# NEW TOOLS FOR COMPARING CLASSICAL AND NEURAL ODE MODELS FOR TUMOR GROWTH

ANTHONY D. BLAOM

*Department of Computer Science*
*University of Auckland*
*New Zealand*


SAMUEL OKON

*German Research Center for Artificial Intelligence*
*Kaiserslautern*
*Germany*

ABSTRACT. A new computational tool `TumorGrowth.jl` for modeling tumor growth is introduced. The tool allows the comparison of standard textbook models, such as General Bertalanffy and Gompertz, with some newer models, including, for the first time, neural ODE models. As an application, we revisit a human meta-study of non-small cell lung cancer and bladder cancer lesions, in patients undergoing two different treatment options, to determine if previously reported performance differences are statistically significant, and if newer, more complex models perform any better. In a population of examples with at least four time-volume measurements available for calibration, and an average of about 6.3, our main conclusion is that the General Bertalanffy model has superior performance, on average. However, where more measurements are available, we argue that more complex models, capable of capturing rebound and relapse behavior, may be better choices.

## 1. INTRODUCTION

We investigate the performance of models for the growth of tumors, as measured by a single parameter, typically the volume. Models under consideration are based on solving ordinary differential equations (ODE's). These ODE's have unknown parameters, which typically means forecasting a tumor's future size is only possible after calibrating the model using the current clinical history. We introduce a new package, `TumorGrowth.jl` [1], to automate this procedure, for a battery of classical models, such the General Bertalanffy model [2, 3], as well more complex models, which include, for the first time, neural ODE's [4]. Custom models can also be implemented.

1.1. **Previous evaluations of classical model performance.** As an application of the new tool, we revisit the study of Laleh et al. [5], the first of its kind, which includes an out-of-sample evaluation of the accuracy of classical "textbook" models in a meta-study of 652 tumors, in humans undergoing chemotherapy or cancer immunotherapy. Patients in that study are either non-small cell lung cancer or bladder cancer sufferers. The specific treatments compared are Atezolizumab (previously known as MPDL3280A) and Docetaxel. Refer to [5] for details. Models in the meta-study are ranked based on the mean absolute error on a holdout test set. The models compared are: exponential, logistic, classical Bertalanffy, General Bertalanffy, classical Gompertz, and the General

---

*E-mail addresses*: `anthony.blaom@gmail.com`, `samuel.okon@dfki.de`.

Gompertz models. Each study and study arm gets a separate treatment in [5] and aggregated scores are not reported. However the conclusion is a general trend favoring the General Bertalanffy and the classical Gompertz models.

That said, no statistical significance is attached to these results. It could be that, by chance, the actual expected performance of these models is different from the reported performance, and the likelihood of this scenario remains unquantified in the study. The present study provides the missing statistical analysis, and adds some newer models to the comparison.
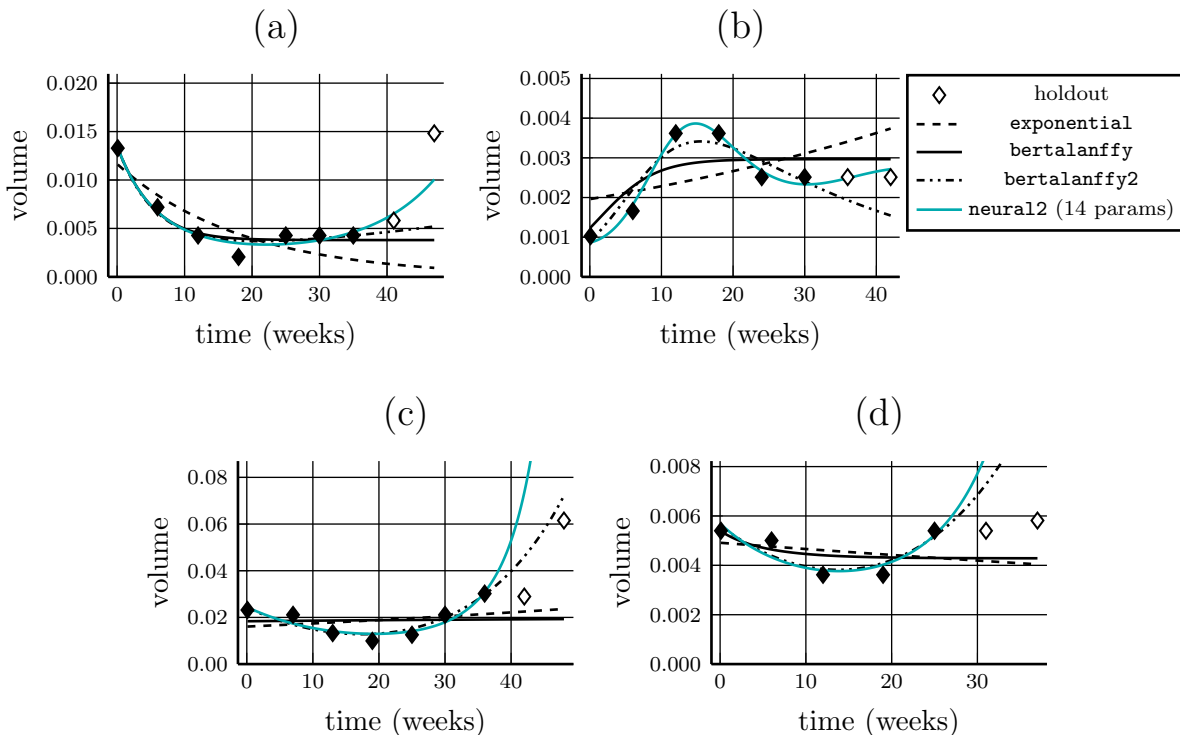


FIGURE 1. Selected model comparisons for observations with relapse or rebound. Solid diamonds indicate calibration data, open diamonds subsequent observations not used in calibration. Two classical models, exponential and General Berta-lanffy (`bertalanffy`), are compared with a 2D generalization of General Bertalanffy (`bertalanffy2`) and a 2D, 14-parameter neural ODE (`neural2`). While the newer models perform better on the holdout data in (a), the classical models do better in (d). Results are mixed in cases (b) and (c).

1.2. **Capturing relapse or rebound behavior.** As they are first order, one-dimensional ODE's, none of the classical models in the Laleh et al. study can capture rebound or relapse phenomena: Solutions are always monotonically increasing, decreasing or constant. This is not surprising, as the classical models were constructed for untreated lesions. In the current study (and in `TumorGrowth.jl`) we consider two models in which the volume is coupled to a second latent variable (making them effectively second-order ODE's) and which do not suffer this limitation:

- A novel but simple 2D generalization of the General Bertalanffy model described in 3.2, with one additional parameter (5 total, including initial condition)
- A basic 2D neural ODE, with 14 parameters, described in 3.3.

A *neural ODE* is an ODE $\dot{x} = F(x, \theta)$ where $F(x, \theta)$ represents the output of an artificial neural network, with input $x$, and a system of internal weights and biases $\theta$ [4]. For simplicity, we focus, in our comparisons, on a particular neural ODE, i.e., on a particular network architecture. Note, however, that by varying the architecture (which is possible in `TumorGrowth.jl`) one can approximate any ODE arbitrarily well. This follows from the Universal Approximation Theorem for neural networks [6]. Although the precise mechanisms of tumor growth or mitigation may be unclear, one tends to believe the underlying process is nevertheless governed by *some* system of ordinary differential equations. That is, one views the use of a neural ODE as more "physically" justified than ordinary curve fitting with, say, polynomials.

As one can see from Figure 1(a), the two new models have the potential to make better predictions in the presence of a rebound or relapse. A neural ODE model can capture even more nuanced behavior, as demonstrated in (b). However, the danger of using models with more free parameters is over-fitting, which is clearly evident in (c) (for the neural ODE model) and (d) (for both new models). The present study tackles the question of which tendency — better prediction or over-fitting – is the more typical in the clinical context.

Obviously, answers to the question just posed depends on the number of observations available for calibration. Following Laleh et al., we restrict attention to observations in their meta-study exceeding six in number. Deviating slightly from Laleh et al., we hold back the last two, rather than three, observations for testing. The average number of observations in this restricted dataset, including hold-outs, is 8.3, with about 73% less than 9.

1.3. **Results.** We now summarize the results of calibration experiments detailed in Section 4. Our analysis treats the observations collated in the Laleh et al. study as a whole, with no attempt to stratify outcomes over individual study arms. We are doubtful statistically significant conclusions can be drawn at the level of study arms, because of small sample size, as it is already difficult in the aggregated case, as we now report.

Where there are differences in the studies, or in individual patients, about the timing of treatment, these are also conflated in the present analysis.

Figure 2 shows a boxplot for the prediction errors for each model. No statistically significant difference in model performance is evident in these results. Although it has dubious probative value, we record in Table 1 rankings based on the mean absolute prediction errors.
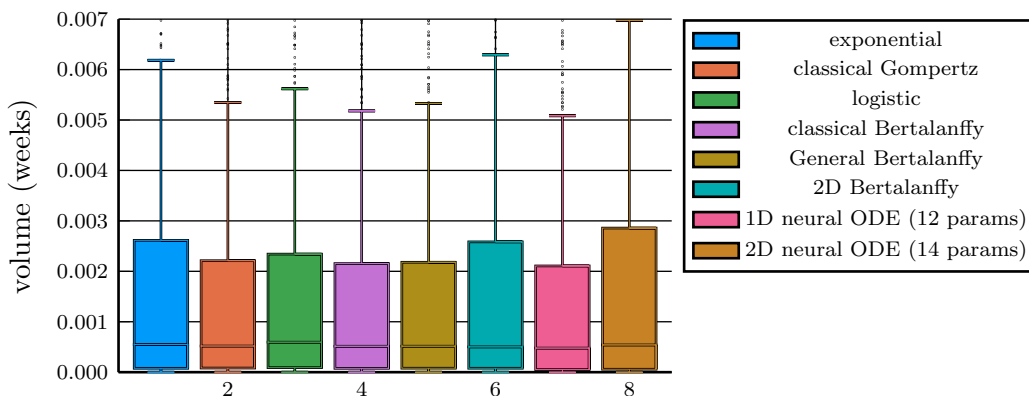


FIGURE 2. Box-and-whisker plots of the absolute prediction errors on a holdout set.

Turning to a more powerful test, we compute 95% bootstrap confidence intervals for the *pairwise difference* in individual lesion prediction errors. When such an interval contains zero, we declare a

| model | mean abs. error | num. parameters | section |
|---|---|---|---|
| General Bertalanffy | 0.0027 | 4 | 3.1 |
| classical Bertalanffy | 0.0028 | 3 | 3.1 |
| classical Gompertz | 0.0028 | 3 | 3.1 |
| logistic | 0.0029 | 3 | 3.1 |
| 1D neural ODE | 0.0031 | 12 | 3.1 |
| 2D General Bertalanffy | 0.0032 | 5 | 3.2 |
| exponential | 0.0033 | 2 | 3.4 |
| 2D neural ODE | 0.0046 | 14 | 3.3 |

TABLE 1. Models ranked by mean absolute error on a two-observation holdout set, but note large standard errors shown in Figure 2.

statistical tie, otherwise a win is declared for the better performing model. The outcomes of these match-ups, presented in Table 2, suggest:

- The General Bertalanffy model (`bertalanffy`) is superior to all other models, except the 1D neural ODE (`neural`), where no statistically significant difference is detectable. However, there is no detectable difference between the 1D neural model and *any* other model, which is a fair basis for rejecting it as an alternative in the current context.
- The 14-parameter 2D neural ODE tested in this study is inferior to all other models.
- The exponential model is inferior to most models and not demonstrably superior to any other classical models.
- Otherwise, there are no statistically significant differences between the models.

| model | c. Gomp. | log. | c. Bert. | G. Bert. | 2D Bert. | 1D neural | 2D neural |
|---|---|---|---|---|---|---|---|
| exponential | ↑ | draw | ↑ | ↑ | draw | draw | ← |
| classical Gompertz | | draw | draw | ↑ | draw | draw | ← |
| logistic | draw | | draw | ↑ | draw | draw | ← |
| classical Bertalanffy | draw | draw | | ↑ | draw | draw | ← |
| General Bertalanffy | ← | ← | ← | | ← | draw | ← |
| 2D Bertalanffy | draw | draw | draw | ↑ | | draw | ← |
| 1D neural ODE | draw | draw | draw | draw | draw | | ← |

TABLE 2. Head-to-head model comparisons. Where "draw" appears, the difference in model performance is not significant at the 5% level. Where there is a significant difference, an arrow points to the superior model.

## 2. THE TUMORGROWTH.JL PACKAGE

The `TumorGrowth.jl` package is a tool for calibrating and comparing models for tumor growth, as measured by a single parameter, typically volume [1]. It features:

- Implementations of the classical and General Bertalanffy, classical Gompertz and logistic models (see 3.1 below) as well as the simple model for exponential decay or growth (3.4)
- A two-dimensional generalization of General Bertalanffy with one extra parameter, for capturing rebound or relapse behavior (3.2)
- One and two-dimensional neural ODE models (3.3)

- Sophisticated control over parameter optimization, such as early stopping criteria
- An option to give higher weight to more recent measurements during calibration (not applied in the study above)
- Comprehensive options for optimization, either by gradient descent, Levenberg-Marquardt, or Powell's dog leg algorithm [7]
- Plotting functions to visualize results
- The ability to specify a custom model
- Convenient access to data from the human meta-study collated in [5].

For further details, and tutorials, the reader is referred to the comprehensive package documentation [1].

`TumorGrowth.jl` is written in pure `Julia`, which has the advantages of speed, high customizability, transparency and reproducibility. Beyond the classical case where analytic solutions to the underlying ODE's are known, model calibration involves differentiating numerically obtained solutions with respect to parameters and initial conditions. In this respect, development of `TumorGrowth.jl` was able to capitalize significantly on Julia's state-of-the-art `SciML` ecosystem [8], which implements automatic differentiation, and in particular, Pontryagin's adjoint method for differentiating ODE solutions.

## 3. Models

In this section we provide the detailed specification of models tested in our study reported above, and implemented in `TumorGrowth.jl`.

### 3.1. The General Bertalanffy model.
In its common formulation, Bertallanfy's model for lesion growth [9, 2] is

$$\frac{dv}{dt} = \eta v^m - \kappa v^n. \tag{1}$$

Here $v$ denotes lesion volume, $t$ is time; $\eta$, $\kappa$, while $m$ and $n$ are parameters. Bertalanffy argued that solutions to (1) are relatively insensitive to variations of $n$ close to unity. In this article, and following [5], we restrict to the special $n = 1$ case, giving what will be called the *General Bertalanffy* model[1]. In this case, Bertalanffy was able to provide an analytic solution, reproduced below. Further specializing to the case $m = 2/3$, as Bertalanffy did, we obtain the *classical Bertalanffy* model[2]. As we shortly recall, the classical logistic or Verhulst[3] and Gompertz[4] models for tumor growth [2, 3] can also be regarded as special cases of the General Bertalanffy model. However, instead of using Equation (1), we prefer a reformulation with new parameters that pays attention to dimensional correctness, described next.

Recall that the Box-Cox transformation with parameter $\lambda$ is defined by

$$B_\lambda(x) = \begin{cases} \frac{x^\lambda - 1}{\lambda} & \text{if } \lambda \neq 0 \\ \log(x) & \text{if } \lambda = 0 \end{cases}.$$

This function is not only continuous at $\lambda = 0$, but is an infinitely differentiable function of $x$ and $\lambda$, defined for all real $\lambda$ and $x > 0$. If $n = 1$, then, after a change of parameters, (1) is equivalent to

$$\frac{dv}{dt} = \omega B_\lambda \left( \frac{v_\infty}{v} \right) v, \qquad v > 0, \tag{2}$$

---

[1]`bertalanffy` in `TumorGrowth.jl`
[2]`classical_bertalanffy`
[3]`logistic`
[4]`gompertz`

as is readily established. Here $\omega$ is a parameter with the units of inverse time, and $v_\infty > 0$ a parameter with the units of volume. The non-dimensional Box-Cox exponent $\lambda$ has the following interpretation:

- If $\lambda = 1/3$, (2) is the classical Bertalanffy model (i.e, with $m = 2/3$ in (1)).
- If $\lambda = -1$, (2) is the logistic model.
- If $\lambda = 0$, (2) is the classical Gompertz model.

Supposing $\omega > 0$, the qualitative dynamics of (2) is always the same. There is a single steady state solution $v(t) = v_\infty$, which is always stable. As a first order, autonomous ODE in one dimension, all its other solutions are either monotonically increasing (true here if $v(0) < v_\infty$) or monotonically decreasing (true if $v(0) > v_\infty$). Writing $v_0 = v(0)$, Bertalanffy's analytic solution is

$$v(t) = \begin{cases} \left(1 + \left(\left(\frac{v_0}{v_\infty}\right)^\lambda - 1\right) e^{-\omega t}\right)^{1/\lambda} v_\infty & \text{if } \lambda \neq 0 \\ \left(\frac{v_0}{v_\infty}\right)^{e^{-\omega t}} v_\infty & \text{if } \lambda = 0 \end{cases}.$$

If we non-dimensionalize volumes by $v_\infty$ and times by $1/\omega$, then the solutions can be written

$$v(t) = \begin{cases} \left(1 + \left(v_0^\lambda - 1\right) e^{-t}\right)^{1/\lambda} & \text{if } \lambda \neq 0 \\ v_0^{e^{-t}} & \text{if } \lambda = 0 \end{cases}.$$

### 3.2. A two-dimensional generalization of the General Bertalanffy model.
In patients undergoing treatment, the growth of lesions is frequently not monotonic. For example, a lesion may initially decrease in size, but ultimately begin increasing again. To capture such behavior in an autonomous ODE we must adopt a higher order model or, equivalently, a first order ODE with more than one dimension. A simple two-dimensional extension[5] of the General Bertalanffy model (2) is obtained by replacing the parameter $v_\infty$ with a new latent variable $u(t)$, the (time-varying) *carrying capacity*, which we allow to evolve independently of $v(t)$, at a rate in proportion to its magnitude:

$$(3) \qquad \frac{dv}{dt} = \omega B_\lambda \left(\frac{u}{v}\right) v, \qquad \frac{du}{dt} = \gamma \omega u.$$

Here $\gamma$ is a new dimensionless parameter, taking any real value, which introduces a second time scale $\frac{1}{|\gamma\omega|}$. Since, $u$ is latent, the value of $u(0)$ is unknown. Recycling notation, we retain $v_\infty$ as a model hyperparameter, but it is now the initial carrying capacity $u(0)$. Then, taking $\gamma = 0$, we get $u(t) = v_\infty$ for all $t$, recovering the solution $v(t)$ to the first order model (2).

Our extension described by Equation (3) is it is not based on any physiological mechanism known to us, but is one of the simplest ways to generalize (2).

### 3.3. Neural ODE's.
A *neural ODE* is an ordinary differential equation of the form $\dot{x} = F(x, \theta)$, where $F(x, \theta)$ is the output of some neural network with input $x$, and some system of weights and biases $\theta$ [4]. We consider two classes of neural ODE models in this article, both implemented in `TumorGrowth.jl`. In the one-dimensional neural ODE[6] a volume scale $v_\infty$ and invertible transform $\phi\colon (0, \infty) \to \mathbb{R}$ are specified, and we declare that the non-dimensionalized, transformed volume $y(t) := \phi(v(t)/v_\infty)$ is to evolve according to the ODE

$$\frac{dy}{dt} = f(y(t), \theta).$$

Here $f(\cdot, \theta)$ is any single input, single output, neural network, as constructed using the `Lux.jl` framework, parameterized by $\theta$ [10]. Typically, $\phi = \log$. In our two-dimensional neural ODE[7], $y(t)$

---

[5]`bertalanffy2` in `TumorGrowth.jl`

[6]`neural`

[7]`neural2`

is defined as for the one-dimensional case, but its evolution is coupled with that of a new latent variable $u(t) \in \mathbb{R}$:

$$\frac{dy}{dt} = f_1(y, u)$$
$$\frac{du}{dt} = f_2(y, u).$$

Here, $f_1$ and $f_2$ are the components of any `Lux.jl` neural network $f(\cdot, \theta)$ with two-dimensional input and output.

3.4. **Exponential decay or growth.** For completeness, we also consider the simple model $dv/dt = -\omega v$, whose solutions are exponential decay or growth.[8]

## 4. COMPARING THE MODELS

We now detail the computations leading to the results reported in 1.3. For full details of the computation, refer to the "Modal Battle" section of the `TumorGrowth.jl` documentation [1].

A total of 652 lesion time series were extracted from the data collated in [5], by discarding any example with less than six measurements, which leaves examples with an average of 8.3 measurements per lesion. These lesions come from distinct patients undergoing chemotherapy or cancer immunotherapy. As detailed below, each of the models listed in Table 1 was calibrated individually using the 652 examples, using all but the last two measurements from each example. Then each calibrated model was used to predict volumes for the two holdout times, and the average mean absolute deviation from the measured volumes was recorded. These deviations were analyzed for statistical significance, as already described in 1.3.

As in [5], calibration is achieved by choosing the initial condition $v_0$ and model parameters minimizing the sum of squares loss for the training observations (not the mean absolute error). Instead of using the Trust Region Reflective algorithm implemented in Python's `scipy` package (a variation on Levenberg-Marquardt optimization) we used Adam gradient descent. `TumorGrowth.jl` also provides Levenberg-Marquardt optimization, but will not work for the neural ODE models tested here, where the number of parameters to be optimized exceeds the number of time-volume pairs.

4.1. **Addressing instability during calibration.** The parameter constraints $v_0, v_\infty > 0$ posed some difficulties, especially in the classical models that have a singularity at zero volume. These issues persist if one instead uses Levenberg-Marquardt or Powell's dog leg optimizers. A first step in mitigating the issue was to arrange that `TumorGrowth.jl` handle parameter bounds in the following (non-standard) way: If an optimization step leads to a parameter moving out of bounds, then instead of moving immediately to the boundary (possibly a singularity) the parameter steps half way towards the boundary. As a second mitigation measure we were led to add a loss penalty discouraging large differences between $v_0$ and $v_\infty$ on a log scale. It has the following form:

$$\text{penalized loss} \quad = \quad \left( \frac{v_0^2 + v_\infty^2}{2 v_0 v_\infty} \right)^\kappa \quad \times \quad \text{least squares loss},$$

where $\kappa$ is the penalty strength (`penalty` in `TumorGrowth.jl`). Experiments led to a default of $\kappa = 0.3$ for the neural network models, and $\kappa = 0.8$ for the other models. Nevertheless, about 2.5% of examples were discarded due to instability issues, mostly associated with the neural ODE models.

Since the computations were not overly long, small learning rates were adopted, and optimizers were run for many thousands of iterations, guaranteeing a high degree of convergence to some local

---

[8]`exponential`

optimum. There was no guarantee that local optima were actually global, which may have led to differences in some calibrations between the two studies.

## 5. Discussion

The `TumorGrowth.jl` package provides fast and convenient tools for calibrating and comparing common "textbook" tumor growth models, as well as newer models, such as neural ODE models.

The main conclusion from a statistical analysis of human treatment data, as first collated and analyzed in [5], is that the superior performance of the General Bertalanffy model, when compared to the other models in Table 1, is statistically significant, except in the case of the one-dimensional neural ODE model, which however does not outperform any model to a statistically detectable degree. The exponential model is the poorest classical model, and the two-dimensional neural ODE model is outperformed by all other models. Otherwise, performance differences between the models are not statistically detectable. It is to be emphasized that these conclusions apply only to performance in the average, for a population of examples including at least 4 calibration measurements, 72% of which have less than 7 measurements. Transferring these conclusions to populations with different characteristics is ill advised.

Where more measurements are available, more complex models, such as those which can account for rebound or relapse behavior, are expected to outperform the classical models. This is already evident in particular cases, as demonstrated by Figure 1. Possible candidates include our two-dimensional generalization of the General Bertalanffy model, and neural ODE models, provided these have a relatively a small of nodes to prevent over-fitting. As more clinical data becomes available, finer, statistically detectable differences in model performance may also be possible.

No attempt to stratify model comparisons over individual cancer types (lung versus bladder) or over treatment types has been attempted here, as we suspect statistically significant results will be elusive without more patient data, but this is another interesting area for further investigation. An interesting exercise also not pursued here is to analyze non-human data, where more statistically significant conclusions may be possible.

An option for improving model performance, provided by the `TumorGrowth.jl` package, is to give greater weight to more recent examples during calibration, and such improvement was indeed observed in individual examples. However, introducing this option essentially adds model complexity, as the degree of weighting arguably needs to be learned along with model parameters; this increases the danger of over-fitting. For this reason, the option was not fully investigated in the present study (and is not considered in [5]) but it may be valuable in other contexts.

**Author Contributions.** Conceptualization: A. B., S. O.; Data curation: S. O.; Investigation: A. B.; Methodology: A. B., S. O.; Software development and documentation: A. B.; Numerical Experiments: A. B.; Writing – original draft: A. B.; Writing – review & editing: A. B., S. O.

## References

[1] Blaom A. TumorGrowth.jl: A Julia library for modelling tumor growth; 2024. Available from: https://github.com/ablaom/TumorGrowth.jl.

[2] Kuang Y, Nagy JD, Eikenberry SE. Introduction to mathematical oncology. Chapman & Hall/CRC Mathematical and Computational Biology Series. CRC Press, Boca Raton, FL; 2016.

[3] Norton L, Simon R, Brereton HD, Bogden AE. Predicting the course of Gompertzian growth. Nature. 1976;264(5586):542–545. doi:10.1038/264542a0.

[4] Chen RTQ, Rubanova Y, Bettencourt J, Duvenaud DK. Neural Ordinary Differential Equations. In: Bengio S, Wallach H, Larochelle H, Grauman K, Cesa-Bianchi N, Garnett R, editors. Advances in Neural Information Processing Systems. vol. 31. Curran Associates, Inc.; 2018.Available from: `https://proceedings.neurips.cc/paper_files/paper/2018/file/69386f6bb1dfed68692a24c8686939b9-Paper.pdf`.

[5] Laleh NG, Loeffler CML, Grajek J, Staňková K, Pearson AT, Muti HS, et al. Classical mathematical models for prediction of response to chemotherapy and immunotherapy. PLoS Computational Biology. 2022;18(2). doi:10.1371/journal.pcbi.1009822.

[6] Hornik K, Stinchcombe M, White H. Multilayer feedforward networks are universal approximators. Neural Networks. 1989;2(5):359–366. doi:10.1016/0893-6080(89)90020-8.

[7] Nocedal J, Wright SJ. Numerical Optimization. 2nd ed. Springer Series in Operations Research and Financial Engineering. Springer, New York; 2006.

[8] Ma Y, Dixit V, Innes MJ, Guo X, Rackauckas C. A Comparison of Automatic Differentiation and Continuous Sensitivity Analysis for Derivatives of Differential Equation Solutions. In: 2021 IEEE High Performance Extreme Computing Conference (HPEC); 2021. p. 1–9.

[9] von Bertalanffy L. Quantitative Laws in Metabolism and Growth. The Quarterly Review of Biology. 1957;32(3):217–231. doi:10.1086/401873.

[10] Pal A. Lux: Explicit Parameterization of Deep Neural Networks in Julia; 2023. Available from: `https://doi.org/10.5281/zenodo.7808904`.