0.665)$  already used in Fig 1 and focusing on a threshold  $\bar{s} = 5$ . We use the example from Fig 4a, in which the stabilization is achieved while incurring the total cost of  $\mathcal{J} = 4.70 < \bar{s}$ , to illustrate the general use of threshold-aware policies. Starting from the initial budget  $s = \bar{s}$ , the optimal decision on whether to use drugs right away is based on the first diagram in Fig 3a. For our initial tumor state, this indicates that  $d_*(q_0, p_0, 5) = 0$  (not prescribing drugs initially) would maximize the probability of stabilizing the tumor without exceeding the threshold  $\bar{s} = 5.0$ . As time passes, we accumulate the cost, thus decreasing the budget, even if the drugs are not used. If we stay in the blue region for the time  $\theta = 1/\delta$ , the second diagram (the "s = 4.0" case) in Fig 3a becomes relevant, with subsequent budget decreases shifting us to lower and lower s slices. Of course, in reality we constantly reevaluate the decision on  $d_*$  (as s changes continuously while Fig 3a presents just a few representative slices) taking into account the changing tumor configuration (Q(t), P(t)). (Movies with additional information for Figs 3 and 4 are available at https:// eikonal-equation.github.io/Stochastic-Cancer/examples.html).

In contrast to the success story in Fig 4a, we note that there are two very different ways of "failing". First, the process can stop if the proportion of GLY cells becomes too high, as in Fig 4b. When VOP is relatively low, the deterministic portion of the dynamics can bring us close to the failure barrier, with random perturbations resulting in a noticeable probability of crossing into  $\Delta_{\text{fail}}$ . Second, even if we stay away from  $\Delta_{\text{fail}}$ , the budget might be exhausted before reaching  $\Delta_{\text{succ}}$ , as in Fig 4c. Threshold-aware policies provide no guidance once s = 0, but it is reasonable to continue (using some different treatment policy) since the patient is still alive. In our numerical simulations, we switch in this case to a deterministic-optimal policy  $d_{\star}$  illustrated in Fig 1. This decision is somewhat arbitrary; e.g., one could choose instead to switch to an MTD-based policy, which in this example maximizes the probability of reaching  $\Delta_{\text{succ}}$  while avoiding  $\Delta_{\text{fail}}$  without any regard to additional cost incurred thereafter. For this initial tumor

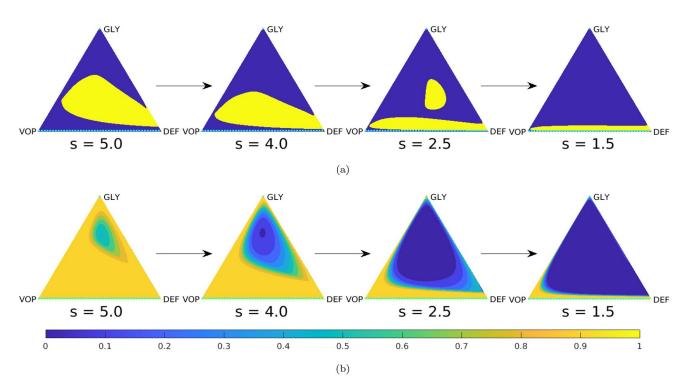


Fig 3. Representative slices of the threshold-aware optimal policy (top row) and the corresponding probability of success (bottom row) for the EGT-based model. Each triangle represents all possible tumor compositions (proportions of GLY/VOP/DEF cells in the population). Top row shows the policy, which prescribes the optimal instantaneous decisions on drug usage given the indicated remaining budget (s) and the current tumor state. Bottom row shows the probability of "stabilization within the budget" if the optimal policy is followed from this time point and onward. Each column corresponds to a specific budget level s, which is shown below each triangle. The arrows indicate the natural decrease of the remaining budget while implementing the policy.

configuration and parameter values, continuing with  $d_{\star}$  typically yields smaller costs while only slightly increasing the chances of eventual failure (e.g.,  $\approx 0.36\%$  of crossings into  $\Delta_{\rm fail}$  using  $d_{\star}$  versus  $\approx 0.07\%$  using the full MTD once the original budget of  $\bar{s}=5.0$  is exhausted). But whatever new policy is chosen for such "unlucky" cases, this choice will only affect the right tail of  $\mathcal{J}$ 's distribution; i.e.,  $\mathbb{P}(\mathcal{J} \geq \tilde{s})$  will be affected only for  $\tilde{s} > \bar{s}$ .

Returning to the optimal probability of success v(q, p, s) shown in Fig 3b, we observe that v has particularly large gradient near the level curves of the deterministic-optimal value function u shown in Fig 1a. (The particular level curve of u near which v changes the most is again s-dependent as the budget decreases.) If the remaining budget is relatively low (e.g., s = 1.5), one can see from Fig 3b that there is no chance to stabilize the tumor within this budget unless the GLY is already low (and a short burst of drug therapy would likely be enough) or VOP is high (and the no-drugs dynamics will bring us to a low GLY concentration later on). Consequently, the optimal policy for s = 1.5 is to not use drugs for the majority of tumor states.

The contrast in threshold-specific performance is easy to explain when the deterministic-optimal and threshold-aware policies prescribe different actions from the very beginning. To illustrate this, we consider  $(q_0, p_0) = (0.27, 0.4)$ , for which  $d_{\star} = d_{\text{max}}$  while  $d_{*}^{\bar{s}} = 0$  for a range of  $\bar{s}$  values; see Fig 5a and 5b for representative paths and Fig 5c for the respective CDFs. Under the deterministic-optimal policy (whose CDF is shown in blue), only 50% of simulations yield the cost not exceeding 4.71. A threshold-aware policy (implemented for  $\bar{s} = 4.71$ , with CDF shown in pink) maximizes this  $\mathbb{P}(\mathcal{J} \leq \bar{s})$  and succeeds in 63.7% of all cases. The potential for

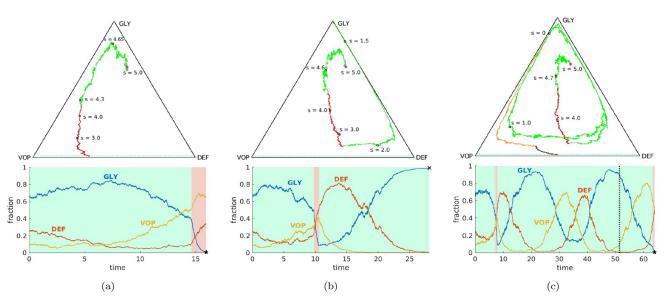


Fig 4. Representative sample paths starting from the same initial state ( $q_0$ ,  $p_0$ ) = (0.26, 0.665) (magenta dot) and the same initial budget  $\bar{s}=5$ . Top row: sample paths on a GLY-DEF-VOP triangle. (a) eventual stabilization with a cost of 4.70 (within the budget); (b) eventual death; (c) failure by running out of budget (eventual stabilization with a total cost of 7.80 by switching to the deterministic-optimal policy after s=0). Some representative tumor states along these paths (with indications of how much budget is left) are marked by *black squares*. In (c), the part where  $\mathcal{J}>5$  is specified in *orange* (no drugs) and *brown* (at MTD level).**Bottom row:** evolution of sub-populations with respect to time based on the sample paths from the top row. Here we use *light green* and *light pink backgrounds* to indicate the time interval(s) of prescribing no drugs and of prescribing drugs at the MTD-rate, respectively. We use *black pentagrams* and *black crosses* to indicate eventual stabilization and death, respectively. In (c), we use a *dashed black* line to indicate the budget depletion time  $t_{\sharp}$ .

improvement is even more significant with lower threshold values. For instance, we see that  $\mathbb{P}(\mathcal{J}(d_\star) \leq 4.35) < 10\%$ , while our threshold-aware policy (implemented for  $\bar{s} = 4.35$ , with CDF shown in orange) ensures that  $\mathbb{P}(\mathcal{J}(d_\star^{\bar{s}}) \leq 4.35) = \nu(q_0, p_0, \bar{s}) \approx 45.6\%$ . This improvement can also be translated to simple medical terms: starting from this initial tumor configuration, the deterministic-optimal policy will likely keep using the drugs at the maximum rate  $d_{\max}$  all the way to stabilization; see Fig 5a. In contrast, our threshold-aware policies tend not to prescribe drugs until GLY is relatively low and VOP is relatively high; see Fig 5b. As a result, the patient would suffer less toxicity from drugs in most scenarios.

It is worth noting that each threshold-aware policy maximizes the probability of success for a single/specific threshold value only. E.g., for all the pink/orange CDFs we have provided, the probability of success is only maximized at those pink/orange dots. Moreover, we clearly see from Fig 5c that the probability of  $\mathcal J$  not exceeding any  $\tilde s \leq 4.35$  is lower on the pink CDF than on the orange CDF (computed for  $\bar s = 4.35$ ). Intuitively, this is not too surprising. In the early stages of treatment, a (pink) policy computed to maximize the chances of not exceeding  $\bar s = 4.71$  is more aggressive in using the drugs and thus spends the "budget" quicker than the (orange) policy, which starts from a lower initial budget  $\bar s = 4.35$ . This is also consistent with the budget-dependent sizes of drugs-on regions in Fig 3a.

#### 3.2 Policies, trajectories, and CDFs for the SR-model

We now turn to the SR model system described in §2.3. Our numerical experiments use  $d_{\text{max}} = 3$ ,  $\gamma_{\text{r}} = 1 - \gamma_{\text{f}} = 10^{-2}$ ,  $\delta = 0.05$ , and volatilities  $\sigma_{\text{R}} = \sigma_{\text{S}} = 0.15$ . For other parameter values, see S1 Text §E.2.

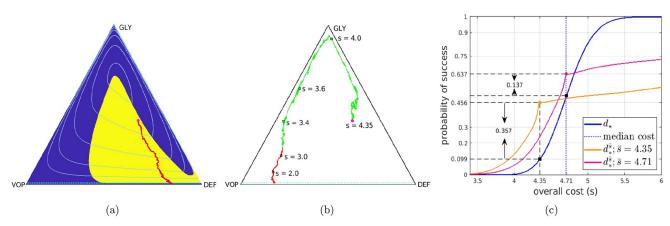


Fig 5. Comparison between threshold-aware policies and the deterministic-optimal policy. Starting from an initial state  $(q_0, p_0) = (0.27, 0.4)$  (magenta dot): (a) a sample path with cost 4.75 under the deterministic-optimal policy; (b) a sample path starting at  $\bar{s} = 4.35$  with a realized total cost of 4.02 under the (orange) threshold-aware policy; (c) CDFs of the cumulative cost  $\mathcal{J}$  approximated using  $10^5$  random simulations. In (c), the *solid blue* curve is the CDF generated with the deterministic-optimal policy. Its median (*dashed blue* line) is 4.71 while its mean conditioning on success is 4.72. The *solid orange* curve is the CDF generated with the threshold-aware policy with  $\bar{s} = 4.35$ ; and the *solid pink* curve is the CDF generated with the threshold-aware policy with  $\bar{s} = 4.71$ .

We show the representative s-slices of threshold aware policies and the corresponding success probabilities in Fig 6. Similarly to the EGT-model, we observe that the drug-on regions (shown in yellow) are strongly budget-dependent and quite different from the ones specified by  $d_{\star}$  in Fig 2b. We note that the drugs-on region generally shrinks in size (toward the Q=1 line, where only S cells are present) as the budget s decreases. For even tighter budgets, this yellow region becomes disconnected, prescribing the drugs for large P values (to substantially decrease the tumor size) and in a thin layer near  $\Delta_{\rm succ}$  (where a short burst of drugs is likely sufficient).

In Fig 7, we provide sample random trajectories and compare the performance of three different policies: the deterministically optimal  $d_{\star}$  and the threshold-aware  $d_{\star}^{\bar{s}}$  implemented for two different thresholds  $\bar{s}=69.45$  and  $\bar{s}=60$ . A suitable choice of the initial tumor configuration is less obvious for this example and deserves a separate comment. For many multi-population models, it is reasonable to assume that the system had approached some drug-free coexistence equilibrium before the tumor was detected and the therapy started. But since the model described in [3] does not include mutations, it also does not have a drug-free coexistence equilibrium. In our testing of various drug policies, we choose the initial tumor with 96% of sensitive cells and the tumor size at 90% of the carrying capacity. Since the resistant cells are much larger [3], this corresponds to initial conditions  $(q_0, p_0) = (0.45, 0.9)$ .

Despite the fact that all three tested policies use no drugs at the very beginning, the deterministic-optimal policy typically starts prescribing drugs much earlier. See the comparison of sample trajectories under  $d_{\star}$  and  $d_{\star}^{\bar{s}}$  in Fig 7a and 7b. As a result, our threshold-aware policy (implemented for  $\bar{s}=69.45$ , with CDF shown in pink) improves  $\mathbb{P}(\mathcal{J}\leq\bar{s})$  to 67.4% from 50% produced by  $d_{\star}$ . This advantage is even more significant with lower thresholds. E.g.,  $\mathbb{P}(\mathcal{J}(d_{\star})\leq 60)$  is only 19.6%, while our threshold-aware policy (implemented for  $\bar{s}=60$ , with CDF shown in orange) more than doubles this probability of under-threshold remission to 39.8%.

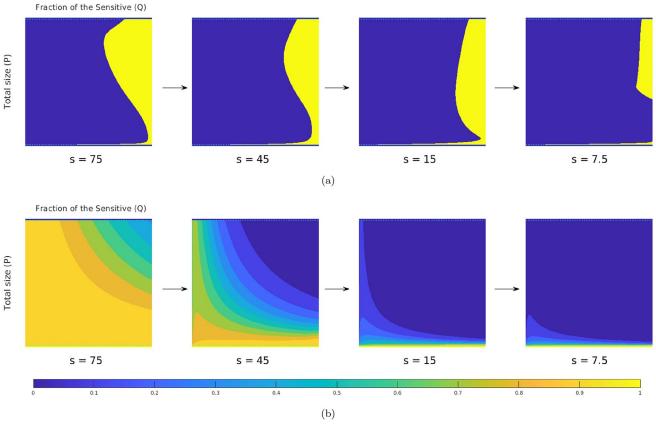


Fig 6. Representative slices of the threshold-aware optimal policy (top row) and the corresponding probability of success (bottom row) for the Carrère example. Each square represents all possible tumor states (sizes and compositions). The horizontal axis is the *fraction of the Sensitive* (Q) and the vertical axis is the *total population* (P). Top row shows the policy, which prescribes the optimal instantaneous decisions on drug usage given the indicated remaining budget (s) and the current tumor state. Bottom row shows the probability of "eradication within the budget" if the optimal policy is followed from this time point and onward. Each column corresponds to a specific budget level s, which is shown below each square. The arrows indicate the natural decrease of the remaining budget while implementing the policy.

#### 4 Discussion

That cancers evolve during therapy is now an accepted fact, and is slowly being incorporated into therapeutic decision making. In some cases, this can be implemented simply by changing from one targeted therapy to another, but in most, where tumors are a heterogeneous mixture of interacting phenotypes, this is not feasible. In these cases, ecological thinking is rising to the fore in the form of adaptive therapy. Until recently, clinical trials, and theoretical investigations, of adaptive therapy have relied on *a priori* assumptions of the underlying interactions, and their effects on tumor composition over time. Several studies, both *in vitro* [20] and *in vivo* [19, 52], however, have begun to provide methods for more rigorous quantification of these interactions. As these tools mature, the next challenges will be to understand these interactions in patients and to exploit them in improving personalized treatment.

The presented approach is a step in this direction, aiming to limit the probability of high-cost outcomes in the presence of stochastic perturbations. It is applicable to a broad class of stochastic cancer models and therapy goals (e.g., tumor eradication or stabilization). While it is standard to tune the treatment plan to maximize the probability of reaching its goal, we go farther and maximize the probability of goal attainment without exceeding a prescribed

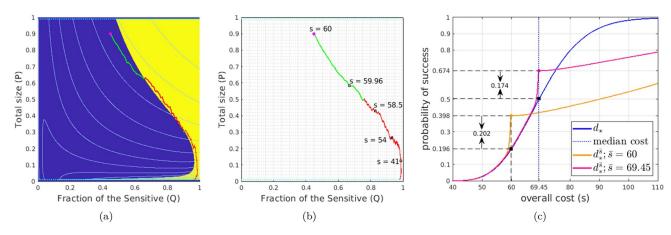


Fig 7. Comparison between threshold-aware policies and the deterministic-optimal policy. Starting from an initial state  $(q_0, p_0) = (0.45, 0.9)$  (magenta dot): (a) a sample path with cost 57.3 under the deterministic-optimal policy; (b) a sample path starting at  $\bar{s} = 60$  with a total cost of 53.63 under the (orange) threshold-aware policy; (c) CDFs of the cumulative cost  $\mathcal{J}$  with  $10^5$  samples. In (c), the *solid blue* curve is the CDF generated with the deterministic-optimal policy. Its median (*dashed blue* line) is 69.45 while its mean conditioning on success is 70.5. The *solid orange* curve is the CDF generated with the threshold-aware policy with  $\bar{s} = 60$ ; and the *solid pink* curve is the CDF generated with the threshold-aware policy with  $\bar{s} = 69.45$ . See S1 Text §E.3 for time-evolution plots associated with sample paths in (a) and (b).

threshold on cumulative cost (interpreted as a combination of the total drugs used, cumulative disease burden, and the time to remission/stabilization). We show that these optimal treatment policies become *threshold-aware*, with the drugs-on/drugs-off regions changing as the treatment progresses and the initial "cost budget" (for meeting the chosen threshold) gradually decreases. The comparison of CDFs generated for the deterministic-optimal policy and threshold-aware policies demonstrates clear advantages of the latter, often resulting in a significant reduction of drugs used to treat the patient.

More generally, dynamic programming provides an excellent framework for finding optimal treatment policies by solving Hamilton-Jacobi-Bellman (HJB) equations. The fact that these policies are recovered in feedback form makes this approach particularly suitable for optimization of adaptive therapies. But even though the use of general optimal control in cancer treatment is by now common [1], the same is not true for the more robust HJB-based methods, which so far have been used in only a handful of cancer-related applications [4, 8, 23, 29, 56–58]. This is partly due to the HJBs' well-known *curse of dimensionality*: the rapid increase in computational costs when the system state becomes higher-dimensional. This is a relevant limitation since our threshold-aware approach introduces the "budget" as an additional component of the state. Similarly to the presented examples, our current implementation would be easy to adopt to any cancer model based on a two-dimensional (q, p) state space, with the budget s adding the third dimension. For cancer models with a larger number of subpopulations, the general approach would remain the same, but the approximate HJB-solver would likely need to rely on sparse grids [59], tensor decompositions [60], or deep neural networks [61].

The presented examples did not model any mutations, but we note that this is not really a limitation of the method itself. E.g., drug-usage-dependent mutations would be easy to incorporate into our EGT-based example by switching to a Replicator-Mutator ODE/SDE eco-evolutionary model [62, 63]. We did not pursue such examples here primarily to make for an easier comparison with prior work [29] and to limit the number of model parameters.

Our SDE models of global (or environmental) stochastic perturbations to subpopulation fitnesses and intrinsic growth rates are based on perspectives well-established in biological applications [42, 51]. While our focus on environmental stochasticity is motivated by "averagingout" the variability within each subpopulation, it is worth noting that this assumption is not always justifiable. Whenever the subpopulation size is sufficiently small, the demographic stochasticity becomes crucially important. (This is also the regime in which the validity of ODE/ SDE models is far less obvious.) Even though we do not deal with this important limitation here, we note that our threshold-aware approach can be used with a variety of perturbation types, including jump-diffusion processes, which could be used to build future models that account for demographic stochasticity in these special small-subpopulation regimes. Such discontinuous jump-transitions (e.g., reflecting possible subpopulation extinctions) can be naturally handled in our framework. For instance, a similar method has been developed in [64] for controlling "piecewise-deterministic" processes, where perturbations happen at discrete points in time and amount to abrupt switches in system dynamics. More recently, our framework was also used to control the hybrid dynamics of a sailboat navigating in stochastically changing wind conditions and trying to reach the destination prior to a specified deadline  $\bar{s}$  [65]. We note that dynamic programming is also used in discrete population models focused on demographic stochasticity [66]. It will be also interesting to investigate the usability of our approach in that discrete setting.

Sensitivity with respect to threshold variation can be tested by comparing CDFs of  $d_*^{\bar{s}}$  for different  $\bar{s}$  values. While it is also possible to perform a similar comparison under perturbation of model parameters, we believe that another approach is more promising: any bounded uncertainty in parameter values can be treated as a "game against nature," leading to a Hamilton-Jacobi-Isaacs PDE, whose solution will yield policies optimizing the threshold-performance in the "worst parameter variation" scenarios [64].

Another important extension will be to move to "partial observability" since the state of the tumor is only occasionally assessed directly through biopsies and some proxy measurements have to be used at all other times [8]. Finally, it will be also interesting to study the multiobjective control problem of optimizing threshold-aware policies for two different threshold values simultaneously.

In summary, we have presented a theoretical and computational advance for the toolbox of evolutionary therapy, a new subfield of medicine focused on using knowledge of evolutionary responses to inform therapeutic scheduling. While there are a number of cancer trials using this type of evolutionary-informed thinking, most are based on heuristic designs and are not formulated to consider the underlying stochasticities. Developing a theoretical foundation for future clinical studies requires EGT models directly grounded in objectively measurable biology [20]. Therapy optimization based on such models requires efficient computational methods, particularly in the presence of stochastic perturbations. We hope that the general approach presented here will be useful for a broad range of increasingly accurate stochastic cancer models.

# **Supporting information**

**S1 Text. Supplementary materials.** Mathematical details and additional computational experiments. (PDF)

# **Acknowledgments**

The authors are grateful to Roberto Ferretti and Lars Grüne for their advice on some aspects of numerical methods used in this project.

#### **Author Contributions**

Conceptualization: MingYi Wang, Alexander Vladimirsky.

Formal analysis: MingYi Wang.

**Funding acquisition:** Jacob G. Scott, Alexander Vladimirsky.

Investigation: MingYi Wang.

Methodology: MingYi Wang, Alexander Vladimirsky.

Software: MingYi Wang.

**Supervision:** Jacob G. Scott, Alexander Vladimirsky.

Validation: MingYi Wang. Visualization: MingYi Wang.

Writing - original draft: MingYi Wang, Jacob G. Scott, Alexander Vladimirsky.

Writing - review & editing: MingYi Wang, Jacob G. Scott, Alexander Vladimirsky.

#### References

- Schättler H, Ledzewicz U. Optimal Control for Mathematical Models of Cancer Therapies. vol. 42 of Interdisciplinary Applied Mathematics. New York, NY: Springer New York; 2015.
- Martin RB, Fisher ME, Minchin RF, Teo KL. Optimal control of tumor size used to maximize survival time when cells are resistant to chemotherapy. Mathematical Biosciences. 1992; 110(2):201–219.
- 3. Carrère C. Optimization of an in vitro chemotherapy to avoid resistant tumours. Journal of Theoretical Biology. 2017; 413:24–33. https://doi.org/10.1016/j.jtbi.2016.11.009
- Carrère C, Zidani H. Stability and reachability analysis for a controlled heterogeneous population of cells. Optimal Control Applications and Methods. 2020; 41(5):1678–1704. <a href="https://doi.org/10.1002/oca.2627">https://doi.org/10.1002/oca.2627</a>
- Day RS. Treatment sequencing, asymmetry, and uncertainty: protocol strategies for combination chemotherapy. Cancer Research. 1986; 46(8):3876–3885. PMID: 3731062
- Coldman AJ, Murray JM. Optimal control for a stochastic model of cancer chemotherapy. Mathematical Biosciences. 2000; 168(2):187–200. <a href="https://doi.org/10.1016/S0025-5564(00)00045-6">https://doi.org/10.1016/S0025-5564(00)00045-6</a> PMID: 11121565
- Katouli AA, Komarova NL. The worst drug rule revisited: mathematical modeling of cyclic cancer treatments. Bulletin of mathematical biology. 2011; 73:549–584. https://doi.org/10.1007/s11538-010-9539-y
   PMID: 20396972
- Fischer A, Vázquez-García I, Mustonen V. The value of monitoring to control evolving populations. Proceedings of the National Academy of Sciences. 2015; 112(4):1007–1012. <a href="https://doi.org/10.1073/pnas.1409403112">https://doi.org/10.1073/pnas.1409403112</a> PMID: 25587136
- Haslam A, Kim M, Prasad V. Updated estimates of eligibility for and response to genome-targeted oncology drugs among US cancer patients, 2006-2020. Annals of Oncology. 2021; 32(7):926–932. https://doi.org/10.1016/j.annonc.2021.04.003 PMID: 33862157
- Gillies RJ, Verduzco D, Gatenby RA. Evolutionary dynamics of carcinogenesis and why targeted therapy does not work. Nature Reviews Cancer. 2012; 12(7):487–493. https://doi.org/10.1038/nrc3298 PMID: 22695393
- Imamovic L, Sommer MO. Use of collateral sensitivity networks to design drug cycling protocols that avoid resistance development. Science translational medicine. 2013; 5(204):204ra132–204ra132. https://doi.org/10.1126/scitranslmed.3006609 PMID: 24068739

- Nichol D, Jeavons P, Fletcher AG, Bonomo RA, Maini PK, Paul JL, et al. Steering evolution with sequential therapy to prevent the emergence of bacterial antibiotic resistance. PLoS computational biology. 2015; 11(9):e1004493. https://doi.org/10.1371/journal.pcbi.1004493 PMID: 26360300
- Nichol D, Rutter J, Bryant C, Hujer AM, Lek S, Adams MD, et al. Antibiotic collateral sensitivity is contingent on the repeatability of evolution. Nature communications. 2019; 10(1):1–10. <a href="https://doi.org/10.1038/s41467-018-08098-6">https://doi.org/10.1038/s41467-018-08098-6</a> PMID: 30659188
- Maltas J, Wood KB. Pervasive and diverse collateral sensitivity profiles inform optimal strategies to limit antibiotic resistance. PLoS biology. 2019; 17(10):e3000515. https://doi.org/10.1371/journal.pbio. 3000515 PMID: 31652256
- Zhao B, Sedlak JC, Srinivas R, Creixell P, Pritchard JR, Tidor B, et al. Exploiting temporal collateral sensitivity in tumor clonal evolution. Cell. 2016; 165(1):234–246. <a href="https://doi.org/10.1016/j.cell.2016.01.045">https://doi.org/10.1016/j.cell.2016.01.045</a>
   PMID: 26924578
- Dhawan A, Nichol D, Kinose F, Abazeed ME, Marusyk A, Haura EB, et al. Collateral sensitivity networks reveal evolutionary instability and novel treatment strategies in ALK mutated non-small cell lung cancer. Scientific reports. 2017; 7(1):1–9. https://doi.org/10.1038/s41598-017-00791-8 PMID: 28450729
- Gatenby RA, Silva AS, Gillies RJ, Frieden BR. Adaptive therapy. Cancer research. 2009; 69(11):4894–4903. https://doi.org/10.1158/0008-5472.CAN-08-3658 PMID: 19487300
- Zhang J, Cunningham JJ, Brown JS, Gatenby RA. Integrating evolutionary dynamics into treatment of metastatic castrate-resistant prostate cancer. Nature communications. 2017; 8(1):1–9. <a href="https://doi.org/10.1038/s41467-017-01968-5">https://doi.org/10.1038/s41467-017-01968-5</a> PMID: 29180633
- Enriquez-Navas PM, Kam Y, Das T, Hassan S, Silva A, Foroutan P, et al. Exploiting evolutionary principles to prolong tumor control in preclinical models of breast cancer. Science translational medicine. 2016; 8(327):327ra24–327ra24. https://doi.org/10.1126/scitranslmed.aad7842 PMID: 26912903
- Kaznatcheev A, Peacock J, Basanta D, Marusyk A, Scott JG. Fibroblasts and alectinib switch the evolutionary games played by non-small cell lung cancer. Nature ecology & evolution. 2019; 3(3):450–456. https://doi.org/10.1038/s41559-018-0768-z PMID: 30778184
- Farrokhian N, Maltas J, Dinh M, Durmaz A, Ellsworth P, Hitomi M, et al. Measuring competitive exclusion in non–small cell lung cancer. Science Advances. 2022; 8(26):eabm7212. <a href="https://doi.org/10.1126/sciadv.abm7212">https://doi.org/10.1126/sciadv.abm7212</a> PMID: 35776787
- Chisholm RH, Lorenzi T, Clairambault J. Cell population heterogeneity and evolution towards drug resistance in cancer: biological and mathematical assessment, theoretical treatment optimisation. Biochimica et Biophysica Acta (BBA)-General Subjects. 2016; 1860(11):2627–2645. https://doi.org/10. 1016/j.bbagen.2016.06.009 PMID: 27339473
- Kuosmanen T, Cairns J, Noble R, Beerenwinkel N, Mononen T, Mustonen V. Drug-induced resistance evolution necessitates less aggressive treatment. PLoS computational biology. 2021; 17(9):e1009418. https://doi.org/10.1371/journal.pcbi.1009418 PMID: 34555024
- Greene JM, Gevertz JL, Sontag ED. Mathematical approach to differentiate spontaneous and induced evolution to drug resistance during cancer treatment. JCO clinical cancer informatics. 2019; 3:1–20. https://doi.org/10.1200/CCI.18.00087 PMID: 30969799
- Cunningham JJ, Brown JS, Gatenby RA, Staňková K. Optimal control to develop therapeutic strategies for metastatic castrate resistant prostate cancer. Journal of theoretical biology. 2018; 459:67–78. https://doi.org/10.1016/j.jtbi.2018.09.022 PMID: 30243754
- **26.** West JB, Dinh MN, Brown JS, Zhang J, Anderson AR, Gatenby RA. Multidrug cancer therapy in metastatic castrate-resistant prostate cancer: an evolution-based strategy. Clinical Cancer Research. 2019; 25(14):4413–4421. https://doi.org/10.1158/1078-0432.CCR-19-0006 PMID: 30992299
- West J, You L, Zhang J, Gatenby RA, Brown JS, Newton PK, et al. Towards multidrug adaptive therapy. Cancer research. 2020; 80(7):1578–1589. <a href="https://doi.org/10.1158/0008-5472.CAN-19-2669">https://doi.org/10.1158/0008-5472.CAN-19-2669</a> PMID: 31948939
- Cunningham J, Thuijsman F, Peeters R, Viossat Y, Brown J, Gatenby R, et al. Optimal control to reach eco-evolutionary stability in metastatic castrate-resistant prostate cancer. PLoS One. 2020; 15(12): e0243386. https://doi.org/10.1371/journal.pone.0243386 PMID: 33290430
- 29. Gluzman M, Scott JG, Vladimirsky A. Optimizing adaptive cancer therapy: dynamic programming and evolutionary game theory. Proceedings of the Royal Society B. 2020; 287(1925):20192454. <a href="https://doi.org/10.1098/rspb.2019.2454">https://doi.org/10.1098/rspb.2019.2454</a> PMID: 32315588
- Gupta PB, Fillmore CM, Jiang G, Shapira SD, Tao K, Kuperwasser C, et al. Stochastic state transitions give rise to phenotypic equilibrium in populations of cancer cells. Cell. 2011; 146(4):633–644. <a href="https://doi.org/10.1016/j.cell.2011.07.026">https://doi.org/10.1016/j.cell.2011.07.026</a> PMID: 21854987
- Kumar N, Cramer GM, Dahaj SAZ, Sundaram B, Celli JP, Kulkarni RV. Stochastic modeling of phenotypic switching and chemoresistance in cancer cell populations. Scientific reports. 2019; 9(1):1–10.

- Lande R, Engen S, Saether BE, et al. Stochastic population dynamics in ecology and conservation. Oxford University Press on Demand; 2003.
- Wilke CO. Quasispecies theory in the context of population genetics. BMC evolutionary biology. 2005;
   5(1):1–8. https://doi.org/10.1186/1471-2148-5-44 PMID: 16107214
- Lauring AS, Andino R. Quasispecies theory and the behavior of RNA viruses. PLoS pathogens. 2010; 6
   (7):e1001005. https://doi.org/10.1371/journal.ppat.1001005 PMID: 20661479
- Engen S, Bakke Ø, Islam A. Demographic and environmental stochasticity—concepts and definitions. Biometrics. 1998; p. 840–846. https://doi.org/10.2307/2533838
- Iram S, Dolson E, Chiel J, Pelesko J, Krishnan N, Güngör Ö, et al. Controlling the speed and trajectory
  of evolution with counterdiabatic driving. Nature Physics. 2021; 17(1):135–142. <a href="https://doi.org/10.1038/s41567-020-0989-3">https://doi.org/10.1038/s41567-020-0989-3</a>
- **37.** Braumann CA. Environmental versus demographic stochasticity in population growth. In: Workshop on Branching Processes and Their Applications. Springer; 2010. p. 37–52.
- Allen E. Modeling with Itô stochastic differential equations. vol. 22. Springer Science & Business Media; 2007.
- Fleming WH, Rishel RW. Deterministic and stochastic optimal control. vol. 1. Springer Science & Business Media; 2012.
- Bardi M, Dolcetta I. Optimal Control and Viscosity Solutions of Hamilton-Jacobi-Bellman Equations. Boston, MA: Birkhäuser; 1997.
- Fleming WH, Soner HM. Controlled Markov processes and viscosity solutions. vol. 25. Springer Science & Business Media; 2006.
- Browning AP, Sharp JA, Mapder T, Baker CM, Burrage K, Simpson MJ. Persistence as an optimal hedging strategy. Biophysical Journal. 2021; 120(1):133–142. <a href="https://doi.org/10.1016/j.bpj.2020.11.2260">https://doi.org/10.1016/j.bpj.2020.11.2260</a> PMID: 33253635
- Nourmohammad A, Eksin C. Optimal evolutionary control for artificial selection on molecular phenotypes. Physical Review X. 2021; 11(1):011044. https://doi.org/10.1103/PhysRevX.11.011044
- Pontryagin L, Boltyanskii V, Gamkrelidze R, Mishchenko E. The mathematical theory of optimal processes. New York: John Wiley & Sons, Inc.; 1962.
- **45.** Zouhri S, El Baroudi M, Saadi S. Optimal Control with Isoperimetric Constraint for Chemotherapy of Tumors. International Journal of Applied and Computational Mathematics. 2022; 8(4):215. https://doi.org/10.1007/s40819-022-01425-y
- Hamdache A, Elmouki I, Saadi S. Optimal control with an isoperimetric constraint applied to cancer immunotherapy. International Journal of Computer Applications. 2014; 94(15). <a href="https://doi.org/10.5120/16421-6073">https://doi.org/10.5120/16421-6073</a>
- Kumar A, Vladimirsky A. An efficient method for multiobjective optimal control and optimal control subject to integral constraints. Journal of Computational Mathematics. 2010; 28(4):517–551. <a href="https://doi.org/10.4208/jcm.1003-m0015">https://doi.org/10.4208/jcm.1003-m0015</a>
- 48. Kaznatcheev A, Vander Velde R, Scott JG, Basanta D. Cancer treatment scheduling and dynamic heterogeneity in social dilemmas of tumour acidity and vasculature. British journal of cancer. 2017; 116 (6):785–792. https://doi.org/10.1038/bjc.2017.5 PMID: 28183139
- **49.** Taylor PD, Jonker LB. Evolutionary stable strategies and game dynamics. Mathematical biosciences. 1978; 40(1-2):145–156. https://doi.org/10.1016/0025-5564(78)90077-9
- **50.** Hofbauer J, Sigmund K. Evolutionary games and population dynamics. Cambridge university press;
- Fudenberg D, Harris C. Evolutionary dynamics with aggregate shocks. Journal of Economic Theory. 1992; 57(2):420–441. https://doi.org/10.1016/0022-0531(92)90044-I
- Marusyk A, Tabassum DP, Altrock PM, Almendro V, Michor F, Polyak K. Non-cell-autonomous driving of tumour growth supports sub-clonal heterogeneity. Nature. 2014; 514(7520):54–58. <a href="https://doi.org/10.1038/nature13556">https://doi.org/10.1038/nature13556</a> PMID: 25079331
- Vander Velde R, Yoon N, Marusyk V, Durmaz A, Dhawan A, Miroshnychenko D, et al. Resistance to targeted therapies as a multifactorial, gradual adaptation to inhibitor specific selective pressures. Nature communications. 2020; 11(1):1–13. https://doi.org/10.1038/s41467-020-16212-w PMID: 32409712
- Schreiber SJ, Benaïm M, Atchadé KA. Persistence in fluctuating environments. Journal of Mathematical Biology. 2011; 62:655–683. https://doi.org/10.1007/s00285-010-0349-5 PMID: 20532555
- Hening A, Nguyen DH, Chesson P. A general theory of coexistence and extinction for stochastic ecological communities. Journal of Mathematical Biology. 2021; 82(6):56. https://doi.org/10.1007/s00285-021-01606-1 PMID: 33963448

- 56. Nowakowski A, Popa A. A Dynamic Programming Approach for Approximate Optimal Control for Cancer Therapy. J Optim Theory Appl. 2013; 156(2):365–379. https://doi.org/10.1007/s10957-012-0137-z
- Alamir M. Robust feedback design for combined therapy of cancer. Optimal Control Applications and Methods. 2014; 35(1):77–88. https://doi.org/10.1002/oca.2057
- Jeong YD, Kim KS, Roh Y, Choi S, Iwami S, Jung IH, et al. Optimal Feedback Control of Cancer Chemotherapy Using Hamilton-Jacobi-Bellman Equation. Complexity. 2022; 2022. <a href="https://doi.org/10.1155/2022/2158052">https://doi.org/10.1155/2022/2158052</a>
- **59.** Miksis ZM, Zhang YT. Sparse-Grid Implementation of Fixed-Point Fast Sweeping WENO Schemes for Eikonal Equations. Communications on Applied Mathematics and Computation. 2022; p. 1–27.
- 60. Dolgov S, Kalise D, Kunisch KK. Tensor decomposition methods for high-dimensional Hamilton— Jacobi–Bellman equations. SIAM Journal on Scientific Computing. 2021; 43(3):A1625–A1650. https://doi.org/10.1137/19M1305136
- **61.** Xuanxi Zhang, Jihao Long, Wei Hu, Weinan E and Jiequn Han. Initial Value Problem Enhanced Sampling for Closed-Loop Optimal Control Design with Deep Neural Networks. preprint: <a href="https://arxiv.org/abs/2209.04078">https://arxiv.org/abs/2209.04078</a>.
- Hadeler KP. Stable polymorphisms in a selection model with mutation. SIAM Journal on Applied Mathematics. 1981; 41(1):1–7. https://doi.org/10.1137/0141001
- Hofbauer J. The selection mutation equation. Journal of mathematical biology. 1985; 23:41–53. https://doi.org/10.1007/BF00276557 PMID: 4078498
- Cartee E, Nellis A, Van Hook J, Farah A, Vladimirsky A. Quantifying and managing uncertainty in piecewise-deterministic Markov processes. SIAM/ASA Journal on Uncertainty Quantification. 2023; 11 (3):814–847. https://doi.org/10.1137/20M1357275
- **65.** Wang M, Patnaik N, Somalwar A, Wu J, Vladimirsky A. Risk-aware stochastic control of a sailboat. ACC-2024; preprint: https://arxiv.org/abs/2309.13436.
- 66. Mononen T, Kuosmanen T, Cairns J, Mustonen V. Understanding cellular growth strategies via optimal control. Journal of the Royal Society Interface. 2023; 20(198):20220744. <a href="https://doi.org/10.1098/rsif.2022.0744">https://doi.org/10.1098/rsif.2022.0744</a> PMID: 36596459